

საქართველოს ტექნიკური უნივერსიტეტი

ელგუჯა ყუბანეიშვილი

ბ ი ო მ ე ტ რ ი ა

დამხმარე სასკოლო მასალა

თბილისი

სახელმძღვანელო განკუთვნილია საქართველოს ტექნიკური უნივერსიტეტის ბიოსამედიცინო ტექნიკის კათედრის სტუდენტებისათვის. მისი გამოყენება შეუძლიათ სამედიცინო უმაღლესი სასწავლებლების სტუდენტებს „ბიოსტატისტიკის“ კურსის შესასწავლად, აგრეთვე ინჟინრებს, ექიმებს, ეკონომისტებს, ბიოლოგებს, მაგისტრანტებსა და დოქტორანტებს.

რეცენზენტები:

ფიზიკა-მათემატიკის მეცნიერებათა დოქტორი, პროფესორი **ზ. ფირანაშვილი**

მედიცინის მეცნიერებათა დოქტორი, პროფესორი **ა. ციბაძე**

მეხუთე შესწორებული
გამოცემა

რედაქტორი მ. გიორგობიანი
კომპიუტერული უზრუნველყოფა ლ. ყუბანიშვილი
ქალაქის ზომა 60X84 $\frac{1}{16}$
ნაბეჭდი თაბახი 15,5 შეკვეთა 35

გამომცემლობა
შპს „თოპალინი“, 2005

შესავალი

თანამედროვე სამედიცინო-ბიოლოგიურ კვლევებში ფართოდ გამოიყენება მონაცემების სტატისტიკური დამუშავება არსებული კანონზომიერების გამოსავლენად, რაც, ძირითადად, განპირობებულია პრაქტიკაში კომპიუტერული ტექნოლოგიების ფართო დანერგვით. აქედან გამომდინარე, სულ უფრო იზრდებოდა ინტერესი ქართულ ენაზე სახელმძღვანელოს ან საცნობაროს მიმართ, სადაც მოცემული იქნებოდა ის მათემატიკურ-სტატისტიკური მეთოდები, რომლებიც გამოიყენება კლინიკურ-ექსპერიმენტალური მონაცემების დამუშავებისთვის.

ტერმინი „ბიომეტრია“ (ბერძნული სიტყვიდან *bios* – სიცოცხლე და *metron* – ზომა) პირველად მე-19 საუკუნის დამდეგს (1889 წ.) გამოჩნდა მეცნიერებაში, როგორც ახალი მიმართულება ბიოლოგიაში და მიზნად ისახავდა მათემატიკური მეთოდების გამოყენებას ბიოლოგიურ კვლევებში. ცოტა მოგვიანებით, კერძოდ, 1899 წელს გაჩნდა სხვა დასახელება „ვარიაციული სტატისტიკა“, ხოლო უფრო მოგვიანებით (1901 წ.) – „ბიოლოგიური სტატისტიკა“, რომელიც დღესაც ფართოდ გამოიყენება ბიომეტრიასთან ერთად.

ფორმალური თვალსაზრისით, ბიომეტრია წარმოადგენს ბიოლოგიაში გამოყენებული მათემატიკური მეთოდების ერთობლიობას, რომელიც ძირითადად ნასესხებია ალბათობის თეორიისა და მათემატიკური სტატისტიკის დისციპლინებიდან. განსაკუთრებით მჭიდრო კავშირი აქვს ბიომეტრიას მათემატიკურ სტატისტიკასთან, რომლის შედეგებსაც იგი, ძირითადად, იყენებს. თავის მხრივ, ბიომეტრია გარკვეულ ზეგავლენას ახდენს მათემატიკური სტატისტიკის განვითარებაზე.

ამრიგად, ბიომეტრია არის გამოყენებითი მეცნიერების დარგი, რომელიც სწავლობს ბიოსამედიცინო ინფორმაციის შეკრების, სისტემატიზაციისა და დამუშავების საკითხებს ალბათობის თეორიისა და მათემატიკური სტატისტიკის მეთოდების გამოყენებით. ბიომეტრია, როგორც დამოუკიდებელი სამეცნიერო დისციპლინა, ჩამოყალიბდა მე-19 საუკუნის მეორე ნახევარში, თუმცა, მისი სათავეები უნდა ვეძიოთ მე-17 საუკუნეში, როცა 1614 წელს გამოჩნდა სანტარიოს წიგნი „სტატისტიკური მედიცინა“, ხოლო 1680 წელს გამოვიდა ბორელის წიგნი „ცხოველების მოძრაობა“.

აღბათობის თეორია და მათემატიკური სტატისტიკა აღმოცენდა მე-17 საუკუნეში ერთმანეთისგან დამოუკიდებლად. აღბათობის თეორიის ჩამოყალიბებას სტიმული მისცა აზარტულმა თამაშებმა. მის განვითარებაში დიდი წვლილი მიუძღვით ხ. ჰიუგენსის, ი. ბერნულის, პ. ლაპლასის, კ. გაუსისა და ს. პუასონის შრომებს, ხოლო შემდგომში, პ. ჩეზიშევის, მარკოვისა და ა. ლიაპუნოვის გამოკვლევებს. განსაკუთრებით დიდია ა. კოლმოგოროვის როლი აღბათობის თეორიის, როგორც მათემატიკური მეცნიერების, ჩამოყალიბების საქმეში. აღბათობის თეორია წარმოადგენს მათემატიკის დარგს, რომელიც შეისწავლის შემთხვევითი მოვლენების კანონზომიერებას.

რაც შეეხება მათემატიკურ სტატისტიკას, იგი შეიქმნა იმ დროისათვის განვითარებულ ქვეყნებში, სადაც წამოიჭრა მოთხოვნილებები დემოგრაფიაში, ვაჭრობაში, ჯანმრთელობის დაცვასა და სხვა სამეურნეო სფეროებში აღწერითი სტატისტიკის მიმართ. ამაში დიდი ღვაწლი მიუძღვის ა. კეტელს, რომელმაც პარალელურად შექმნა ბიომეტრიის საფუძველი.

ბიომეტრიის განვითარებაში დიდი ღვაწლი მიუძღვით ფ. გალტონსა და კ. პირსონს, რომლებზედაც დიდი გავლენა იქონია დარვინის მოძღვრებამ. ფ. გალტონმა გამოაქვეყნა მთელი რიგი ორიგინალური ნაშრომები ანთროპოლოგიასა და გენეტიკაში, სადაც გამოყენებულია მათემატიკური სტატისტიკის მეთოდები. ლონდონის უნივერსიტეტის პროფესორმა კ. პირსონმა განავითარა ფ. გალტონის მიერ დანყებული საქმე და შექმნა ბიომეტრიის მათემატიკური აპარატი. კერძოდ, მან შემოიტანა ისეთი ცნებები, როგორცაა „ხი-კვადრატი“, „საშუალო კვადრატული გადახრა“, „ვარიაციის კოეფიციენტი“. მას ეკუთვნის კორელაციისა და რეგრესიის გაუმჯობესებული მეთოდები. 1901 წელს პირსონმა დაიწყო ყურნალ „ბიომეტრიის“ გამოცემა.

ფ. გალტონის და კ. პირსონის დანყებული საქმე გააგრძელეს ვ. გოსეტმა (სტიუდენტი) და რ. ფიშერმა, რომლებმაც ბიომეტრიის განვითარებაში უდიდესი ღვაწლი შეიტანეს. განსაკუთრებით დიდია რ. ფიშერის როლი მათემატიკური სტატისტიკისა და ბიომეტრიის განვითარებაში ახალი მეთოდებისა და ცნებების შემოტანით.

სტატისტიკური მეთოდების გამოყენება ბიოსამედიცინო კვლევებში სულ უფრო და უფრო იზრდება კომპიუტერული ტექნიკისა და ტექნოლოგიების სრულყოფასთან ერთად. დღეისათვის არსებობს ინფორმაციის სტატისტიკური დამუშავების მრავალი კომპიუტერული სისტემები, რომლებიც ორიენტირებულია Win-

dows-ის ოპერაციულ სისტემაზე. მათ შორის შეიძლება გამოვყოთ ისეთი პაკეტები, როგორცაა *Statistica*, *SPSS*, *Statgraphics*, *SAS*, *Excel* და სხვა, რომლებიც ფართოდ გამოიყენებიან ნებისმიერი სახის ინფორმაციის დამუშავებისათვის.

აღბათობის თეორიის საფუძვლები

1. ხდომილობა და მისი აღბათობა

1.1. ხდომილობის კლასიფიკაცია

აღბათობის თეორიის ერთ-ერთი ძირითადი ცნებაა ხდომილობის ცნება. მოვლენას ან ფაქტს, რომელიც წარმოიშობა ცდის შედეგად, **ხდომილობა** ეწოდება. **ცდა** (ექსპერიმენტი) ეწოდება ისეთ მოქმედებას, რომლის დროსაც სრულდება გარკვეულ პირობათა კომპლექსი, რაც აუცილებელია რაიმე მოვლენის წარმოსაქმნელად. ხდომილობის უმარტივეს მაგალითს წარმოადგენს მონეტის აგდების შედეგი. ამ მაგალითში ცდას წარმოადგენს მონეტის აგდება, ხოლო ხდომილობას – ამ აგდების შედეგი. ამ ცდაში გვაქვს გერბისა და ციფრის გამოჩენის ხდომილობები. ხდომილობებს, როგორც წესი, აღნიშნავენ ლათინური ანბანის დიდი ასოებით *A, B, C, ...*

ხდომილობა შეიძლება იყოს შეთავსებადი და შეუთავსებადი. ორ ხდომილობას ეწოდება **თავსებადი**, თუ ერთი მათგანის მოხდენა იწვევს ან არ გამორიცხავს მეორე ხდომილობის მოხდენას იმავე ცდაში. წინააღმდეგ შემთხვევაში, ხდომილობები იქნებიან **შეუთავსებადი**. შეუთავსებადი ხდომილობების მაგალითია გერბისა და ციფრის გამოჩენა მონეტის აგდების ცდაში.

ხდომილობას ეწოდება **უტყუარი**, თუ იგი აუცილებლად მოხდება მოცემული ცდის პირობებში. მაგალითად, ნორმალური ატმოსფერული წნევის პირობებში წყალი იყინება ნულოვანი ტემპერატურის დროს. ხდომილობას ეწოდება **შეუძლებელი**, თუ იგი არ შეიძლება მოხდეს მოცემული ცდის პირობებში. მაგალითად, თოვლის მოსვლა, როცა ატმოსფეროს ტემპერატურა დადებითია.

ხდომილობას ეწოდება **შემთხვევითი**, თუ ცდის შედეგად იგი შეიძლება მოხდეს ან არ მოხდეს. მაგალითად, სროლის შედეგად მიზანში მოხვედრა. ერთ-ერთი მნიშვნელოვანი ცნებაა ხდომილობების სრული ჯგუფის ცნება. ცდის დროს რამდენიმე ხდომილობა ადგენს **სრულ ჯგუფს**, თუ ცდის შედეგად აუცილებლად მოხდება ერთ-ერთი მათგანი მაინც. მაგალითად, ციფრისა და გერბის გამოჩენა მონეტის აგდების დროს ადგენენ ხდომილობათა სრულ ჯგუფს.

A ხდომილობის მოპირდაპირე ხდომილობა ეწოდება \bar{A} ხდომილობას, რომელიც აუცილებლად მოხდება, თუ არ მოხდება A ხდომილობა. მოპირდაპირე ხდომილობები შეუთავსებადი არიან და ადგენენ ხდომილობათა სრულ ჯგუფს. მაგ. თუ დამზადებული პროდუქციის პარტია შედგება ვარგისი და წუნი პროდუქციისგან, მაშინ ერთ-ერთი პროდუქციის შემონმებისას, იგი შეიძლება იყოს ვარგისი ან წუნი (A ან \bar{A}).

1.2. მოქმედებები ხდომილობაზე

ალბათობის თეორიაში, შემთხვევითი ხდომილობების შესწავლისას, მეტად მნიშვნელოვან ცნებას წარმოადგენს შეკრების (გაერთიანების) და ნამრავლის (თანაკვეთის) ცნებები. ორი A და B ხდომილობების ჯამი ეწოდება ისეთ ხდომილობას, რომელიც ხდება მაშინ და მხოლოდ მაშინ, როცა ხდება A და B ხდომილობებიდან ერთი მაინც:

$$C = A + B \quad \text{ან} \quad C = A \cup B.$$

ანალოგიურად, რამდენიმე ხდომილობების ჯამი ეწოდება ხდომილობას, რომელიც მდგომარეობს თუნდაც ერთ-ერთის მოხდენაში:

$$E = A + B + C + \dots + N.$$

A და B ხდომილობების ნამრავლი ეწოდება Q ხდომილობას, რომელიც მდგომარეობს A და B ხდომილობების ერთდროულად მოხდენაში:

$$Q = A \cdot B \quad \text{ან} \quad Q = A \cap B.$$

ანალოგიურად, რამდენიმე ხდომილობების ნამრავლი ეწოდება ხდომილობას, რომელიც მდგომარეობს ამ ხდომილობათა ერთდროულად მოხდენაში:

$$W = A \cdot B \cdot C \cdot \dots \cdot N$$

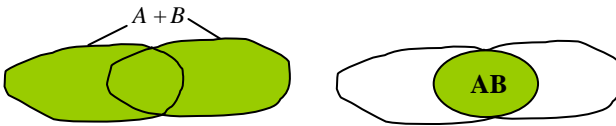
ხდომილობების ჯამისა და ნამრავლის განსაზღვრიდან გამომდინარეობს:

$$A + A = A \quad \text{და} \quad AA = A.$$

ზოგჯერ ერთი, მაგ. B ხდომილობის მოხდენა იწვევს მეორე A ხდომილობის მოხდენას. მაშინ ამბობენ, რომ B ხდომილობა შედის A ხდომილობაში და აღინიშნება შემდეგნაირად: $B \subset A$. თუ B ხდომილობა ეკუთვნის A ხდომილობას, მაშინ ადგილი აქვს შემდეგ ტოლობებს:

$$A + B = A \quad \text{და} \quad AB = B.$$

ხდომილობათა ჯამისა და ნამრავლის ცნებებს გააჩნიათ თვალსაჩინო გეომეტრიული ინტერპრეტაცია. მაგალითად, ისეთი, როგორიც წარმოდგენილია შემდეგ ნახაზზე:



1.3. ხდომილობის ალბათობა

ხდომილობათა რაოდენობრივი შედარებისათვის, იმის მიხედვით, თუ რამდენად შესაძლებელია მათი მოხდენა, შემოღებულია გარკვეული ზომა, რომელსაც ეწოდება ხდომილობის ალბათობა. ამრიგად, **ხდომილობის ალბათობა** ეწოდება ამ ხდომილობის მოხდენის შესაძლებლობის ობიექტური ხარისხის რიცხვით ზომას. არსებობს შემთხვევითი ხდომილობის ალბათობის განსაზღვრის ორი მეთოდი: უშუალო და სტატისტიკური.

A შემთხვევითი ხდომილობის მოხდენის ალბათობა მოცემული ცდისათვის უშუალოდ გამოითვლება შემდეგნაირად:

$$P(A) = \frac{m}{n}, \tag{1.1}$$

სადაც, $P(A)$ – A ხდომილობის ალბათობა; n – შემთხვევათა საერთო რაოდენობა; m – A ხდომილობის ხელშემწყობ შემთხვევათა რაოდენობა.

შემთხვევას ეწოდება **ხელშემწყობი**, თუ ამ შემთხვევის მოხდენა იწვევს მოცემული ხდომილობის მოხდენას.

ალბათობის განსაზღვრიდან გამომდინარეობს შემდეგი შედეგები:

1. შემთხვევითი A ხდომილობის ალბათობა დადებითი სიდიდეა, რომელიც აკმაყოფილებს შემდეგ პირობებს:

$$0 \leq P(A) \leq 1.$$

მართლაც, რადგან ხელშემწყობ შემთხვევათა რაოდენობა m მოთავსებულია 0 -სა და n -ს შორის, ამიტომ ალბათობა, გამოთვლილი (1.1) ფორმულით ყოველთვის აკმაყოფილებს ამ პირობას.

2. უტყუარ ხდომილობათა ალბათობა ერთის ტოლია.

მართლაც, რადგან $m = n$, ამიტომ $P(A) = \frac{m}{n} = \frac{n}{n} = 1$.

3. შეუძლებელი ხდომილობის ალბათობა ნულია ტოლია.

მართლაც, რადგან $m = 0$, ამიტომ $P(A) = 0$.

ალბათობის განსაზღვრის (1.1) ფორმულას ხშირად უწოდებენ „კლასიკურს“. განვიხილოთ მაგალითი. ყუთში მოთავსებულია 2 თეთრი და 3 შავი ბურთულა. რას უდრის იმის ალბათობა, რომ შემთხვევით ამოღებული ბურთულა თეთრი ფერისა იქნება?

თეთრი ბურთულას ამოღების ხდომილობა აღვნიშნოთ A -თი. რადგან, $n = 5$ და $m = 2$ ამიტომ $P(A) = \frac{2}{5}$.

ალბათობის კლასიკური გამოთვლის დადებითი თვისება ისაა, რომ იგი არ მოითხოვს ცდის ჩატარებას და ეფუძნება ლოგიკურ მსჯელობას. ალბათობის უშუალო განსაზღვრა შეიძლება გამოყენებულ იქნეს მხოლოდ იმ შემთხვევაში, როცა შემთხვევითი ხდომილობები ადგენენ შეუთავსებად ტოლშესაძლო ხდომილობათა სრულ ჯგუფს. ასეთი სიტუაციები პრაქტიკაში იშვიათად გვხვდება, ამიტომ ხდომილობის ალბათობის განსაზღვრისათვის უფრო ხშირად მიმართავენ სტატისტიკურ მეთოდს.

ჯერ განვიხილოთ ხდომილობის **ფარდობითი სიხშირის** ცნება. ამისათვის ჩავატაროთ რაიმე ცდა, რომლის დროს შეიძლება მოხდეს ან არ მოხდეს A ხდომილობა. ვთქვათ, ჩატარდა ასეთი n ცდა. A ხდომილობის ფარდობითი სიხშირე ეწოდება სიდიდეს, რომელიც მიიღება A ხდომილობის მოხდენის რაოდენობის შეფარდებით ცდათა მთელ რაოდენობაზე, ე.ი.

$$P^*(A) = \frac{m}{n},$$

სადაც, $P^*(A)$ – A ხდომილობის ფარდობითი სიხშირეა;

n – ცდათა საერთო რაოდენობა;

m – ცდების რაოდენობა, როდესაც მოხდა A ხდომილობა.

ცდათა მცირე რაოდენობის შემთხვევაში, ხდომილობების ფარდობით სიხშირეს შემთხვევითი ხასიათი გააჩნია. მაგრამ, ცდების რაოდენობის ზრდასთან ერთად, ხდომილობების ფარდობითი სიხშირე თანდათან კარგავს თავის შემთხვევით ხასიათს და უახლოვდება ხდომილობის ალბათობას. ამ ფაქტში კარგად ჩანს დიდ რიცხვთა კანონის მოქმედება, რომლის თეორიული დასაბუთება მოგვცა ცნობილმა შვეიცარიელმა მეცნიერმა იაკობ ბერნულიმ. დიდ რიცხვთა კანონი ამტკიცებს, რომ A ხდომილობის ფარდობითი სიხშირე მით უფრო ახლოს იქნება მისივე ალბათობასთან, თუ ცდათა რაოდენობას უსასრულოდ გავზრდით.

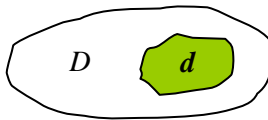
მაგალითი. სროლის ჩატარების შემდეგ მიზანში მოხვედრის ფარდობითი სიხშირე 0,8-ის ტოლია. რამდენ გასროლას ჰქონდა ადგილი, თუ ცნობილია, რომ იგი მიზანს 12-ჯერ მოხვდა?

თუ მიზანში მოხვედრის ხდომილებას A -თი აღვნიშნავთ,

მაშინ $P^*(A) = 0,8$, $m = 12$, ე.ი. $0,8 = \frac{12}{n}$, აქედან $n = 15$.

ალბათობის სტატისტიკური განსაზღვრის დადებითი თვისება ისაა, რომ ის ეფუძნება რეალურ ექსპერიმენტს. ნაკლი ისაა, რომ ალბათობის საიმედო განსაზღვრისათვის, საჭიროა ცდების დიდი რაოდენობის ჩატარება, რაც ძალზე ხშირად ეკონომიურად მიუღებელია.

ალბათობის გეომეტრიული განსაზღვრა. ვთქვათ, სიბრტყეზე გვაქვს D არე, რომლის ფართობია S_D და მასში მოთავსებულია რაიმე d არე S_d ფართობით. D არეში ალაღბედზე ისვრიან წერტილს.



საძიებელია იმის ალბათობა, რომ ეს წერტილი მოხვდება d არეში. აქვე იგულისხმება, რომ წერტილის მოხვედრის ალბათობა D არის ნებისმიერ წერტილში თანაბარია. ამ შემთხვევაში წერტილის d არეში მოხვედრის ალბათობა გამოითვლება შემდეგნაირად:

$$P = \frac{S_d}{S_D}.$$

ამრიგად, თუ რაიმე W სივრცე წარმოადგენს სასრულ მონაკვეთს (ფიგურას, სხეულს), მაშინ რაიმე A ხდომილობის გეომეტრიული ალბათობა ეწოდება ამ ხდომილობის შესაბამისი მონაკვეთის სიგრძის (ფიგურის ფართობის, სხეულის მოცულობის) შეფარდებას W სივრცის შესაბამისი მონაკვეთის სიგრძესთან (ფიგურის ფართობთან, სხეულის მოცულობასთან).

მაგალითი. R რადიუსიან წრეში ჩახაზულია ტოლფერდა სამკუთხედი. რას უდრის იმის ალბათობა, რომ წერტილი მოხვდება სამკუთხედში?

$$P = \frac{3\sqrt{3}R^2}{4\pi R^2} = \frac{3\sqrt{3}}{4\pi} \approx 0,4137.$$

1.4. კომბინატორიკის ძირითადი ფორმულები

1. გამრავლების წესი. ვთქვათ საჭიროა მიყოლებით შევასრულოთ k რაოდენების მოქმედება. ამასთან, პირველი მოქმედება შეიძლება შევასრულოთ n_1 ხერხით, მეორე – n_2 ხერხით და ა. შ. k მოქმედებამდე, რომელიც შეიძლება შევასრულოთ n_k ხერხით. მაშინ m რაოდენების ხერხი, რომლებიც შეიძლება შევასრულოთ ყველა k მოქმედებით ტოლია:

$$m = n_1 \cdot n_2 \cdot \dots \cdot n_r.$$

მაგალითი. ვთქვათ გვაქვს 5 სხვადასხვა ჰალსტუხი, 8 პერანგი და 3 კოსტუმი. ერთ კომპლექტში მათი გამოყენების საერთო ვარიანტების რაოდენობა ტოლია: $m = 5 \cdot 8 \cdot 3 = 120$

გამრავლების წესი ხშირად გამოიყენება ერთი და იგივე სიმრავლიდან ელემენტების მრავალჯერადი გამოყენების რაოდენ-

ნობის გამოსათვლელად, მაგალითად, ციფრებიანი ნომრების შესადგენად. ამ შემთხვევაში სიმრავლის ელემენტი გამოყენების შემდეგ ისევ უბრუნდება სიმრავლეს. მაშინ ამბობენ, რომ სრულდება **ამორჩევა დაბრუნებით**. k სიმრავლის ერთი და იგივე n რაოდენობის ელემენტისთვის რაოდენობის მოქმედების შესრულების ხერხის რაოდენობა m ტოლია:

$$m = n^k .$$

მაგალითი. დავუშვათ საკეტების კოდირებისათვის იყენებენ ხუთ ციფრს 0-დან 5-მდე, მაშინ კოდების რაოდენობა ტოლია:
 $m = 5^3 = 125$.

2. გადანაცვლება. განვიხილოთ n ელემენტთა ერთობლიობა. ნებისმიერად დავალაგოთ ელემენტები ერთი მეორეს მიყოლებით. შედეგად მივიღებთ ელემენტთა თანმიმდევრობის ერთ-ერთ შესაძლო ვარიანტს, რომელსაც გადანაცვლება ეწოდება. n ელემენტის ყველა შესაძლო გადანაცვლებათა რიცხვი, რომელიც აღინიშნება P_n სიმბოლოთი, განისაზღვრება შემდეგი ფორმულით:

$$P_n = n! ,$$

სადაც $n! = n(n-1)(n-2) \cdots 1$ და მას ფაქტორიალი ეწოდება.

მაგალითი. მაგიდის გარშემო ზის 7 ადამიანი. მათი სხვადასხვა გადანაცვლების რაოდენობა ტოლია: $P_7 = 7! = 5040$.

3. წყობა. დავუშვათ, ყუთში მოთავსებულია n რაოდენობის გადანომრილი ბურთულა და შემდეგ ყუთიდან ვიღებთ ერთმანეთის მიყოლებით m რაოდენობის ბურთულას. მივიღებთ m ელემენტისგან შემდგარ მოწესრიგებულ სიმრავლეს. თუ გავიმეორებთ ცდას, მაშინ ვღებულობთ სხვადასხვა სიმრავლეებს. n -ელემენტიანი სიმრავლის ყოველ m -ელემენტიან დალაგებულ ქვესიმრავლეს ($m \leq n$) ეწოდება წყობა n ელემენტისაგან m -ად. n -ელემენტიანი სიმრავლის ყველა შესაძლო m -ელემენტიან წყობათა რიცხვი აღინიშნება A_n^m სიმბოლოთი და განისაზღვრება ფორმულით:

$$A_n^m = \frac{n!}{(n-m)!} \quad 0 \leq m \leq n .$$

გადანაცვლება წარმოადგენს განლაგების კერძო შემთხვევას, როცა $m = n$ ამიტომ $A_n^m = P_n$.

მაგალითი. ჯგუფში 9 გოგონაა და 11 ბიჭი. წარმომადგენლობით ფორუმზე ამ ჯგუფიდან ირჩევენ 3 პიროვნებას, რომლებიც შერჩევის შემდეგ ლეზულობენ რიგით ნომერს და ქმნიან მწკრივს. ასეთი მწკრივების შექმნის რაოდენობა ტოლია:

$$A_{20}^3 = \frac{20!}{(20-3)!} = \frac{20!}{17!} = 20 \cdot 19 \cdot 18 = 6840.$$

4. ჯუფთება. მოცემულია n რაოდენობის ელემენტთა ერთობლიობა. გვანტერესებს, რამდენი ხერხით შეიძლება ამ ერთობლიობიდან შევარჩიოთ m რაოდენობის ელემენტები. ელემენტის შერჩევის რიგითობას, განლაგებასთან შედარებით, ამ შემთხვევაში მნიშვნელობა არ აქვს. ე.ი. ვლუბულობთ არამონესრიგებულ სიმრავლევებს. გარდა ამისა, ცდაში ამოღებული ელემენტი უკან არ ბრუნდება, ამიტომ საქმე გვაქვს ამორჩევასთან დაბრუნების გარეშე.

n -ელემენტიანი სიმრავლის ყოველ m -ელემენტიან არამონესრიგებულ ქვესიმრავლეს ($m \leq n$) ეწოდება ჯუფთება n ელემენტიდან m -ად. n -ელემენტიანი სიმრავლის ყველა შესაძლო m -ელემენტიან ჯუფდებათა რიცხვი აღინიშნება C_n^m სიმბოლოთი და განისაზღვრება შემდეგი ფორმულით:

$$C_n^m = \frac{n!}{m!(n-m)!} \quad 0 \leq m \leq n .$$

მაგალითი. ყუთში მოთავსებულია 10 ბურთულა. ყუთიდან იღებენ ორ-ორ ბურთულას. რას უდრის ორი ბურთულას ამოღების ყველა შესაძლო მნიშვნელობათა რაოდენობა?

$$C_{10}^2 = \frac{10!}{2!8!} = \frac{10 \cdot 9}{1 \cdot 2} = 45 .$$

ჯუფთებათა რიცხვი, გადანაცვლებათა რიცხვი და წყობათა რიცხვი დაკავშირებულია ერთმანეთთან შემდეგი ფორმულით:

$$A_n^m = P_n \cdot C_n^m .$$

ალბათობის გამოსათვლელად ხშირად გამოიყენება ჯუფთებათა რიცხვი. განვიხილოთ მაგალითი.

მაგალითი. ყუთში მოთავსებულია 7 თეთრი და 3 შავი ბურთულა. რას უდრის იმის ალბათობა, რომ შემთხვევით ამოღებული ორი ბურთულიდან, ორივე თეთრი იქნება?

ორი თეთრი ბურთულას ამოღების ხდომილობა აღვნიშნოთ A -თი. მაშინ

$$P(A) = \frac{m}{n}.$$

ორი ბურთულას ამოღების ყველა შესაძლო შემთხვევათა რაოდენობა ტოლია:

$$n = C_{10}^2 = \frac{10!}{2!(10-2)!} = \frac{10 \cdot 9}{1 \cdot 2} = 45.$$

თეთრი ბურთულების ამოღების შესაძლო შემთხვევათა რაოდენობა იქნება:

$$m = C_7^2 = \frac{7 \cdot 6}{1 \cdot 2} = 21.$$

$$\text{ამრიგად, } P(A) = \frac{21}{45} = \frac{7}{15}.$$

1.5. ალბათობის თეორიის ძირითადი თეორემები

პრაქტიკაში ხდომილობების ალბათობის გამოსათვლელად, ძირითადად, იყენებენ არაპირდაპირ მეთოდებს, რომლებიც საშუალებას გვაძლევს ერთი ცნობილი ხდომილობის ალბათობით გამოვთვალოთ სხვა ხდომილობების ალბათობები. არაპირდაპირი მეთოდების გამოყენებისას ვსარგებლობთ ალბათობის თეორიის ძირითადი თეორემებით. ასეთი თეორემა ორია: ალბათობების შეკრებისა და ალბათობების გამრავლების თეორემები.

ალბათობების შეკრების თეორემა. ორი შეუთავსებადი ხდომილობების ჯამის ალბათობა ტოლია ამ ხდომილობების ალბათობების ჯამისა.

$$P(A+B) = P(A) + P(B). \quad (1.2)$$

მართლაც, თუ n ცდიდან A ხდომილობა მოხდა m -ჯერ, ხოლო B – k -ჯერ, მაშინ

$$P(A) = \frac{m}{n}, \quad P(B) = \frac{k}{n}.$$

რადგან A და B ხდომილობები შეუთავსებადი არიან, ამიტომ

$$P(A+B) = \frac{m+k}{n}.$$

თუ ამ მნიშვნელობებს ჩავსვამთ (1.2) ფორმულაში, მივიღებთ იგივე-ეობას. თეორემა დამტკიცებულია.

ეს შედეგი შეგვიძლია განვაზოგადოდ ნებისმიერი რაოდენობის შეუთავსებად ხდომილობებზე. მაშინ გვექნება:

$$P\left(\sum_{i=1}^l A_i\right) = \sum_{i=1}^l P(A_i).$$

ალბათობების შეკრების თეორემიდან გამომდინარეობს შემდეგი შედეგები:

1. თუ A_1, A_2, \dots, A_n ხდომილობები ადგენენ შეუთავსებად ხდომილობათა სრულ ჯგუფს, მაშინ

$$\sum_{i=1}^n P(A_i) = 1.$$

2. ურთიერთმოპირდაპირე ხდომილობათა ალბათობების ჯამი ერთის ტოლია.

$$P(A) + P(\bar{A}) = 1.$$

ზოგჯერ, პრაქტიკული ამოცანების ამოხსნისას, უფრო ადვილია რაიმე A ხდომილობის ალბათობის განსაზღვრა მისი მოპირდაპირე \bar{A} ხდომილობის ალბათობით: $P(A) = 1 - P(\bar{A})$.

თუ საქმე გვაქვს თავსებად ხდომილობებთან, მაშინ ალბათობის შეკრების თეორემა ასე ჩამოყალიბდება: ორი თავსებადი ხდომილობათა ალბათობების ჯამი ტოლია:

$$P(A+B) = P(A) + P(B) - P(AB).$$

თუ გვაქვს სამი ხდომილობა, მაშინ

$$P(A+B+C) = P(A) + P(B) + P(C) - P(AB) - P(AC) - P(BC) + P(ABC).$$

ზოგადად, l ხდომილობის დროს გვექნება:

$$P\left(\sum_i A_i\right) = \sum_i P(A_i) - \sum_{i,j} P(A_i A_j) + \sum_{i,j,k} P(A_i A_j A_k) - \dots + (-1)^{l-1} P(A_1 A_2 \dots A_l)$$

მაგალითი. ვთქვათ, ლატარიაში თამაშდება 1000 ბილეთი. აქედან, ერთი იგებს 500 ლარს, 10 ბილეთი – 100 ლარს, 50 ბილეთი – 20 ლარს და 100 ბილეთი – 5 ლარს. დანარჩენი ბილეთები არაფერს არ იგებენ. ვყიდულობთ ერთ ბილეთს. გვინტერესებს, რას უდრის იმის ალბათობა, რომ ჩვენ მოვიგებთ არა უმეტეს 20 ლარსა. მოგების ხდომილობა აღვნიშნოთ A -თი. 20 ლარი მოგებისა – A_1 -ით, 100 ლარი მოგებისა – A_2 -ით და 500 ლარი მოგებისა A_3 -ით. ცხადია, რომ $A = A_1 + A_2 + A_3$. ალბათობის შეკრების თეორემის თანახმად,

$$P(A) = P(A_1) + P(A_2) + P(A_3) = \frac{50}{1000} + \frac{10}{1000} + \frac{1}{1000} = 0,05 + 0,01 + 0,001 = 0,061.$$

ალბათობების გამრავლების თეორემის ჩამოყალიბებამდე განვიხილოთ რამდენიმე განსაზღვრება.

1. A ხდომილობას ეწოდება **დამოუკიდებელი** B ხდომილობისაგან, თუ A ხდომილობა არ არის დამოკიდებული იმაზე, მოხდა თუ არა B ხდომილობა.

2. A ხდომილობას ეწოდება B ხდომილობაზე **დამოკიდებული**, თუ A ხდომილობის ალბათობა დამოკიდებულია იმაზე, მოხდა თუ არა B ხდომილობა.

3. A ხდომილობის ალბათობას, გამოთვლილს იმ პირობით, რომ ადგილი ჰქონდა B ხდომილობას, ეწოდება A ხდომილობის **პირობითი ალბათობა** და აღინიშნება $P(A|B)$ ან $P_B(A)$ სიმბოლოთი.

ცხადია, რომ თუ A და B ხდომილებები ერთმანეთის მიმართ დამოუკიდებელნი არიან, მაშინ $P(A|B) = P(A)$, ხოლო თუ ისინი დამოკიდებულნი არიან, მაშინ $P(A|B) \neq P(A)$.

ალბათობების გამრავლების თეორემა. ორი A და B ხდომილობის ნამრავლის ალბათობა ტოლია ერთ-ერთი მათგანის მოხდენის ალბათობის ნამრავლისა მეორე ხდომილობის პირობით ალბათობაზე იმ პირობით, რომ ადგილი ჰქონდა პირველ ხდომილობას. ე.ი.

$$P(AB) = P(A) P(B | A), \text{ ან } P(AB) = P(B) P(A | B). \quad (1.3)$$

მართლაც, ვთქვათ n ცდიდან A ხდომილობა მოხდა m -ჯერ, ხოლო B – k -ჯერ. რადგან A და B ხდომილობები არ არიან თავსებადი, ამიტომ მათი ერთდროული მოხდენის რაოდენობა აღვნიშნოთ l -ით. მაშინ გვექნება:

$$P(AB) = \frac{l}{n}; P(A) = \frac{m}{n}.$$

გამოვთვალოთ B ხდომილობის ალბათობა, როცა A ხდომილობა უკვე მოხდა, ე.ი.

$$P(B | A) = \frac{l}{m}.$$

თუ ამ მნიშვნელობებს ჩავსვამთ (1.3) ფორმულაში, მივიღებთ იგივეობას. თეორემა დამტკიცებულია.

(1.3) ფორმულა შეიძლება განზოგადდეს თანამამრავლთა ნებისმიერი სასრული რაოდენობისათვის:

$$P\left(\prod_{i=1}^n A_i\right) = P(A_1)P(A_2 | A_1)P(A_3 | A_1A_2) \dots P(A_n | A_1A_2 \dots A_{n-1}).$$

ალბათობების გამრავლების თეორემიდან გამომდინარეობს შემდეგი შედეგები:

1. ორი A და B თავსებადი ხდომილობების ერთდროული მოხდენის ალბათობა ნულის ტოლია $P(AB) = 0$.

2. თუ A ხდომილობა დამოკიდებულია B ხდომილობაზე, მაშინ B ხდომილობაც დამოკიდებულია A ხდომილობაზე.

3. დამოუკიდებელ ხდომილობათა ნამრავლის ალბათობა ტოლია ამ ხდომილობათა ალბათობების ნამრავლისა.

$$P\left(\prod_{i=1}^n A_i\right) = \prod_{i=1}^n P(A_i).$$

მაგალითი 1. ყუთში არის ორი თეთრი და სამი შავი ბურთულა. ყუთიდან მიმდევრობით იღებენ ორ ბურთულას. ვიპოვოთ ალბათობა იმისა, რომ ორივე ბურთულა თეთრია.

აღვნიშნოთ ორი თეთრი ბურთულას ამოღების ხდომილობა A -თი. თავის მხრივ, $A = A_1 \cdot A_2$, სადაც, A_1 არის თეთრი ბურთულას პირველად ამოღების ხდომილობა, A_2 – თეთრი ბურთულას მეორედ ამოღების ხდომილობა. ალბათობების გამრავლების თეორემის თანახმად გვექნება:

$$P(A) = P(A_1)P(A_2 | A_1) = \frac{2}{5} \cdot \frac{1}{4} = 0,1.$$

მაგალითი 2. პირველი მსროლელის მიზანში მოხვედრის ალბათობაა $0,8$, ხოლო მეორესი – $0,6$. მსროლელებმა მოახდინეს თითო გასროლა. რას უდრის იმის ალბათობა, რომ მიზანში მოახვედრებს რომელიმე მათგანი?

A -თი აღვნიშნოთ პირველი მსროლელის მიზანში მოხვედრის ხდომილობა, B -თი – მეორე მსროლელის მიზანში მოხვედრის ხდომილობა. C -თი აღვნიშნოთ რომელიმე მათგანის მიზანში მოხვედრის ხდომილობა. ცხადია, რომ $C = A + B$.

რადგან A და B ხდომილობები თავსებადი და დამოუკიდებელი არიან, ამიტომ

$$P(C) = P(A) + P(B) - P(A)P(B) = 0,8 + 0,6 - 0,8 \cdot 0,6 = 0,92.$$

1.6. სრული ალბათობის ფორმულა. ბაიესის ფორმულა

ალბათობების შეკრებისა და გამრავლების თეორემებიდან გამომდინარეობს სრული ალბათობის ფორმულა. თუ A ხდომილობა შეიძლება მოხდეს მხოლოდ იმ პირობით, რომ მოხდა თუნდაც ერთი ხდომილობა H_1, H_2, \dots, H_n ხდომილობებისაგან შემდგარ შეუთავსებად ხდომილობათა სრული ჯგუფიდან, მაშინ A ხდომილობის მოხდენის ალბათობა გამოითვლება ფორმულით:

$$P(A) = \sum_{i=1}^n P(H_i)P(A|H_i),$$

რომელსაც სრული ალბათობის ფორმულა ეწოდება. გამოვიყვანოთ ეს ფორმულა.

რადგან H_1, H_2, \dots, H_n ხდომილობათა ერთობლიობა ქმნის ხდომილობათა სრულ ჯგუფს, ამიტომ A ხდომილობა შეიძლება მოხდეს ერთ-ერთი ამ ხდომილობასთან ერთად. ე.ი.

$$A = AH_1 + AH_2 + \dots + AH_n.$$

რადგან H_1, H_2, \dots, H_n შეუთავსებადი ხდომილობებია, ამიტომ AH_1, AH_2, \dots, AH_n ხდომილობებიც შეუთავსებადია. თუ გამოვიყენებთ ალბათობების შეკრების თეორემას, მივიღებთ:

$$P(A) = P(AH_1) + P(AH_2) + \dots + P(AH_n) = \sum_{i=1}^n P(AH_i). \quad (1.4)$$

A და H_i ხდომილობების ნამრავლი იქნება:

$$P(AH_i) = P(H_i)P(A | H_i).$$

თუ ამ გამოსახულებას ჩავსვამთ (1.4) ფორმულაში, მივიღებთ სრული ალბათობის ფორმულას:

$$P(A) = \sum_{i=1}^n P(H_i)P(A | H_i).$$

აქ მოყვანილ H_1, H_2, \dots, H_n ხდომილობებს ჰიპოთეზებს უწოდებენ.

მაგალითი. გვაქვს სამი ყუთი. პირველ ყუთში არის ორი თეთრი და ერთი შავი ბურთულა, მეორე ყუთში – სამი თეთრი და ერთი შავი, მესამე ყუთში – ორი თეთრი და ორი შავი ბურთულა. რომელიმე ერთ-ერთი ყუთიდან ამოვიღოთ ერთი ბურთულა. ვიპოვოთ რას უდრის იმის ალბათობა, რომ ეს ბურთულა თეთრია. განვიხილოთ სამი ჰიპოთეზა. H_1 – პირველი ყუთის შერჩევა, H_2 – მეორე ყუთის შერჩევა, H_3 – მესამე ყუთის შერჩევა. თეთრი ბურთულას ამოღების ხდომილობა აღვნიშნოთ A -თი. რადგან ჰიპოთეზები ტოლშესაძლო ხდომილობები არიან, ამიტომ $P(H_1) = P(H_2) = P(H_3) = \frac{1}{3}$. გამოვთვალოთ A ხდომილობის პირობითი ალბათობები

$$P(A | H_1) = \frac{2}{3}; \quad P(A | H_2) = \frac{3}{4}; \quad P(A | H_3) = \frac{1}{2}.$$

სრული ალბათობის ფორმულის გამოყენებით მივიღებთ:

$$P(A) = \frac{1}{3} \cdot \frac{2}{3} + \frac{1}{3} \cdot \frac{3}{4} + \frac{1}{3} \cdot \frac{1}{2} = \frac{23}{36}.$$

ბაიესის ფორმულა. პრაქტიკაში ფართოდ გამოიყენება ბაიესის ფორმულა. ვთქვათ, მოცემულია შეუთავსებად ხდომილობათა სრული ჯგუფი H_1, H_2, \dots, H_n და ამ ხდომილობების აპრიორული (ცდამდე) ალბათობები $P(H_1), P(H_2), \dots, P(H_n)$. ჩატარდა ცდა, რომლის შედეგად მოხდა A ხდომილობა. საჭიროა გამოვთვალოთ აპოსტერიორული (ცდის შემდგომი) ალბათობები, ე.ი. შემდეგი პირობითი ალბათობები $P(H_1|A), P(H_2|A), \dots, P(H_n|A)$. ალბათობების გამრავლების თეორემის თანახმად,

$$P(AH_j) = P(A) P(H_j | A) \quad \text{და} \quad P(AH_j) = P(H_j) P(A | H_j).$$

ამ გამოსახულებათა მარცხენა მხარეები ტოლია, მაშინ

$$P(A) P(H_j | A) = P(H_j) P(A | H_j),$$

აქედან

$$P(H_j | A) = \frac{P(H_j) P(A | H_j)}{P(A)}, \quad \text{სადაც, } P(A) \neq 0.$$

თუ $P(A)$ -ს გამოვსახავთ სრული ალბათობის ფორმულით, მივიღებთ:

$$P(H_j | A) = \frac{P(H_j) P(A | H_j)}{\sum_{j=1}^n P(H_j) P(A | H_j)}.$$

ამ ფორმულას ეწოდება ბაიესის ფორმულა.

მაგალითი. სამ ყუთში მოთავსებულია ერთი და იგივე სახის პროდუქცია. პირველ ყუთში არის 10 ცალი, რომელთაგან 3 არასტანდარტულია. მეორეშია 15 ცალი, აქედან 5 არასტანდარტულია და მესამეში 20 ცალი, აქედან 6 არასტანდარტულია. ერთ-ერთი ყუთიდან ალაღბედზე იღებენ ერთ პროდუქციას და იგი, ვთქვათ, აღმოჩნდა არასტანდარტული. რას უდრის იმის ალბათობა, რომ ეს პროდუქტი ამოღებულია მეორე ყუთიდან?

აღვნიშნოთ H_1, H_2, H_3 -ით ჰიპოთეზები იმის შესახებ, რომ ალაღბედზე ამოღებული პროდუქტი მიეკუთვნება 1, 2 და 3 ყუთს. მაშინ მათი აპრიორული ალბათობები ერთმანეთის ტოლია

$$P(H_1) = P(H_2) = P(H_3) = \frac{1}{3}.$$

ამოღებული არასტანდარტული პროდუქციის მოხდენის ხდომილობა აღვნიშნოთ A -თი. გამოვთვალოთ A ხდომილობის პირობითი ალბათობები H_1, H_2, H_3 ჰიპოთეზების დროს

$$P(A|H_1) = \frac{3}{10}; \quad P(A|H_2) = \frac{5}{15} = \frac{1}{3}; \quad P(A|H_3) = \frac{3}{10}.$$

ბაიესის ფორმულის გამოყენებით მივიღებთ:

$$P(H_2|A) = \frac{\frac{1}{3} \cdot \frac{1}{3}}{\frac{1}{3} \left(\frac{3}{10} + \frac{1}{3} + \frac{3}{10} \right)} = \frac{5}{14}.$$

1.7. ცდების გამეორება. ბერნულის ფორმულა

ვთქვათ, ვატარებთ ცდებს და ვაკვირდებით რაიმე ხდომილობის მოხდენის ფაქტს. თუ ყოველი ცდისათვის ხდომილობის მოხდენის ალბათობა არ არის დამოკიდებული სხვა ცდების შედეგებზე, მაშინ ისეთ ცდათა მიმდევრობას, რომლის ნებისმიერი ცდის შედეგი გავლენას არ ახდენს მომდევნო ცდების შესაძლო შედეგთა ალბათობაზე, დამოუკიდებელ ცდათა მიმდევრობა ეწოდება.

დავუშვათ, რომ ერთი და იგივე პირობებში ჩატარდა n დამოუკიდებელი ცდა, რომელთა შედეგად შეიძლება მოხდეს A ხდომილობა $P(A) = p$ ალბათობით ან მისი მოპირდაპირე \bar{A} ხდომილობა $P(\bar{A}) = 1 - p$ ალბათობით. ასეთ ცდას ორშედეგიანი, ანუ **ბინარული ცდა** ეწოდება.

განვიხილოთ B ხდომილობა, რომლის დროსაც A ხდომილობა მოხდა m -ჯერ:

$$B = A_1 A_2 \cdots A_m \bar{A}_1 \bar{A}_2 \cdots \bar{A}_n. \quad (1.5)$$

m -ჯერ $(n-m)$ -ჯერ

ცდების დამოუკიდებლობის გამო, (1.5) ფორმულის მარჯვენა მხარეში გვაქვს შეუთავსებლად ხდომილობათა ერთობლიობა. ამიტომ, თუ გამოვიყენებთ ალბათობების ნამრავლის თეორემას, მაშინ

$$P(B) = P(A_1)P(A_2) \cdots P(A_m)P(\bar{A}_{m+1}) \cdots P(\bar{A}_n) = P^m (1 - p)^{n-m}.$$

ასეთი სიდიდე გვექნება თითოეული ცდისათვის. რადგან ცდები ურთიერთდამოუკიდებელია, ამიტომ ალბათობების შეკრების თეორემის თანახმად გვექნება:

$$P_n = C_n^m p^m (1-p)^{n-m} = C_n^m p^m q^{n-m}, \quad (1.6)$$

სადაც, $C_n^m = \frac{n!}{m!(n-m)!}$ არის ჯუფთებათა რიცხვი n ელემენტი-

დან m ელემენტად. მიღებულ (1.6) ფორმულას ბერნულის ფორმულა ეწოდება.

ბერნულის ფორმულა პრაქტიკულად გამოუსადეგარი ხდება, როცა ცდათა n რაოდენობა საკმაოდ დიდია. არსებობს **მუარ-ლაპლასის** მიახლოებითი ფორმულა, რომელსაც აქვს შემდეგი სახე:

$$P_n \approx \frac{1}{\sqrt{npq}} f(x),$$

სადაც,

$$f(x) = \frac{1}{\sqrt{2\pi}} e^{-\frac{x^2}{2}}; \quad x = \frac{m-np}{\sqrt{npq}}.$$

სტანდარტული ნორმალური განაწილების სიმკვრივის $f(x)$ ფუნქციის მნიშვნელობები დადებითი არგუმენტისათვის მოცემულია სპეციალურ ცხრილებში (იხ. დანართი). x -ის უარყოფითი მნიშვნელობის დროს $f(x)$ მნიშვნელობის მოძებნისას უნდა გავითვალისწინოთ ამ ფუნქციის ლუნობა $f(-x) = f(x)$.

მაგალითი. ალბათობა იმისა, რომ ხარატის მიერ დამზადებული დეტალი სტანდარტული იქნება, 0,8-ის ტოლია. ვიპოვოთ იმის ალბათობა, რომ 10000 დამზადებული დეტალიდან 8020 აღმოჩნდება სტანდარტული.

გვაქვს $n = 10000$; $p = 0,8$; $q = 1 - 0,8 = 0,2$; $m = 8020$. მაშინ:

$$P_{10000} = \frac{1}{\sqrt{10000 \cdot 0,8 \cdot 0,2}} f(x) = \frac{1}{40} f(x);$$

$$x = \frac{8020 - 0,8 \cdot 10000}{40} = 0,5; f(0,5) = 0,3521;$$

$$P_{10000} = \frac{1}{40} \cdot 0,3521 = 0,0088$$

2. შემთხვევითი სიდიდეები

2.1. შემთხვევითი სიდიდის ცნება და მისი განაწილების კანონი

შემთხვევითი სიდიდე ეწოდება სიდიდეს, რომელსაც შეუძლია ცდის შედეგად მიიღოს ესა თუ ის მნიშვნელობა, ამასთან, წინასწარ არ არის ცნობილი, კერძოდ, რომელი. შემთხვევითი სიდიდის ცნება წარმოადგენს ალბათობის თეორიაში ფუნდამენტალურ ცნებას და თამაშობს მეტად მნიშვნელოვან როლს მათემატიკურ სტატისტიკაში. შემთხვევითი სიდიდეები აღინიშნებიან დიდი ლათინური ასოებით X, Y, Z, \dots , ხოლო მათი შესაბამისი მნიშვნელობები – მცირე ლათინური ასოებით x, y, z, \dots

შემთხვევითი სიდიდეები შეიძლება დაიყოს ორ კლასად: დისკრეტულად და უწყვეტად. **დისკრეტული** ეწოდება ისეთ შემთხვევით სიდიდეს, რომელიც ლეზულობს სასრულ ან უსასრულო თვლად კონკრეტულ მნიშვნელობებს. მაგალითად, n პარტიაში დეფექტური დეტალების რაოდენობა, სატელეფონო სადგურში გამოძახების რაოდენობა, მიზანში მოხვედრამდე სროლის რაოდენობა და სხვ. **უწყვეტი** ეწოდება ისეთ შემთხვევით სიდიდეს, რომელიც გაზომვის შედეგად ლეზულობს ყველა შესაძლო მნიშვნელობას რალაც სასრულ ან უსასრულო ინტერვალიდან. მაგალითად, ფიზიკური სიდიდის გაზომვის ცდომილება, ჩარხზე დეტალის დამზადების სიზუსტე, ელ. აპარატურაში ტრანზისტორების უტყუარი მუშაობის დრო და სხვ.

შემთხვევითი სიდიდე წარმოადგენს შემთხვევითი ხდომილობების აბსტრაქტულ გამოსახულებას. ამასთან, ადვილი შესაძლებელია ხდომილობიდან შემთხვევით სიდიდეზე გადასვლა. მაგალითად, დაფუშვათ, ტარდება ცდა, რომლის შედეგადაც შეიძლება მოხდეს ან არ მოხდეს A ხდომილობა. A ხდომილობის

ნაცვლად განვიხილოთ X შემთხვევითი სიდიდე, რომელიც უდრის ერთს, თუ A ხდომილობა მოხდა და უდრის ნულს, თუ A ხდომილობა არ მოხდა. ცხადია, ამ შემთხვევაში, X წარმოადგენს დისკრეტულ შემთხვევით სიდიდეს, რომელსაც გააჩნია ორი შესაძლო მნიშვნელობა 0 და 1.

შემთხვევითი სიდიდის დასახასიათებლად არ არის საკმარისი მისი ყველა შესაძლო მნიშვნელობების ჩამონათვალი. აუცილებელია აგრეთვე იმისი ცოდნა, თუ რამდენად ხშირად გვხვდება მისი ესა თუ ის მნიშვნელობა ცდების შედეგად, ე.ი. საჭიროა მოცემული იყოს მისი მოხდენის ალბათობა.

ვთქვათ, X დისკრეტული შემთხვევითი სიდიდეა, რომელიც ცდების შემდეგ იღებს შემდეგ შესაძლო მნიშვნელობებს:

$$X = x_1, X = x_2, X = x_3, \dots, X = x_n. \quad (2.1)$$

შემთხვევითი სიდიდის ეს ჩამონათვალი შეიძლება განვიხილოთ, როგორც ცდის შედეგად მიღებული ხდომილობები. ამ ხდომილობების ალბათობები აღვნიშნოთ P_i -ით. მაშინ გვექნება:

$$P(X = x_1) = p_1, P(X = x_2) = p_2, \dots, P(X = x_n) = p_n.$$

(2.1) ხდომილობები ქმნიან შეუთავსებად ხდომილობათა სრულ ჯგუფს. აქედან გამომდინარე, X შემთხვევითი სიდიდის ყველა შესაძლო მნიშვნელობათა ალბათობების ჯამი ერთის ტოლია, ე.ი.

$$\sum_{i=1}^n P(X = x_i) = \sum_{i=1}^n p_i = 1. \quad (2.2)$$

ეს ჯამური ალბათობა რალაცნაირად განაწილებულია ცალკეულ მნიშვნელობათა შორის. შემთხვევითი სიდიდე სრულად იქნება ალბათურად აღწერილი, თუ მოცემული იქნება ეს განაწილება, ე.ი. თუ ზუსტად მივუთითებთ, რა ალბათობისაა ნებისმიერი ხდომილობა (2.1)-დან. ამით ჩვენ დავადგენთ შემთხვევითი სიდიდის ე.წ. განაწილების კანონს.

შემთხვევითი სიდიდის **განაწილების კანონი** ეწოდება ნებისმიერ თანაფარდობას, რომელიც ამყარებს კავშირს შემთხვევითი სიდიდის შესაძლო მნიშვნელობებსა და მათ შესაბამის ალბათობებს შორის. შემთხვევითი სიდიდის შესახებ ჩვენ ვიტყვით, რომ იგი ემორჩილება მოცემულ განაწილების კანონს.

შემთხვევითი სიდიდის განაწილება შეიძლება მოცემული იყოს ცხრილის სახით, განაწილების ფუნქციის სახით და განაწილების სიმკვრივის სახით.

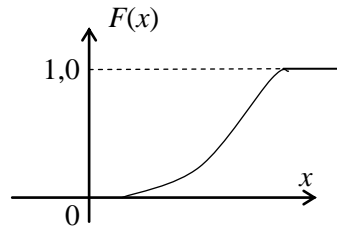
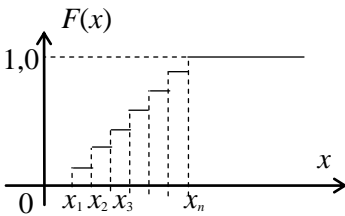
განაწილების ცხრილის სახით წარმოდგენა შესაძლებელია მხოლოდ დისკრეტული შემთხვევითი სიდიდეებისათვის, რომელთაც აქვთ სასრული რაოდენობის მნიშვნელობები. უწყვეტ შემთხვევით სიდიდეს გააჩნია შესაძლო მნიშვნელობათა უსასრულო სიმრავლე. ამიტომ მისთვის ასეთი ცხრილის შედგენა შეუძლებელია. ასეთ შემთხვევაში მოსახერხებელია განაწილების ფუნქციის გამოყენება.

განაწილების ფუნქცია წარმოადგენს შემთხვევითი სიდიდის განაწილების კანონის წარმოდგენის ყველაზე ზოგად სახეს. იგი გამოიყენება როგორც დისკრეტული, ასევე უწყვეტი შემთხვევითი სიდიდეებისათვის. ჩვეულებრივ, იგი აღინიშნება $F(x)$ სიმბოლოთა.

განაწილების ფუნქცია განსაზღვრავს იმის ალბათობას, რომ X შემთხვევითი სიდიდე მიიღებს დაფიქსირებულ x -ზე ნაკლებ მნიშვნელობას, ე.ი.

$$F(x) = P(X < x) .$$

აქედან გამომდინარე, $F(x)$ ფუნქცია დამოკიდებულია x -ზე, ამიტომ, რომ მას განაწილების ფუნქციას უწოდებენ. განაწილების ფუნქცია მიიღება ალბათობების თანმიმდევრული აჯამებით და დისკრეტული შემთხვევითი სიდიდისთვის წარმოადგენს ნყვეტილ საფეხუროვან ტეხილ წირს, ხოლო უწყვეტი შემთხვევითი სიდიდისათვის – უწყვეტ წირს.



განვიხილოთ განაწილების ფუნქციის ზოგადი თვისებები.

1. განაწილების $F(x)$ ფუნქცია დადებითია და (2.2) ფორმულის თანახმად, მოთავსებულია ნულსა და ერთს შორის

$$0 \leq F(x) \leq 1 .$$

2. $F(x)$ ფუნქცია ზრდადია. ე.ი. როცა $x_2 \geq x_1$, მაშინ $F(x_2) \geq F(x_1)$.

3. მიწუს უსასრულობაში განაწილების ფუნქცია ნულის ტოლია $F(-\infty) = 0$.

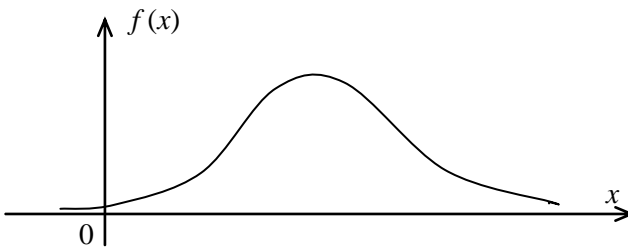
4. პლიუს უსასრულობაში განაწილების ფუნქცია ერთის ტოლია $F(\infty) = 1$.

განაწილების ფუნქციას ზოგჯერ უწოდებენ განაწილების ინტეგრალურ ფუნქციას ან განაწილების ინტეგრალურ კანონს.

უწყვეტი შემთხვევითი სიდიდე შეიძლება წარმოდგენილი იყოს არა მხოლოდ განაწილების ინტეგრალური ფუნქციით, არამედ განაწილების დიფერენციალური ფუნქციითაც, რომელიც წარმოადგენს ინტეგრალური ფუნქციის პირველ წარმოებულს:

$$f(x) = F'(x).$$

ზოგჯერ $f(x)$ ფუნქციას უწოდებენ განაწილების დიფერენციალურ კანონს ან **განაწილების სიმკვრივის ფუნქციას**. მრუდს, რომელიც გამოსახავს შემთხვევითი სიდიდის განაწილების სიმკვრივეს, ეწოდება განაწილების მრუდი.



განვიხილოთ განაწილების სიმკვრივის ფუნქციის თვისებები:

1. განაწილების სიმკვრივე დადებითი ფუნქციაა
 $f(x) > 0$.

2. ინტეგრალი უსასრულო საზღვრებით განაწილების სიმკვრივიდან ერთის ტოლია

$$\int_{-\infty}^{+\infty} f(x) dx = 1.$$

გეომეტრიულად განაწილების სიმკვრივის ძირითადი თვისებები ნიშნავს იმას, რომ იგი ყოველთვის იმყოფება აბსცისათა ღერძის ზემოთ და მის მიერ შემოსაზღვრული სრული ფართობი ერთის ტოლია.

მხედველობაში უნდა ვიქონიოთ, რომ განაწილების ფუნქციას არა აქვს განზომილება, ხოლო განაწილების სიმკვრივის

ფუნქციის განზომილება შემთხვევითი სიდიდის განზომილების შეზღუდვით სიდიდეა.

განაწილების ფუნქცია შეიძლება გამოვსახოთ განაწილების სიმკვრივის ფუნქციით შემდეგი ფორმულით:

$$F(x) = \int_{-\infty}^x f(x)dx.$$

განაწილების ფუნქციის საშუალებით შეიძლება გამოვთვალოთ X შემთხვევითი სიდიდის რაიმე $[\alpha, \beta]$ ინტერვალში მოხვედრის ალბათობა. ამისათვის განვიხილოთ შემდეგი სამი ხდომილობა: ხდომილობა A , როდესაც $X < \beta$, ხდომილობა B , როდესაც $X < \alpha$ და ხდომილობა C , როდესაც $\alpha < X < \beta$. ცხადია, რომ A ხდომილობა წარმოადგენს ორ B და C შეუთავსებად ხდომილობათა ჯამს, ე.ი. $A = B + C$. ალბათობის შეკრების თეორემის თანახმად

$$P(A) = P(B) + P(C).$$

რადგან

$$P(A) = P(X < \beta) = F(\beta), \quad P(B) = P(X < \alpha) = F(\alpha),$$

$$P(C) = P(\alpha < X < \beta),$$

ამიტომ

$$F(\beta) = F(\alpha) + P(\alpha < X < \beta),$$

აქედან

$$P(\alpha < X < \beta) = F(\beta) - F(\alpha).$$

თუ მოცემულია განაწილების სიმკვრივის ფუნქცია, მაშინ გვექნება:

$$F(\beta) = \int_{-\infty}^{\beta} f(x)dx; \quad F(\alpha) = \int_{-\infty}^{\alpha} f(x)dx.$$

$$P(\alpha < X < \beta) = \int_{-\infty}^{\beta} f(x)dx - \int_{-\infty}^{\alpha} f(x)dx = \int_{\alpha}^{\beta} f(x)dx.$$

მაგალითი. მოცემულია X შემთხვევითი სიდიდე $f(x) = 0,5 \sin x$ განაწილების სიმკვრივის ფუნქციით. ვიპოვოთ მისი $[0; \pi/2]$ ინტერვალში მოხვედრის ალბათობა

$$P\left(0 < X < \frac{\pi}{2}\right) = 0,5 \int_0^{\frac{\pi}{2}} \sin x dx = -0,5 \cos x \Big|_0^{\frac{\pi}{2}} = -0,5 \left(\cos \frac{\pi}{2} - \cos 0\right) = 0,5$$

2.2. შემთხვევით სიდიდეთა სისტემა

შემთხვევითი მოვლენების შესწავლისას საკმე გვაქვს არა ერთ, არამედ რამდენიმე შემთხვევით სიდიდესთან, რომლებიც ქმნიან შემთხვევით სიდიდეთა სისტემას. შემთხვევით სიდიდეთა სისტემის განხილვისას ხელსაყრელია გეომეტრიული ინტერპრეტაციის გამოყენება. ასე მაგალითად, ორი შემთხვევითი სიდიდე შეიძლება განვიხილოთ, როგორც შემთხვევითი წერტილი სიბრტყეზე x და y კოორდინატებით, სამი შემთხვევითი სიდიდე, როგორც წერტილი სამგანზომილებიან სივრცეში და ა.შ. ზოგადად, n -რაოდენობის შემთხვევით სიდიდეთა სისტემა შეიძლება განვიხილოთ, როგორც წერტილი n -განზომილებიან სივრცეში ან როგორც n -განზომილებიანი ვექტორი.

შემთხვევით სიდიდეთა სისტემის განხილვისას შემოვიფარგლოთ ორი სიდიდით, რადგან ყველა ის შედეგი, რომელიც მიიღება ორი სიდიდის შემთხვევაში, ადვილად ვრცელდება ნებისმიერი რაოდენობის შემთხვევით სიდიდეებზე. განვიხილოთ ორი შემთხვევითი სიდიდის განაწილების ფუნქცია და განაწილების სიმკვრივის ფუნქცია.

განაწილების ფუნქცია. ორი შემთხვევითი X და Y სიდიდის განაწილების ფუნქცია ეწოდება ორარგუმენტიან $F(x, y)$ ფუნქციას, რომელიც ტოლია ორი $X < x$ და $Y < y$ უტოლობის ერთდროულად შესრულების ალბათობისა. ე.ი.

$$F(x, y) = P(X < x, Y < y)$$

და მას **ერთობლივი განაწილების ფუნქცია** ეწოდება. დაუმტკიცებლად მოვიყვანოთ ერთობლივი განაწილების ფუნქციის ძირითადი თვისებები.

1. $F(x, y)$ ფუნქცია ორივე არგუმენტისათვის ზრდადი ფუნქციაა, ე.ი. როცა $x_2 > x_1$, მაშინ $F(x_2, y) \geq F(x_1, y)$. ასევე, როცა $y_2 > y_1$, მაშინ $F(x, y_2) \geq F(x, y_1)$.

2. თუ $F(x, y)$ ფუნქციის ერთ-ერთი არგუმენტი მიისწრაფვის პლიუს უსასრულობისკენ, მაშინ ერთობლივი განაწილების ფუნქცია მიისწრაფვის მეორე არგუმენტის შესაბამისი შემთხვევითი სიდიდის განაწილების ფუნქციისკენ:

$$\lim_{y \rightarrow \infty} F(x, y) = F(x, \infty) = F_x(x), \quad \lim_{x \rightarrow \infty} F(x, y) = F(\infty, y) = F_y(y),$$

სადაც $F_x(x)$ და $F_y(y)$ უწოდებენ კერძო განაწილების ფუნქციებს.

3. თუ ორივე არგუმენტი მიისწრაფვის პლიუს უსასრულობისკენ, მაშინ ერთობლივი განაწილების ფუნქცია მიისწრაფვის ერთისკენ.

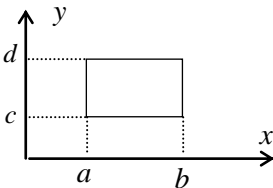
$$\lim_{x, y \rightarrow \infty} F(x, y) = 1 \quad \text{ან} \quad F(\infty, \infty) = 1.$$

4. თუ რომელიმე ერთი ან ორივე არგუმენტი მიისწრაფვის მინუს უსასრულობისკენ, მაშინ ერთობლივი განაწილების ფუნქცია მიისწრაფვის ნულისკენ.

$$\lim_{x \rightarrow -\infty} F(x, y) = \lim_{y \rightarrow -\infty} F(x, y) = \lim_{x, y \rightarrow -\infty} F(x, y) = 0,$$

$$\text{ან} \quad F(-\infty, y) = F(x, -\infty) = F(-\infty, -\infty) = 0.$$

5. კოორდინატთა ღერძების პარალელურგვერდებიან ნებისმიერ ოთკუთხედში ნერტილის მოხვედრის ალბათობა განისაზღვრება ფორმულით:



$$P(a \leq X < b, c \leq Y < d) =$$

$$= F(b, d) - F(a, d) - F(b, c) + F(a, c).$$

სიმკვრივის ფუნქცია. ზემოთ განხილული განაწილების ფუნქცია წარმოადგენს შემთხვევით სიდიდეთა სისტემის უნივერსალურ მახასიათებელს, რომელიც გამოიყენება როგორც დისკრეტული, ასევე უწყვეტი შემთხვევითი სიდიდეებისთვის. უნდა აღინიშნოს, რომ პრაქტიკული გამოყენება უფრო გააჩნია უწყვეტ შემთხვევით სისტემას, რომლის განაწილება ხასიათდება არა განაწილების ფუნქციით, არამედ განაწილების სიმკვრივით. განაწილების სიმკვრივის ფუნქცია წარმოადგენს სისტემის ამომწურავ მახასიათებელს, რომლის საშუალებითაც სისტემის

განაწილების კანონის აღწერა გაცილებით თვალსაჩინოა, ვიდრე განაწილების ფუნქციით.

ორგანზომილებიანი სისტემის განაწილების სიმკვრივის ფუნქცია განისაზღვრება ისევე, როგორც ერთი შემთხვევითი სიდიდის დროს. კერძოდ, თუ ერთობლივი განაწილების ფუნქცია უწყვეტია და ორჯერ დიფერენცირებადი, მაშინ ერთობლივი განაწილების სიმკვრივის ფუნქცია განისაზღვრება შემდეგნაირად:

$$f(x, y) = \frac{\partial^2 F(x, y)}{\partial x \partial y} = F''(x, y).$$

განვიხილოთ ერთობლივი განაწილების სიმკვრივის ფუნქციის ძირითადი თვისებები:

1. ერთობლივი განაწილების სიმკვრივე დადებითი ფუნქციაა $f(x, y) \geq 0$.

2. ორმაგი ინტეგრალი უსასრულო ზღვრებით ერთობლივი განაწილების სიმკვრივის ფუნქციიდან ერთის ტოლია.

$$\int_{-\infty}^{\infty} \int_{-\infty}^{\infty} f(x, y) dx dy = 1.$$

3. თუ ცნობილია ერთობლივი განაწილების სიმკვრივის ფუნქცია $f(x, y)$, მაშინ შემთხვევითი წერტილის (X, Y) ნებისმიერ D არეში მოხვედრის ალბათობა განისაზღვრება ფორმულით:

$$P[(X, Y) \in D] = \iint_D f(x, y) dx dy.$$

ერთობლივი განაწილების ფუნქცია შეიძლება გამოვსახოთ განაწილების სიმკვრივის ფუნქციით შემდეგნაირად:

$$F(x, y) = \int_{-\infty}^x \int_{-\infty}^y f(x, y) dx dy.$$

გეომეტრიულად $f(x, y)$ ფუნქცია შეიძლება წარმოვადგინოთ როგორც რაიმე ზედაპირი, რომელსაც განაწილების ზედაპირი ეწოდება.

მაგალითი. მოცემულია ორგანზომილებიანი სისტემის განაწილების სიმკვრივის ფუნქცია

$$f(x, y) = \frac{a}{1 + x^2 + x^2 y^2 + y^2}.$$

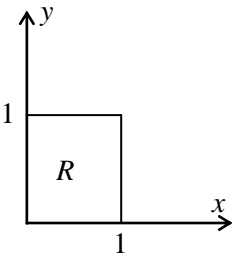
ვიპოვოთ a , ერთობლივი განაწილების ფუნქცია $f(x, y)$ და R კვადრატში შემთხვევითი წერტილის მოხვედრის ალბათობა.

ერთობლივი განაწილების სიმკვრივის მეორე თვისებიდან გამომდინარე გვქვია:

$$\begin{aligned} \int_{-\infty}^{\infty} \int_{-\infty}^{\infty} \frac{a}{1+x^2+x^2y^2+y^2} dx dy &= a \int_{-\infty}^{\infty} \int_{-\infty}^{\infty} \frac{dx dy}{(1+x^2)(1+y^2)} = \\ &= a \int_{-\infty}^{\infty} \frac{dx}{1+x^2} \int_{-\infty}^{\infty} \frac{dy}{1+y^2} = a \cdot \arctg x \Big|_{-\infty}^{\infty} \cdot \arctg y \Big|_{-\infty}^{\infty} = a\pi^2 = 1. \end{aligned}$$

აქედან $a = \frac{1}{\pi^2}$ ერთობლივი განაწილების ფუნქცია ტოლია:

$$F(x, y) = \frac{1}{\pi^2} \int_{-\infty}^x \int_{-\infty}^y \frac{dx dy}{(1+x^2)(1+y^2)} = \left(\frac{1}{\pi} \arctg x + \frac{1}{2} \right) \left(\frac{1}{\pi} \arctg y + \frac{1}{2} \right).$$



შემთხვევითი წერტილის R კვადრატში მოხვედრის ალბათობა ტოლია:

$$P[(X, Y) \in D] = \frac{1}{\pi^2} \int_0^1 \int_0^1 \frac{dx dy}{(1+x^2)(1+y^2)} =$$

$$= \frac{1}{\pi^2} \int_0^1 \frac{dx}{1+x^2} \int_0^1 \frac{dy}{1+y^2} = \frac{1}{\pi^2} \arctg x \Big|_0^1 \cdot \arctg y \Big|_0^1 = \frac{1}{\pi^2} \frac{\pi}{4} \frac{\pi}{4} = \frac{1}{16}.$$

2.3. პირობითი განაწილების სიმკვრივის ფუნქცია

შემთხვევით სიდიდეთა სისტემის დახასიათებისთვის არ არის საკმარისი ვიცოდეთ თითოეული შემთხვევითი სიდიდის განაწილების კანონი, საჭიროა ვიცოდეთ მათ შორის დამოკი-

დებულება ც. ეს დამოკიდებულება შეიძლება დახასიათდეს ე.წ. პირობითი განაწილების კანონით.

სისტემის ერთი, X შემთხვევითი სიდიდის განაწილების კანონს, გამოთვლილს იმ პირობით, რომ მეორე, Y შემთხვევითმა სიდიდემ მიიღო გარკვეული $Y = y$ მნიშვნელობა, ეწოდება პირობითი განაწილების კანონი. იგი შეიძლება წარმოვადგინოთ, როგორც პირობითი განაწილების სიმკვრივის ფუნქცია $f(x | y)$, ან როგორც პირობითი განაწილების ფუნქცია $F(x | y)$. ერთობლივი განაწილების ფუნქციის თვისებიდან გამომდინარე, $F_1(x) = F(x, \infty)$ და $F_2(y) = F(\infty, y)$. აქედან გამომდინარე, კერძო განაწილების სიმკვრივის ფუნქციები ტოლია:

$$f_x(x) = F_1'(x) = \int_{-\infty}^{\infty} f(x, y) dy,$$

$$f_y(y) = F_2'(y) = \int_{-\infty}^{\infty} f(x, y) dx.$$

ადვილად მტკიცდება, რომ სისტემის ერთობლივი განაწილების სიმკვრივის ფუნქცია $f(x, y)$ მიიღება სისტემაში შემავალი ერთი შემთხვევითი სიდიდის კერძო განაწილების სიმკვრივისა $f_x(x)$ და მეორე შემთხვევითი სიდიდის პირობითი განაწილების სიმკვრივის $f(y | x)$ ერთმანეთზე გადამრავლებით. ე.ი.

$$f(x, y) = f_x(x)f(y | x) \text{ და } f(x, y) = f_y(y)f(x | y).$$

ამ ტოლობებს ხშირად განაწილების კანონების გამრავლების თეორემას უწოდებენ. პირობითი განაწილების სიმკვრივის ფუნქციები ასე განისაზღვრება:

$$f(x | y) = \frac{f(x, y)}{f_y(y)} = \frac{f(x, y)}{\int_{-\infty}^{\infty} f(x, y) dx},$$

$$f(y | x) = \frac{f(x, y)}{f_x(x)} = \frac{f(x, y)}{\int_{-\infty}^{\infty} f(x, y) dy}.$$

პირობითი განაწილების სიმკვრივეს ახასიათებს იგივე თვისებები, რაც ჩვეულებრივი განაწილების სიმკვრივის ფუნქციას, კერძოდ:

$$\int_{-\infty}^{\infty} f(x|y)dx = \int_{-\infty}^{\infty} f(y|x)dy = 1.$$

პირობითი განაწილების სიმკვრივის ცნებიდან გამომდინარე, შეგვიძლია შემოვიტანოთ ალბათობის თეორიაში ერთ-ერთი უმნიშვნელოვანესი ცნება – შემთხვევით სიდიდეთა დამოუკიდებლობის ცნება.

X შემთხვევითი სიდიდის დამოუკიდებლობა Y შემთხვევით სიდიდესთან ნებისმიერი y -თვის შეიძლება ასე ჩაინეროს:

$$f(x/y) = f_x(x).$$

თუ შემთხვევითი სიდიდე X დამოკიდებულია Y -ზე, მაშინ:

$$f(x/y) \neq f_x(x).$$

თუ X სიდიდე არ არის დამოკიდებული Y -ზე, მაშინ არც Y სიდიდეა დამოკიდებული X -ზე.

ამრიგად, უწყვეტ X და Y შემთხვევით სიდიდეებს ეწოდებათ დამოუკიდებელი, თუ თითოეული მათგანის განაწილების კანონი არ არის დამოკიდებული იმაზე, თუ რა მნიშვნელობა მიიღო მეორე შემთხვევითმა სიდიდემ. წინააღმდეგ შემთხვევაში, შემთხვევითი სიდიდეები დამოკიდებულია. დამოუკიდებელი X და Y შემთხვევითი სიდიდეთა ერთობლივი განაწილების სიმკვრივის ფუნქცია ტოლია:

$$f(x, y) = f_x(x) f_y(y).$$

მაგალითი. ვთქვათ, მოცემულია ერთობლივი განაწილების სიმკვრივის ფუნქცია:

$$f(x, y) = \frac{1}{\pi^2(x^2 + x^2y^2 + y^2 + 1)}.$$

ეს ტოლობა წარმოვადგინოთ შემდეგნაირად:

$$f(x, y) = \frac{1}{\pi(x^2 + 1)} \cdot \frac{1}{\pi(y^2 + 1)}.$$

აქედან ჩანს, რომ შემთხვევითი სიდიდეები X და Y დამოუკიდებელია, რადგან $f(x, y)$ გამოსახულების პირველი მამრავლი დამოკიდებულია მხოლოდ X -ზე, ხოლო მეორე მამრავლი – Y -ზე. ე.ი. $f(x, y) = f_x(x) f_y(y)$.

2.4. შემთხვევითი სიდიდის რიცხვითი მახასიათებლები

ჩვენ ვიცით, რომ განაწილების კანონი შემთხვევით სიდიდეს სრულად ახასიათებს. მაგრამ, ძალიან ხშირად, განაწილების კანონი უცნობია. ამიტომ გაცილებით უფრო მოსახერხებელია ზოგიერთი რაოდენობრივი მაჩვენებლების განსაზღვრა, რომლებიც შეკუმშული ფორმით მოგვანვდიან ინფორმაციას შემთხვევით სიდიდეზე. ასეთ მაჩვენებლებს უწოდებენ შემთხვევითი სიდიდის რიცხვით მახასიათებლებს. მათგან ძირითადია მათემატიკური ლოდინი, დისპერსია, სხვადასხვა რიგის მომენტები, მოდა და მედიანა.

მათემატიკური ლოდინი ახასიათებს შემთხვევით სიდიდეს რიცხვით ღერძზე და განსაზღვრავს მის განაწილების ცენტრს, ამიტომ ხშირად მათემატიკურ ლოდინს უწოდებენ შემთხვევითი სიდიდის საშუალო მნიშვნელობას. მათემატიკური ლოდინი აღინიშნება M ან E სიმბოლოთი. განვიხილოთ შემთხვევითი სიდიდე X , რომელიც იღებს x_1, x_2, \dots, x_n მნიშვნელობებს p_1, p_2, \dots, p_n ალბათობით, მაშინ მათემატიკური ლოდინი ტოლია:

$$M(X) = \frac{x_1 p_1 + x_2 p_2 + \dots + x_n p_n}{p_1 + p_2 + \dots + p_n}.$$

თუ გავითვალისწინებთ, რომ $p_1 + p_2 + \dots + p_n = 1$, მაშინ

$$M(X) = \sum_{i=1}^n x_i p_i.$$

ამრიგად, დისკრეტული შემთხვევითი სიდიდის მათემატიკური ლოდინი ეწოდება მისი ყველა შესაძლო მნიშვნელობების მათსავე ალბათობებზე ნამრავლის ჯამს. უწყვეტი შემთხვევითი სიდიდისათვის გვაქვს:

$$M(x) = \int_{-\infty}^{\infty} x f(x) dx.$$

მოვიყვანოთ დაუმტკიცებლად მათემატიკური ლოდინის ზოგიერთი თვისებები:

1. მუდმივი სიდიდის მათემატიკური ლოდინი უდრის თვით ამ სიდიდეს

$$M(c) = c, \quad c = \text{const.}$$

2. შემთხვევითი სიდიდის მუდმივი მამრავლი შეიძლება გატანილ იქნეს მათემატიკური ლოდინის აღნიშვნის გარეთ

$$M(cX) = cM(X).$$

3. ორი შემთხვევითი სიდიდის ჯამის მათემატიკური ლოდინი ტოლია მათი მათემატიკური ლოდინების ჯამისა

$$M(X + Y) = M(X) + M(Y).$$

შედეგი. რამდენიმე შემთხვევითი სიდიდის ჯამის მათემატიკური ლოდინი ტოლია მათი მათემატიკური ლოდინების ჯამისა

$$M(X + Y + Z + \dots + N) = M(X) + M(Y) + M(Z) + \dots + M(N).$$

4. ორი დამოუკიდებელი შემთხვევითი სიდიდის ნამრავლის მათემატიკური ლოდინი ტოლია მათი მათემატიკური ლოდინების ნამრავლისა

$$M(X \cdot Y) = M(X) \cdot M(Y).$$

შედეგი. რამდენიმე ურთიერთდამოუკიდებელი შემთხვევითი სიდიდის ნამრავლის მათემატიკური ლოდინი ტოლია მათი მათემატიკური ლოდინების ნამრავლისა

$$M(X \cdot Y \cdot Z \cdot \dots \cdot N) = M(X) \cdot M(Y) \cdot M(Z) \cdot \dots \cdot M(N).$$

5. შემთხვევითი სიდიდის თავის მათემატიკურ ლოდინთან გადახრის მათემატიკური ლოდინი ნულის ტოლია

$$M[X - M(X)] = 0.$$

შემთხვევითი სიდიდის მდებარეობის სხვა მახასიათებლებია მოდა და მედიანა.

მოდა M_0 ეწოდება დისკრეტული შემთხვევითი სიდიდის იმ მნიშვნელობას, რომელსაც გააჩნია უდიდესი ალბათობა. უწყვეტი შემთხვევითი სიდიდისთვის მოდა ეწოდება მის იმ მნიშვნელობას, რომლის დროსაც განაწილების სიმკვრივე მაქსიმალურია. არსებობს ორმოდიანი და მრავალმოდიანი განაწილებები. გვხვდება ისეთი განაწილებები, რომელთაც აქვთ მინიმუმი და არ გააჩნიათ მაქსიმუმი. ასეთ განაწილებებს უწოდებენ ანტიმოდალურს.

მედიანა M_e არის იმ წერტილის აბსცისა, რომელზედაც განაწილების მრუდით შემოსაზღვრული ფართობი იყოფა შუაზე. უნდა აღინიშნოს, რომ თუ განაწილება ერთმოდიანია და სიმეტრიული, მაშინ მათემატიკური ლოდინი, მოდა და მედიანა ერთმანეთს ემთხვევა.

დისპერსია და საშუალო კვადრატული გადახრა. ამ მახასიათებლებით შეიძლება ვიმსჯელოთ შემთხვევითი სიდიდის გაფანტვის შესახებ მისი მათემატიკური ლოდინის მიმართ. შემ-

თხვევითი სიდიდის დისპერსია ეწოდება X შემთხვევითი სიდიდის მისი მათემატიკურ ლოდინთან გადახრის კვადრატს.

$$D(X) = M[X - M(X)]^2 .$$

დისკრეტული შემთხვევითი სიდიდეებისათვის გვექნება:

$$D(X) = \sum_{i=1}^n [x_i - M(X)]^2 p_i ,$$

ხოლო უწყვეტი შემთხვევითი სიდიდეებისათვის:

$$D(X) = \int_{-\infty}^{+\infty} [x - M(X)]^2 f(x) dx .$$

დისპერსიის გამოსათვლელად უფრო მიზანშეწონილია შემდეგი ფორმულის გამოყენება:

$$\begin{aligned} D(X) &= M[X - M(X)]^2 = M[X^2 - 2X M(X) + M^2(X)] = M(X^2) - \\ &- 2M(X)M(X) + M^2(X) = M(X^2) - 2M^2(X) + M^2(X) = \\ &= M(X^2) - M^2(X) . \end{aligned}$$

დისპერსიის ნაკლი ისაა, რომ მისი განზომილება შემთხვევითი სიდიდის განზომილების კვადრატის ტოლია და გაფანტვის დახასიათებისთვის ერთგვარად უხერხულობას ქმნის. ამ ნაკლისაგან თავისუფალია საშუალო კვადრატული გადახრა, რომელიც წარმოადგენს დადებით კვადრატულ ფესვს დისპერსიიდან

$$s_x = \sqrt{D(X)} .$$

განვიხილოთ დისპერსიის ზოგიერთი თვისებები:

1. მუდმივი სიდიდის დისპერსია ნულის ტოლია

$$D(c) = 0 , \quad c = \text{const.}$$

2. შემთხვევითი სიდიდის მუდმივი მამრავლი შეიძლება გატანილ იქნეს დისპერსიის ნიშნის გარეთ, რომელიც კვადრატში უნდა იყოს აყვანილი

$$D(cX) = c^2 D(X) .$$

3. ორი დამოუკიდებელი შემთხვევითი სიდიდის ჯამის დისპერსია ტოლია ამ სიდიდეთა დისპერსიების ჯამისა

$$D(X + Y) = D(X) + D(Y) .$$

შედეგი. რამდენიმე ურთიერთდამოუკიდებელი შემთხვევითი სიდიდის ჯამის დისპერსია ამ სიდიდეთა დისპერსიების ჯამის ტოლია.

$$D(X + Y + \dots + N) = D(X) + D(Y) + \dots + D(N) .$$

4. ორი დამოუკიდებელი შემთხვევითი სიდიდის სხვაობის დისპერსია ამ სიდიდეთა დისპერსიების ჯამის ტოლია

$$D(X - Y) = D(X) + D(Y) .$$

მართლაც, $D(X - Y) = D(X) + D(-Y) = D(X) + (-1)^2 D(Y) = D(X) + D(Y)$.

5. ორი დამოუკიდებელი შემთხვევითი სიდიდის ნამრავლის დისპერსია განისაზღვრება შემდეგი ფორმულით:

$$D(XY) = D(X) D(Y) + M^2(X) D(Y) + M^2(Y) D(X) .$$

მაგალითი. ვთქვათ, მოცემულია X შემთხვევითი სიდიდის განაწილების სიმკვრივის ფუნქცია:

$$f(x) = \begin{cases} 0, & \text{როცა } x < 0 \\ ax^2, & \text{როცა } 0 < x < 2 \\ 0, & \text{როცა } x > 2 \end{cases}$$

ვიპოვოთ a კოეფიციენტი, განაწილების ფუნქცია, მათემატიკური ლოდინი, დისპერსია და საშუალო კვადრატული გადახრა. რადგან,

$$\int_0^2 f(x) dx = 1,$$

ამიტომ

$$\int_0^2 f(x) dx = \int_0^2 ax^2 dx = \frac{a}{3} x^3 \Big|_0^2 = \frac{8}{3} a = 1, \text{ აქედან } a = \frac{3}{8}$$

$$\text{და } f(x) = \frac{3}{8} x^2. F(x) = \int_0^x f(x) dx = \frac{3}{8} \int_0^x x^2 dx = \frac{3}{8} \frac{x^3}{3} \Big|_0^x = \frac{x^3}{8}$$

განვსაზღვროთ მათემატიკური ლოდინი:

$$M(x) = \int_0^2 xf(x)dx = \frac{3}{8} \int_0^2 x^3 dx = \frac{3}{8} \frac{x^4}{4} \Big|_0^2 = 1,5$$

დისპერსია და საშუალო კვადრატული გადახრა:

$$D(X) = M(X^2) - M^2(X).$$

$$M(X^2) = \int_0^2 x^2 f(x)dx = \frac{3}{8} \int_0^2 x^4 dx = \frac{3}{8} \frac{x^5}{5} \Big|_0^2 = 2,4;$$

$$D(X) = 2,4 - (1,5)^2 = 0,15; \quad s_x = \sqrt{0,15} = 0,39.$$

შემთხვევითი სიდიდის ძირითად რიცხვით მახასიათებელთა განზოგადებას წარმოადგენს შემთხვევითი სიდიდის მომენტის ცნება. არჩევენ ორი სახის მომენტს: საწყისს და ცენტრალურს. საწყისი მომენტებიდან განსაკუთრებული მნიშვნელობა აქვს პირველი რიგის მომენტს, რომელიც წარმოადგენს მათემატიკურ ლოდინს. მაღალი რიგის საწყისი მომენტები, ძირითადად, გამოიყენება ცენტრალური მომენტების გამოსათვლელად. დისკრეტული შემთხვევითი სიდიდისათვის k -ური რიგის ცენტრალური მომენტი გამოითვლება ფორმულით:

$$\mu_k = \sum_{i=1}^n [x_i - M(X)]^k p_i ,$$

ხოლო უწყვეტი შემთხვევითი სიდიდისათვის გვექნება:

$$\mu_k = \int_{-\infty}^{+\infty} [x - M(X)]^k f(x)dx .$$

პირველი რიგის ცენტრალური მომენტი, როგორც ფორმულიდან ჩანს, ნულის ტოლია, ხოლო მეორე რიგის მომენტი წარმოადგენს დისპერსიას.

მესამე რიგის ცენტრალური მომენტი ახასიათებს განაწილების ასიმეტრიას. თუ შემთხვევითი სიდიდე თავისი მათემატიკური ლოდინის მიმართ სიმეტრიულადაა განაწილებული, მაშინ მესამე რიგის ცენტრალური მომენტი ნულის ტოლია. ამ მომენტის საშუალებით გამოითვლება ასიმეტრიის კოეფიციენტი

$$A_x = \frac{\mu_3}{s_x^3},$$

სადაც, s_x^3 – საშუალო კვადრატული გადახრაა აყვანილი კუბში. განაწილების მრუდს გააჩნია დადებითი ასიმეტრია, როცა $A_x > 0$ და უარყოფითი, როცა $A_x < 0$.

მეოთხე რიგის ცენტრალური მომენტი გამოიყენება განაწილების მრუდის წამახვილების ხარისხის დასახასიათებლად. ეს თვისება აღინიშნება ე.წ. ექსცესის კოეფიციენტის საშუალებით, რომელიც გამოითვლება შემდეგი ფორმულით:

$$E_x = \frac{\mu_4}{s_x^4} - 3,$$

სადაც, s_x^4 – საშუალო კვადრატული გადახრაა აყვანილი მეოთხე ხარისხში.

მიღებულია, რომ ნორმალურად განაწილებული მრუდისთვის ექსცესის კოეფიციენტი ნულის ტოლია და ეს შედეგი მიღებულია ეტალონად, რომელთანაც შედარდებიან სხვა განაწილების მრუდეები. მრუდს, რომელსაც უფრო მაღალი წვერო აქვს, ვიდრე ნორმალურს, ე.ი. უფრო მახვილწვეროიანია, შეესაბამება დადებითი ექსცესა, ხოლო მრუდს, რომელსაც უფრო დაბალი და ბრტყელი წვერო აქვს – უარყოფითი ექსცესა.

2.5. ორგანოზომილუბიანი შემთხვევითი სისტემის რიცხვითი მახასიათებლები

ორი X და Y შემთხვევით ვექტორთა (k, s) რიგის საწყისი მომენტი ეწოდება $X^k Y^s$ სიდიდეთა ნამრავლის მათემატიკურ ლოდინს და აღინიშნება $a_{k,s}$ სიმბოლოთი.

$$a_{ks} = M[X^k Y^s]$$

დისკრეტულ შემთხვევით სიდიდეთა სისტემისათვის გვექნება:

$$a_{ks} = \sum_i \sum_j x_i^k y_j^s p_{ij},$$

სადაც, $p_{ij} = P(X = x_i, Y = y_j)$.

უნყვეტ შემთხვევით სიდიდეთა სისტემისთვის გვექნება:

$$a_{ks} = \int_{-\infty}^{\infty} \int_{-\infty}^{\infty} x^k y^s f(x, y) dx dy.$$

პრაქტიკაში ხშირად გამოიყენება პირველი რიგის სანყისი მომენტები, რომლებიც X და Y შემთხვევითი სიდიდეების მათემატიკურ ლოდინებს წარმოადგენენ.

$$a_{10} = M[X^1 Y^0] = M[X] = m_x,$$

$$a_{01} = M[X^0 Y^1] = M[Y] = m_y.$$

(k, s) რიგის ცენტრალური მომენტი ეწოდება $(X - m_x)^k (Y - m_y)^s$ ნამრავლის მათემატიკურ ლოდინს და აღინიშნება μ_{ks} სიმბოლოთი.

$$\mu_{ks} = M[(X - m_x)^k (Y - m_y)^s].$$

დისკრეტული შემთხვევითი სიდიდეებისათვის გვექნება:

$$\mu_{ks} = \sum_i \sum_j (x_i - m_x)^k (y_j - m_y)^s p_{ij},$$

ხოლო უწყვეტი შემთხვევითი სიდიდეებისათვის:

$$\mu_{ks} = \int_{-\infty}^{\infty} \int_{-\infty}^{\infty} (x - m_x)^k (y - m_y)^s f(x, y) dx dy.$$

პრაქტიკაში ყველაზე ხშირად გამოიყენება მეორე რიგის ცენტრალური მომენტი, რომელიც დისპერსიას წარმოადგენს.

$$D(x) = \mu_{20} = M[(X - m_x)^2 (Y - m_y)^0] = M[(X - m_x)^2].$$

$$D(y) = \mu_{02} = M[(X - m_x)^0 (Y - m_y)^2] = M[(Y - m_y)^2].$$

პრაქტიკულ კვლევებში განსაკუთრებულ როლს თამაშობს მეორე რიგის შერეული ცენტრალური მომენტი μ_{11} , რომელსაც X და Y შემთხვევითი სიდიდეთა **კორელაციურ მომენტს** ან კავშირის მომენტს უწოდებენ და მას k_{xy} სიმბოლოთი აღნიშნავენ.

$$k_{xy} = \mu_{11} = M[(X - m_x)(Y - m_y)].$$

დისკრეტული შემთხვევითი სიდიდეებისათვის გვექნება:

$$k_{xy} = \sum_i \sum_j (x_i - m_x)(y_j - m_y)p_{ij},$$

ხოლო უწყვეტი შემთხვევითი სიდიდეებისათვის:

$$k_{xy} = \int_{-\infty}^{\infty} \int_{-\infty}^{\infty} (x - m_x)(y - m_y)f(x, y)dx dy.$$

კორელაციური მომენტი ახასიათებს ორი შემთხვევითი სიდიდის კავშირს. კავშირის ხარისხის შეფასებისთვის გამოიყენება არა თვით კორელაციური მომენტი, არამედ უგანზომილებო სიდიდე

$$r_{xy} = \frac{k_{xy}}{s_x s_y},$$

რომელსაც **კორელაციის კოეფიციენტი** ეწოდება. s_x და s_y ნარმოდგენენ საშუალო კვადრატულ გადახრებს.

ადვილად მტკიცდება, რომ დამოუკიდებელ შემთხვევით სიდიდეთა კორელაციური მომენტი და კორელაციის კოეფიციენტი ნულის ტოლია. წინააღმდეგ შემთხვევაში, შემთხვევითი სიდიდეები დამოკიდებულნი ანუ კორელირებულნი არიან.

3. შემთხვევითი სიდიდეების ძირითადი განაწილების კანონები

3.1. ბინომური განაწილების კანონი

ბინომური კანონით განაწილებულია ძირითადად თვისებრივი პარამეტრები. ვთქვათ, A და B ხდომილობებს აქვთ მხოლოდ დიქტომიური შედეგი (მაგ. „კი“, „არა“, „+“, „-“). A ხდომილობის ალბათობა $P(A) = p$ ყოველი ცდისათვის იყოს მუდმივი. აქედან გამომდინარე, B ხდომილობის ალბათობაც $P(B) = q$ იქნება მუდმივი სიდიდე. ცხადია, $p + q = 1$. ვთქვათ, ჩატარდა n ცდა და A ხდომილობა მოხდა m -ჯერ, მაშინ B ხდომილობა მოხდებოდა $(n - m)$ -ჯერ. შედეგის ხდომილობის ალბათობა გამოითვლება ბერნულის ფორმულით:

$$P_n(m) = C_n^m p^m q^{n-m} = \frac{n!}{m!(n-m)!} p^m (1-p)^{n-m}.$$

რადგან $P_n(m)$ ალბათობა დაკავშირებულია ბინომის C_n^m გაშლასთან, ამიტომ m შემთხვევითი სიდიდის ამ განაწილებას ბინომური განაწილება ეწოდება. იგი წარმოადგენს დისკრეტული შემთხვევითი სიდიდის განაწილებას, რადგან m იღებს გარკვეულ მთელ მნიშვნელობებს $m = 0, 1, 2, 3, \dots, n$. ბინომურად არის განაწილებული არა მარტო ალტერნატიული მაჩვენებლები, არამედ ისეთი მაჩვენებლებიც, რომლებთაც გააჩნიათ არა ორი, არამედ უფრო მეტი შედეგი.

ბინომური კანონით განაწილებული X შემთხვევითი სიდიდის მათემატიკური ლოდინი ტოლია $M(X) = np$, ხოლო საშუალო კვადრატული გადახრა – $s_p = \sqrt{npq}$.

ბინომური განაწილების სიმკვრივის გრაფიკი დამოკიდებულია ცდების n რაოდენობაზე და მოსალოდნელი შედეგის p ალბათობაზე. როცა $p = 0,5$, ბინომური მრუდი სიმეტრიულია. თუ $p \neq q$, მაშინ ბინომურ მრუდს გააჩნია ასიმეტრია, რომელიც $|p - q|$ სხვაობის ზრდასთან ერთად იზრდება.

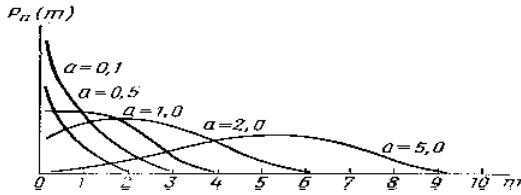
3.2. პუასონის განაწილება

როდესაც მოსალოდნელი შედეგის ალბათობა ძალზედ მცირეა, მაგალითად, ერთის მეასედი ან მეათასედი, მაშინ განაწილების მრუდი ძალიან ასიმეტრიულია. ასეთი იშვიათი ხდომილობის განაწილების სიხშირე შეიძლება გამოითვალოს პუასონის ფორმულით:

$$P_n(m) = \frac{a^m}{m!} e^{-a},$$

სადაც, m არის n დამოუკიდებელი ცდის მოსალოდნელი ხდომილობათა სიხშირე, a – პუასონის განაწილების პარამეტრია ($a > 0$).

პუასონის განაწილების სიმკვრივის გრაფიკები, რომლებიც a სიდიდეზეა დამოკიდებული, ნაჩვენებია შემდეგ ნახაზზე:



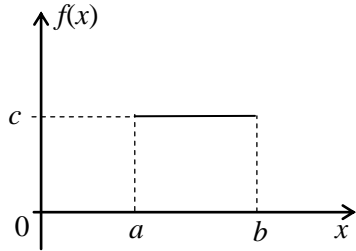
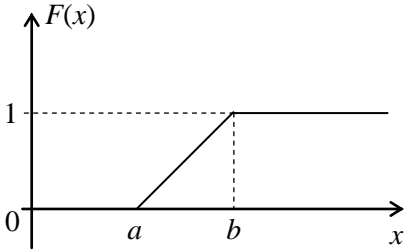
ირკვევა, რომ პუასონის კანონით განაწილებული X შემთხვევითი სიდიდის მათემატიკური ლოდინი და დისპერსია პუასონის განაწილების a პარამეტრის ტოლია, ე.ი. $M(X) = D(X) = a$. პუასონის განაწილების ეს თვისება ხშირად გამოიყენება პრაქტიკაში X შემთხვევითი სიდიდის პუასონის განაწილებაზე ჰიპოთეზის შესამოწმებლად. ამისათვის საჭიროა გამოვთვალოთ შემთხვევითი სიდიდის მათემატიკური ლოდინი და დისპერსია. თუ მათი მნიშვნელობები ერთმანეთის ტოლია ან ახლოს არიან, მაშინ მიახლოებით შეგვიძლია ჩავთვალოთ, რომ ამ შემთხვევით სიდიდეს გააჩნია პუასონის განაწილება.

3.3. თანაბარი განაწილება

უნყვეტ X შემთხვევით სიდიდეს $[a, b]$ ინტერვალში აქვს თანაბარი განაწილება (თანაბრადაა განაწილებული), თუ მისი განაწილების სიმკვრივე ამ ინტერვალში მუდმივია, ხოლო მის გარეთ ნულის ტოლია, ე.ი.

$$f(x) = \begin{cases} 0, & \text{როცა } x < a, \\ c, & \text{როცა } a \leq x \leq b, \\ 0, & \text{როცა } x > b, \quad c = \text{const.} \end{cases}$$

განაწილების სიმკვრივეს და განაწილების ფუნქციას აქვთ შემდეგი სახე:



ვიზოვით c მუდმივია. რადგან განაწილების მრუდით შემოსაზღვრული ფართობი ერთის ტოლია, ამიტომ

$$\int_a^b f(x)dx = \int_a^b cdx = 1. \text{ აქედან } c = \frac{1}{b-a}. \text{ ე.ი. } f(x) = \frac{1}{b-a}.$$

განაწილების ფუნქცია ტოლია:

$$F(x) = \int_a^x f(x)dx = \int_a^x \frac{1}{b-a} dx = \frac{x}{b-a} \Big|_a^x = \frac{x-a}{b-a}, \text{ ე.ი. } F(x) = \frac{x-a}{b-a}.$$

მათემატიკური ლოდინი:

$$M(x) = \int_a^b xf(x)dx = \int_a^b \frac{x}{b-a} dx = \frac{1}{b-a} \frac{x^2}{2} \Big|_a^b = \frac{b^2 - a^2}{2(b-a)} = \frac{a+b}{2},$$

ბოლო დისპერსია:

$$D(x) = \int_a^b \left[x - \frac{a+b}{2} \right]^2 \frac{1}{b-a} dx = \frac{1}{b-a} \frac{1}{3} \left[x - \frac{a+b}{2} \right]^3 \Big|_a^b = \frac{(b-a)^2}{12};$$

$$s = \frac{b-a}{2\sqrt{3}}.$$

განაწილების სიმეტრიის გამო ასიმეტრიის კოეფიციენტი ნულის ტოლია. ვიზოვით ექსცესის კოეფიციენტი.

$$\mu_4 = \int_a^b \left(x - \frac{a+b}{2} \right)^4 \frac{1}{b-a} dx = \frac{(b-a)^4}{80}.$$

$$E_x = \frac{\mu_4}{s^4} - 3 = -1,2.$$

განვიხილოთ თანაბრად განაწილებული X შემთხვევითი სიდიდის $[a, b]$ ინტერვალის რაიმე $[\alpha, \beta]$ შუალედში მოხვედრის ალბათობა

$$P(\alpha < X < \beta) = \int_{\alpha}^{\beta} \frac{dx}{b-a} = \frac{x}{b-a} \Big|_{\alpha}^{\beta} = \frac{\beta - \alpha}{b-a}.$$

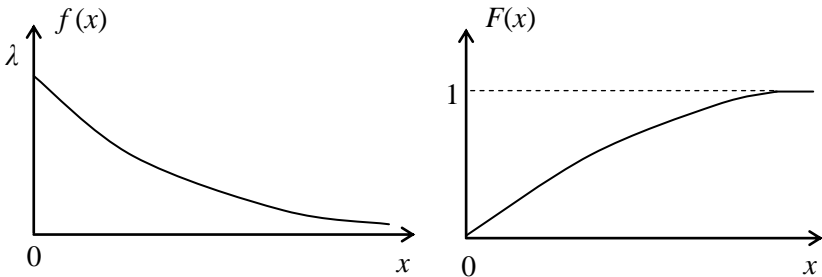
3.4. მაჩვენებლიანი განაწილება

უნყვეტი შემთხვევითი X სიდიდის მაჩვენებლიანი (ექსპონენციალური) განაწილება ეწოდება ისეთ განაწილებას, რომლის სიმკვრივის ფუნქციას აქვს შემდეგი სახე:

$$f(x) = \begin{cases} 0 & , \text{ როცა } x < 0 \\ \lambda e^{-\lambda x} & , \text{ როცა } x \geq 0 \end{cases}$$

სადაც, λ – მუდმივი დადებითი რიცხვია.

ჩვენ ვხედავთ, რომ მაჩვენებლიანი განაწილება განისაზღვრება მხოლოდ ერთი პარამეტრით λ , რაც მიგვანიშნებს მის დადებით მხარეზე. განაწილებისა და სიმკვრივის ფუნქციებს აქვთ შემდეგი სახე:



ვიპოვოთ განაწილების ფუნქცია:

$$F(x) = \int_{-\infty}^x f(x) dx = \int_{-\infty}^x \lambda e^{-\lambda x} dx = 1 - e^{-\lambda x}, \text{ როცა } x \geq 0.$$

განვსაზღვროთ მათემატიკური ლოდინი და დისპერსია:

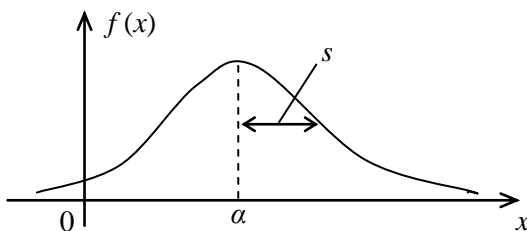
$$M(X) = \int_0^{\infty} xf(x)dx = \int_0^{\infty} x\lambda e^{-\lambda x} dx = \frac{1}{\lambda}.$$

$$D(X) = M(X^2) - [M(X)]^2 = \frac{2}{\lambda^2} - \frac{1}{\lambda^2} = \frac{1}{\lambda^2}; \quad s = \frac{1}{\lambda}.$$

ამრიგად, მაჩვენებლიანი განაწილებისათვის მათემატიკური ლოდინი და საშუალო კვადრატული გადახრა ერთმანეთის ტოლია. მაჩვენებლიანი განაწილება ხშირად გვხვდება საიმედოობის თეორიისა და მასობრივი მომსახურების ამოცანებში.

3.5. ნორმალური განაწილება

ნორმალური განაწილება ყველაზე უფრო გავრცელებული განაწილების კანონია. ნორმალური განაწილების ძირითადი თვისება იმაში მდგომარეობს, რომ სხვა განაწილებები მიისწრაფვიან ნორმალური განაწილებისაკენ. ნორმალური განაწილების კანონს ექვემდებარებიან მხოლოდ უწყვეტი შემთხვევითი სიდიდეები. ამიტომ ნორმალური განაწილება, რომელიც აღინიშნება $N(a,s)$ სიმბოლოთი, შეიძლება მოცემული იყოს განაწილების სიმკვრივის ან განაწილების ფუნქციის სახით.



$$f(x) = \frac{1}{s\sqrt{2\pi}} e^{-\frac{(x-a)^2}{2s^2}} \quad -\infty \leq x \leq \infty$$

$$F(x) = \frac{1}{s\sqrt{2\pi}} \int_{-\infty}^x e^{-\frac{(x-a)^2}{2s^2}} dx$$

ნორმალური განაწილების სიმკვრივის გრაფიკს ეწოდება ნორმალური მრუდი. მას აქვს ზარისებული ფორმა და სიმეტრიულია $x = a$ წერტილზე გამავალი წრფისა, ხოლო აბსცისათა ღერძებს ასიმპტოტურად უახლოვდება, როცა $x \rightarrow \infty$.

როგორც ნორმალური განაწილების სიმკვრივის ფორმულიდან ჩანს, ნორმალური განაწილება განისაზღვრება a და s პარამეტრებით. განვსაზღვროთ ეს პარამეტრები. ვიპოვოთ ნორმალურად განაწილებული შემთხვევითი სიდიდის მათემატიკური ლოდინი და დისპერსია

$$M(X) = \int_{-\infty}^{\infty} xf(x)dx = \frac{1}{s\sqrt{2\pi}} \int_{-\infty}^{\infty} xe^{-\frac{(x-a)^2}{2s^2}} dx.$$

შემოვიტანოთ ახალი ცვლადი $\frac{x-a}{s\sqrt{2}} = t$ (3.1), მაშინ

$$x - a = ts\sqrt{2}; \quad x = a + s\sqrt{2}t; \quad dx = s\sqrt{2}dt$$

$$M(X) = \frac{s\sqrt{2}}{s\sqrt{2}\sqrt{\pi}} \int_{-\infty}^{\infty} (a + s\sqrt{2}t)e^{-t^2} dt = \frac{a}{\sqrt{\pi}} \int_{-\infty}^{\infty} e^{-t^2} dt + \frac{s\sqrt{2}}{\sqrt{\pi}} \int_{-\infty}^{\infty} te^{-t^2} dt$$

მეორე ინტეგრალი ნულის ტოლია, როგორც ინტეგრალი კენტი ფუნქციიდან სიმეტრიულ ზღვრებში. პირველი ინტეგრალი წარმოადგენს პუასონის ცნობილ ინტეგრალს

$$\int_{-\infty}^{\infty} e^{-t^2} dt = \sqrt{\pi},$$

ამიტომ,

$$M(X) = \frac{a}{\sqrt{\pi}} \sqrt{\pi} = a.$$

ე.ი. a წარმოადგენს მათემატიკურ ლოდინს. ვიპოვოთ დისპერსია

$$D(X) = \int_{-\infty}^{\infty} (x-a)^2 f(x)dx = \frac{1}{\sigma\sqrt{2\pi}} \int_{-\infty}^{\infty} (x-a)^2 \exp\left\{-\frac{(x-a)^2}{2\sigma^2}\right\} dx.$$

აქაც, თუ გამოვიყენებთ ახალ ცვლადს, (3.1) ფორმულიდან მივიღებთ

$$(x-a) = s\sqrt{2}t, \quad D(X) = \frac{2s^2}{\sqrt{\pi}} \int_{-\infty}^{\infty} t^2 e^{-t^2} dt.$$

თუ გამოვიყენებთ ნაწილობითი ინტეგრირების ხერხს, საბოლოოდ მივიღებთ $D(X) = s^2$, ე.ი. s^2 პარამეტრი წარმოადგენს დისპერსიას. ნორმალური განაწილების სტანდარტული გადახრა $s = \sqrt{D(X)}$ წარმოადგენს მანძილს საშუალოსა და განაწილების მრუდის გადაღუნვის წერტილს შორის, ანუ ისეთ წერტილს შორის, სადაც მრუდი ამოზნექილობას ჩაზნექილობით ცვლის.

ამრიგად, ნათელი ხდება ნორმალური განაწილების მრუდის a და s პარამეტრების სტატისტიკური აზრი და ისინი სრულად განსაზღვრავენ მრუდის მდებარეობას რიცხვით ღერძზე და მის ფორმას.

განვიხილოთ ნორმალური განაწილების ზოგიერთი თვისებები:

1. სიმკვრივის ფუნქცია განსაზღვრულია მთელ ax ღერძზე, ე.ი. x -ის ყოველ მნიშვნელობას შეესაბამება ფუნქციის გარკვეული მნიშვნელობა.

2. x -ის ყველა მნიშვნელობისათვის (როგორც დადებითი, ასევე უარყოფითისთვის) სიმკვრივის ფუნქცია დადებითია, ე.ი. ნორმალური მრუდი მოთავსებულია ax ღერძის ზედა ნახევარსიბრტყეში.

3. სიმკვრივის ფუნქციის ზღვარი x -ის ზრდისას უდრის ნულს

$$\lim_{x \rightarrow \infty} f(x) = 0.$$

4. სიმკვრივის ფუნქციას $x = a$ წერტილში გააჩნია მაქსიმუმი, რომელიც ტოლია:

$$f(a) = \frac{1}{s\sqrt{2\pi}}.$$

5. სიმკვრივის ფუნქციის გრაფიკი სიმეტრიულია $x = a$ წერტილზე გამავალი წრფის მიმართ.

6. კენტი ცენტრალური მომენტები ნულის ტოლია.

7. ასიმეტრიისა და ექსცესის კოეფიციენტები ნულის ტოლია.

8. ნორმალური მრუდის ფორმა a (მათემატიკური ლოდინის) პარამეტრის სიდიდის შეცვლისას არ იცვლება. მათემატიკური ლოდინის გაზრდისას ან შემცირებისას, მრუდის გრაფიკი შესაბამისად გადაინაცვლებს მარცხნივ ან მარჯვნივ.

9. s პარამეტრის შეცვლისას იცვლება მრუდის ფორმა. s -ის ზრდისას განაწილების მრუდის მაქსიმალური ორდინატა მცირდება და პირიქით, s -ის შემცირებისას – იზრდება.

რადგან ნორმალური განაწილება დამოკიდებულია ორ პარამეტრზე (a, s), ამიტომ მისი ცხრილის სახით წარმოდგენა საკმაოდ რთულია. აქედან გამომდინარე, პრაქტიკაში ფართოდ გამოიყენება **სტანდარტიზირებული (სტანდარტული) ნორმალური განაწილება**, რომელიც მიიღება X შემთხვევითი სიდიდის Z – გარდაქმნით:

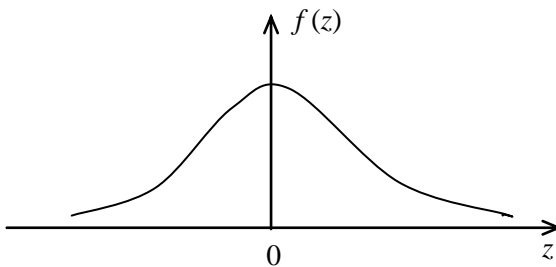
$$Z = \frac{X - a}{s}.$$

მიღებული Z სიდიდის მათემატიკური ლოდინი და დისპერსია ტოლია:

$$M[Z] = M\left[\frac{X - a}{s}\right] = \frac{1}{s}[M(X) - M(X)] = 0,$$

$$D(Z) = D\left[\frac{X - a}{s}\right] = \frac{1}{s^2}D(X - a) = 1.$$

თუ Z მნიშვნელობას ჩავსვამთ ნორმალური განაწილების სიმკვრივის ფუნქციის გამოსახულებაში, მაშინ მივიღებთ სტანდარტიზირებულ ნორმალურ განაწილებას $N(0,1)$ არამეტრებით, რომლის განაწილების სიმკვრივის ფუნქციას აქვს შემდეგი სახე:



$$f(z) = \frac{1}{\sqrt{2\pi}} e^{-\frac{z^2}{2}}, \quad F(z) = \frac{1}{\sqrt{2\pi}} \int_{-\infty}^z e^{-\frac{z^2}{2}} dz.$$

განაწილების $F(z)$ ფუნქცია დაკავშირებულია ლაპლასის $\Phi(z)$ ფუნქციასთან შემდეგი ტოლობით:

$$F(z) = \frac{1}{2} [1 + \Phi(z)], \quad (3.2)$$

სადაც

$$\Phi(z) = \frac{2}{\sqrt{2\pi}} \int_0^z e^{-\frac{t^2}{2}} dt.$$

3.6. ნორმალური განაწილება სიბრტყეზე

რადგან ორი შემთხვევითი სიდიდისგან შემდგარი სისტემა წარმოდგენილია შემთხვევითი წერტილით სიბრტყეზე, ამიტომ ორგანზომილებიანი სისტემის ნორმალური განაწილების კანონს ხშირად უწოდებენ ნორმალურ განაწილებას სიბრტყეზე.

დავუშვათ, X და Y დამოუკიდებელი შემთხვევითი სიდიდეებია, რომელთა ნორმალური განაწილების სიმკვრივის ფუნქციებია:

$$f_x(x) = \frac{1}{s_x \sqrt{2\pi}} \exp \left\{ -\frac{(x - m_x)^2}{2s_x^2} \right\},$$

$$f_y(y) = \frac{1}{s_y \sqrt{2\pi}} \exp \left\{ -\frac{(y - m_y)^2}{2s_y^2} \right\}.$$

მაშინ მათი ერთობლივი განაწილების სიმკვრივის ფუნქცია განისაზღვრება შემდეგი გამოსახულებით:

$$f(x, y) = f_x(x)f_y(y) = \frac{1}{2\pi s_x s_y} \exp \left\{ -\frac{1}{2} \left[\frac{(x - m_x)^2}{s_x^2} + \frac{(y - m_y)^2}{s_y^2} \right] \right\}. \quad (3.3)$$

თუ ორგანზომილებიანი სისტემის გაფანტვის ცენტრი ემთხვევა კოორდინანტთა სათავეს, მაშინ $m_x = m_y = 0$. შესაბამისად გვექნება:

$$f(x, y) = \frac{1}{2\pi s_x s_y} \exp\left\{-\frac{1}{2}\left(\frac{x^2}{s_x^2} + \frac{y^2}{s_y^2}\right)\right\}$$

და მას ერთობლივი ნორმალური განაწილების კანონიკური ფორმა ეწოდება.

თუ X და Y შემთხვევითი სიდიდეები დამოკიდებულნი არიან, მაშინ ერთობლივი ნორმალური განაწილების სიმკვრივის ფუნქციას აქვს შემდეგი სახე:

$$f(x, y) = \frac{1}{2\pi s_x s_y \sqrt{1-r_{xy}^2}} \exp\left\{-\frac{1}{2(1-r_{xy}^2)}\left[\frac{(x-m_x)^2}{s_x^2} - 2r_{xy} \frac{(x-m_x)(y-m_y)}{s_x s_y} + \frac{(y-m_y)^2}{s_y^2}\right]\right\} \quad (3.4)$$

და იგი დამოკიდებულია ხუთ პარამეტრზე: $m_x, m_y, \sigma_x, \sigma_y$ და r_{xy} .

პირობითი ნორმალური განაწილების სიმკვრივის ფუნქცია განისაზღვრება შემდეგი ფორმულით:

$$f(y|x) = \frac{1}{\sqrt{2\pi} s_y \sqrt{1-r_{xy}^2}} \exp\left\{-\frac{1}{2s_y^2(1-r_{xy}^2)}\left[y - m_y - r_{xy} \frac{s_y}{s_x} (x - m_x)\right]^2\right\},$$

$$f(x|y) = \frac{1}{\sqrt{2\pi} s_x \sqrt{1-r_{xy}^2}} \exp\left\{-\frac{1}{2s_x^2(1-r_{xy}^2)}\left[x - m_x - r_{xy} \frac{s_x}{s_y} (y - m_y)\right]^2\right\}.$$

თუ X და Y სიდიდეები არაკორელირებულია ($r_{xy} = 0$), მაშინ (3.4) გამოსახულებიდან მივიღებთ (3.3) ფორმულას, რომელიც წარმოადგენს დამოუკიდებელი შემთხვევითი სიდიდეების ერთობლივი

განაწილების სიმკვრივის ფუნქციას. აქედან გამომდინარე, თუ ნორმალურად განაწილებული შემთხვევითი სიდიდეები არაკორელირებულნი არიან, მაშინ ისინი დამოუკიდებელნი არიან.

3.7. მრავალგანზომილებიანი სისტემის ნორმალური განაწილების კანონი

ბიოსამედიცინო კვლევებში, როგორც წესი, საქმე გვაქვს არა ერთ ან ორ, არამედ საკმაოდ ბევრ პარამეტრთან. ვთქვათ მოცემულია მრავალგანზომილებიანი სისტემა, რომელიც აღიწერება X_1, X_2, \dots, X_n ნორმალურად განაწილებული შემთხვევითი სიდიდეებით. მაშინ განაწილების სიმკვრივის ფუნქციას აქვს შემდეგი სახე:

$$f(X_1, X_2, \dots, X_n) = \frac{1}{(2\pi)^{\frac{n}{2}} |S|^{\frac{1}{2}}} \exp\left\{-\frac{1}{2}(X - \bar{X})' S^{-1} (X - \bar{X})\right\}, \quad (3.5)$$

სადაც, X სანყისი მონაცემების მატრიცაა

$$X = (X_1, X_2, \dots, X_n)' = \begin{bmatrix} x_{11} & x_{12} & \dots & x_{1n} \\ x_{21} & x_{22} & \dots & x_{2n} \\ \dots & \dots & \dots & \dots \\ x_{m1} & x_{m2} & \dots & x_{mn} \end{bmatrix},$$

\bar{X} – საშუალო სიდიდეთა ვექტორია $\bar{X} = (\bar{x}_1, \bar{x}_2, \dots, \bar{x}_n)'$,

S – კოვარიაციული მატრიცაა, რომელიც განისაზღვრება შემდეგნაირად:

$$S = \frac{1}{m-1} \sum_{i=1}^m (X_i - \bar{X})(X_i - \bar{X})' = \begin{bmatrix} s_{11} & s_{12} & \dots & s_{1n} \\ s_{21} & s_{22} & \dots & s_{2n} \\ \dots & \dots & \dots & \dots \\ s_{n1} & s_{n2} & \dots & s_{nn} \end{bmatrix}$$

და რომლის მთავარ დიაგონალზე ($s_{11}, s_{22}, \dots, s_{nn}$) იმყოფება დისპერსიების მნიშვნელობები. მიღებული მატრიცა სიმეტრიულია, ე.ი. $s_{ij} = s_{ji}$.

თუ მოვახდენთ მოცემული სისტემის X_1, X_2, \dots, X_n შემთხვევითი სიდიდეების სტანდარტიზირებას (ნორმირებას)

$$Z_i = \frac{X_i - \bar{X}_i}{\sqrt{s_{ii}}},$$

მაშინ მივიღებთ კორელაციურ მატრიცას

$$R = \begin{bmatrix} 1 & r_{12} & \dots & r_{1n} \\ r_{21} & 1 & \dots & r_{2n} \\ \dots & \dots & \dots & \dots \\ r_{n1} & r_{n2} & \dots & 1 \end{bmatrix},$$

ხოლო განაწილების სიმკვრივის ფუნქციას ექნება შემდეგი სახე:

$$f(Z_1, Z_2, \dots, Z_n) = \frac{1}{(2\pi)^{\frac{n}{2}} |R|^{\frac{1}{2}}} \exp\left\{-\frac{1}{2} Z' R^{-1} Z\right\}.$$

ზოგადად, მრავალგანზომილებიანი სისტემის ნორმალურ განაწილებას $N_n(\bar{X}, S)$ სიმბოლოთი აღნიშნავენ. უნდა შევნიშვნოთ, რომ თუ $n=1$, მაშინ (3.5) გამოსახულება გარდაიქმნება ერთგანზომილებიან ნორმალურად განაწილებულ შემთხვევითი სიდიდის განაწილების სიმკვრივის გამოსახულებად. ამრიგად, (3.5) არის ნორმალური განაწილების განზოგადებული ფორმულა.

4. ალბათობის თეორიის ზღვარითი თეორემა

ალბათობის თეორია შეისწავლის კანონზომიერებებს, რომლითაც ხასიათდება მასიური შემთხვევითი მოვლენები. თუ ცდათა რაოდენობა საკმაოდ დიდია, მაშინ შემთხვევითი მოვლენებისა და შემთხვევითი სიდიდეების მახასიათებლები კარგავენ შემთხვევით

ხასიათს და თითქმის არაშემთხვევითები ხდებიან. მაგ. ხდომილობის სიხშირე ცდათა დიდი რაოდენობისას ხდება სტაბილური, ასევე ითქმის შემთხვევითი სიდიდის საშუალო მნიშვნელობების მიმართ. ეს გარემოება საშუალებას იძლევა, შემთხვევითი სიდიდეების დაკვირვებათა შედეგები გამოვიყენოთ მომავალში ცდათა შედეგების პროგნოზირებისთვის.

თეორემათა ჯგუფი, რომლებიც აკავშირებენ შემთხვევითი სიდიდეებისა და შემთხვევითი მოვლენების თეორიულ და ექსპერიმენტალურ მახასიათებლებს, როცა ცდათა რაოდენობა დიდია, გაერთიანებული არიან ალბათობის თეორიის ზღვარითი თეორემების დასახელებით. აქვე განვიხილოთ დიდ რიცხვთა კანონი და ცენტრალური ზღვარითი თეორემა.

4.1. დიდ რიცხვთა კანონი

როგორც ვიცით, შემთხვევითი ხდომილობის ფარდობითი სიხშირე, ცდათა მრავალგზის გამეორებისას, იჩენს სტატისტიკურ მდგრადობას, ანუ იგი მცირედ განსხვავდება რაიმე მუდმივი სიდიდისაგან. დიდ რიცხვთა კანონი ამყარებს კავშირს საშუალო სიდიდესა და რაიმე მუდმივ სიდიდეს შორის. დიდ რიცხვთა კანონი პირველად ჩამოაყალიბა ი.ბერნულიმ. შემდგომში ეს კანონი განავითარეს ჩებიშევიმა, მარკოვმა და სხვებმა. ჩვენ განვიხილავთ ჩებიშევისა და ბერნულის თეორემებს. ამისათვის ჯერ განვიხილოთ ჩებიშევის უტოლობა.

ჩებიშევის უტოლობა. განვიხილოთ X შემთხვევითი სიდიდე, რომლის მათემატიკური ლოდინია m_x და დისპერსია D_x . ჩამოვყალიბოთ ჩებიშევის უტოლობა. თუ შემთხვევითი სიდიდის თავის მათემატიკურ ლოდინთან გადახრის ალბათობა აბსოლუტური სიდიდით ნებისმიერ დადებით ε რიცხვზე მეტია, მაშინ იგი ზემოდან შემოსაზღვრულია $\frac{D_x}{\varepsilon^2}$ სიდიდით, ე.ი.

$$P(|X - m_x| \geq \varepsilon) \leq \frac{D_x}{\varepsilon^2}.$$

ჩებიშევის უტოლობა შეიძლება მეორენაირად ჩაიწეროს. თუ გამოვიყენებთ მოპირდაპირე ხდომილობის ცნებას, მაშინ გვექნება:

$$P(|X - m_x| < \varepsilon) \geq 1 - \frac{D_x}{\varepsilon^2}.$$

ჩებიშევის თეორემა. ჩებიშევის თეორემა წარმოადგენს დიდი რიცხვთა კანონის ერთ-ერთ მნიშვნელოვან ფორმას. იგი ამყარებს კავშირს შემთხვევითი სიდიდის საშუალო არითმეტიკულსა და მათემატიკურ ლოდინს შორის. ჩამოვყალიბოთ ჩებიშევის თეორემა: დამოუკიდებელ ცდათა უსასრულო გაზრდით შემთხვევითი სიდიდის საშუალო არითმეტიკული, რომელსაც გააჩნია სასრული დისპერსია, ალბათურად კრებადია შესაბამისი მათემატიკური ლოდინისაკენ.

განვიხილოთ ტერმინი „ალბათურად კრებადი“. X_1, X_2, \dots, X_n შემთხვევითი სიდიდეების მიმდევრობას ეწოდება ალბათურად კრებადი რაიმე a სიდიდისაკენ, თუ ნებისმიერი $\varepsilon > 0$ რიცხვისათვის სრულდება შემდეგი ტოლობა

$$\lim_{n \rightarrow \infty} P(|X_n - a|) = 1.$$

ამრიგად, ჩებიშევის თეორემა აღნიშნავს, რომ თუ X_1, X_2, \dots, X_n ერთნაირად განაწილებული დამოუკიდებელი შემთხვევითი სიდიდეებია m_x მათემატიკური ლოდინითა და D_x დისპერსიით, მაშინ ნებისმიერი $\varepsilon > 0$ სიდიდისათვის გვექნება:

$$\lim_{n \rightarrow \infty} P\left(\left|\frac{1}{n} \sum_{i=1}^n x_i - m_x\right| < \varepsilon\right) = 1.$$

ჩებიშევის თეორემის არსი მდგომარეობს შემდეგში: მიუხედავად იმისა, რომ ცალკეულ დამოუკიდებელ შემთხვევით სიდიდეებს შეუძლიათ მიიღონ თავისი მათემატიკური ლოდინისაგან საგრძნობლად განსხვავებული მნიშვნელობები, დიდი რაოდენობის შემთხვევით სიდიდეთა საშუალო არითმეტიკული, გარკვეული აზრით, კარგავს შემთხვევითი სიდიდის ხასიათს და მათი საშუალო არითმეტიკულის გადახრა საკმაოდ მცირეა, როცა $n \rightarrow \infty$.

ჩებიშევის თეორემის კერძო შემთხვევას წარმოადგენს ი.ბერნულის თეორემა, რომელიც შემდეგში მდგომარეობს: ვთქვათ, m არის A ხდომილობის მოხდენის რიცხვი n დამოუკიდებელ ცდაში. ვიგულისხმობთ, რომ თითოეულ ცდაში ხდომილობის მოხდენის ალბათობა არის P , მაშინ ნებისმიერი $\varepsilon > 0$ რიცხვისათვის ადგილი აქვს ტოლობას:

$$\lim_{n \rightarrow \infty} P \left[\left| \frac{m}{n} - p \right| < \varepsilon \right] = 1,$$

სადაც, $\frac{m}{n}$ სიდიდე A ხდომილობის ფარდობითი სიხშირეა, ხოლო $P - A$ ხდომილობის მოხდენის ალბათობაა. ამ თეორემების თანახმად, თუ ცდათა რიცხვი დიდია, მაშინ ფარდობითი სიხშირე თეორიული ალბათობის შეფასებად შეიძლება იქნას მიღებული. ეს შედეგი მეტად მნიშვნელოვანია ბევრი პრაქტიკული ამოცანების გადასაწყვეტად.

თუ ცდები ტარდება არაერთგვაროვან პირობებში, მაშინ პუასონის თეორემის თანახმად, რომელიც წარმოადგენს ბერნულის თეორემის განზოგადებულ თეორემას, დამოუკიდებელ ცდათა რაოდენობის ზრდისას, A ხდომილობის სიხშირე ალბათურად კრებადია შესაბამისი ალბათობების საშუალო არითმეტიკულისკენ.

ამრიგად, ჩებიშევის თეორემის თანახმად, შემთხვევით სიდიდეთა საშუალო არითმეტიკული ალბათურად კრებადია შესაბამის მათემატიკურ ლოდინთა საშუალო არითმეტიკულისკენ, ხოლო ბერნულის თეორემის თანახმად, ხდომილობის ფარდობითი სახშირე ალბათურად კრებადია ამ ხდომილობის ალბათობისკენ.

4.2. ცენტრალური ზღვართი თეორემა

წინა პარაგრაფში განხილული თეორემები წარმოადგენენ დიდ რიცხვთა კანონის სხვადასხვა ფორმით. როგორც აღვნიშნეთ, დიდ რიცხვთა კანონი გამოხატავს შემთხვევითი სიდიდის მისწრაფებას რაიმე მუდმივ მახასიათებლისაკენ, თანაც ჩვენ საქმე არა გვაქვს შემთხვევითი სიდიდის განაწილების კანონთან.

განვიხილოთ შემთხვევით სიდიდეთა მიმდევრობის X_1, X_2, \dots, X_n ჯამის ზღვრული განაწილების კანონი, როდესაც შესაკრებთა რაოდენობა უსასრულოდ იზრდება. აღმოჩნდა, რომ თუ შესაკრებ შემთხვევათა სიდიდეები გარკვეულ პირობებს აკმაყოფილებენ და მათი რიცხვი საკმარისად დიდია, მაშინ ჯამის განაწილების კანონი უახლოვდება ნორმალურს. თეორემებს, რომლებიც ამ ფაქტს ასახავენ, ეწოდებათ ზღვართი თეორემები.

ჩამოვყალიბოთ (დაუმტკიცებლად) უმარტივესი ზღვართი თეორემა, რომელიც ლინდბერგლევის თეორემით არის ცნობილი. თუ X_1, X_2, \dots, X_n ერთნაირად განაწილებული დამოუკიდებელი შემთხვევითი სიდიდეებია m მათემატიკური ლოდინითა და σ^2

დისპერსიით, მაშინ $Y_n = \sum_{i=1}^n X_i$ გამოსახულების განაწილების კანონი

მიისწრაფვის ნორმალურისაკენ, როცა $n \rightarrow \infty$.

ა. ლიპუნოვმა ეს თეორემა განაზოგადა ნებისმიერად განაწილებული შემთხვევითი სიდიდის მიმართ. მოვიყვანოთ ეს თეორემა დაუმტკიცებლად. თუ X_1, X_2, \dots, X_n ნებისმიერად განაწილებული დამოუკიდებელი შემთხვევითი სიდიდეებია, არსებობს მესამე რიგის აბსოლუტური მომენტი

$$M\left(|X_k|^3\right), k=1,2,\dots \text{ და სრულდება პირობა: } \lim_{n \rightarrow \infty} \frac{\sum_{k=1}^n M\left[|X_k|^3\right]}{\left(\sum_{k=1}^n D(X_k)\right)^{3/2}} = 0,$$

მაშინ $Y_n = \sum_{i=1}^n X_i$ შემთხვევითი სიდიდეთა განაწილების კანონი

მიისწრაფვის ნორმალურისაკენ.

თუ თანაბრად განაწილებული X_1, X_2, \dots, X_n დისკრეტული შემთხვევითი სიდიდეები იღებენ 0 ან 1 მნიშვნელობებს, მაშინ მუავრი-ლაპლასის თეორემა, რომელიც წარმოადგენს ცენტრალური ზღვართი თეორემის უმარტივეს სახეს, შეიძლება ასე ჩამოვყალიბოთ: თუ ჩატარდა n დამოუკიდებელი ცდა, სადაც A ხდომილობა მოხდა p ალბათობით, მაშინ ნებისმიერი $[\alpha, \beta]$ ინტერვალისთვის სამართლიანია შემდეგი თანაფარდობა:

$$P\left(\alpha < \frac{m - np}{\sqrt{npq}} < \beta\right) = \frac{1}{2} \left[\Phi\left(\frac{\beta}{\sqrt{2}}\right) - \Phi\left(\frac{\alpha}{\sqrt{2}}\right) \right],$$

სადაც, $m - A$ ხდომილობის მოხდენის რაოდენობაა, $q = 1 - p$, $\Phi(\cdot)$ – ლაპლასის ფუნქციაა. მუავრი-ლაპლასის თეორემა აღწერს ბინომური განაწილების მოქმედებას, როცა n დიდია.

ბიოსტატისტიკის მეთოდები

5. ბიოსტატისტიკის არსი

ბიოსტატისტიკა, ისევე როგორც მათემატიკური სტატისტიკა, არის მათემატიკის ნაწილი, რომელიც სწავლობს შემთხვევით სიდიდეზე წარმოებული დაკვირვებების შედეგების შეკრების, სისტემატიზაციისა და დამუშავების მეთოდებს არსებული კანონზომიერების გამოვლენის მიზნით. სტატისტიკურ მეთოდებს საფუძვლად უდევს ექსპერიმენტალური მონაცემები, რომლებთაც სტატისტიკურ მონაცემებს უწოდებენ.

როგორც ცნობილია, შემთხვევით სიდიდეზე ყველაზე ზუსტი ცნობები შეიძლება მივიღოთ, თუ ჩავატარებთ ამ შემთხვევითი სიდიდის მაქსიმალურად ბევრ გაზომვას. შემოვიტანოთ გენერალური ერთობლიობის ცნება.

გენერალური ერთობლიობა ანუ სტატისტიკური პოპულაცია ეწოდება ყველა შესაძლო დაკვირვებათა ერთობლიობას, რომლებიც შეიძლება მივიღოთ შემთხვევითი სიდიდის გაზომვით. გენერალური ერთობლიობის შემადგენელ წევრთა რაოდენობას ეწოდება ამ გენერალური ერთობლიობის მოცულობა ანუ განზომილება.

პრაქტიკულად, გენერალური ერთობლიობის მიღება თითქმის შეუძლებელია სხვადასხვა მიზეზის გამო, გარდა გამონაკლის შემთხვევაში, როცა გენერალური ერთობლიობა შედგება წევრთა სასრული რაოდენობისგან. აქედან გამომდინარე, როგორც წესი, ჩვენ საქმე გვაქვს სასრული რაოდენობის დაკვირვებებთან, რომლებსაც ამონარჩევი ეწოდებათ.

ამონარჩევი ერთობლიობა ან უბრალოდ ამონარჩევი (შერჩევა, ამოკრეფა) ეწოდება იმ დაკვირვებების (ობიექტების) ერთობლიობას, რომლებიც ამოღებულია გენერალური ერთობლიობიდან. დაკვირვებათა მწკრივი x_1, x_2, \dots, x_n განიხილება, როგორც n მოცულობის ამონარჩევი სასრული ან უსასრულო გენერალური ერთობლიობიდან. ამონარჩევი ყოველთვის სასრული რაოდენობისაა.

მათემატიკური სტატისტიკის კვლევის ერთ-ერთ ძირითად მეთოდს წარმოადგენს ე.წ. **შერჩევითი მეთოდი**. მეთოდს, რომელიც აკეთებს დასკვნას ამონარჩევის მახასიათებლებისა და თვისებების საფუძველზე, გენერალური ერთობლიობის რიცხვით

მახასიათებლებზე და განაწილების კანონის შესახებ, ეწოდება შერჩევითი მეთოდი.

იმისათვის, რომ X შემთხვევითი სიდიდის განაწილების კანონი ან მისი რიცხვითი მახასიათებლების შესახებ მსჯელობა იყოს ეფექტური, აუცილებელია, რომ ამონარჩევი იყოს წარმომადგენლობითი (რეპრეზენტატიული), ე.ი. საკმაოდ კარგად წარმომადგენდეს შესასწავლ შემთხვევით სიდიდეს.

არსებობენ სპეციალური მეთოდები წარმომადგენლობითი ამონარჩევის მისაღებად, რომელთა არსი დაიყვანება იმაზე, რომ გენერალური ერთობლიობის ყოველ ელემენტს ჰქონდეს თანაბარი ალბათობა იმისა, რომ მოხვდეს ამონარჩევიში. ამონარჩევის რეპრეზენტატიულობა მიიღწევა რანდომიზაციით (ინგლისური სიტყვიდან *random*-შემთხვევითი) ანუ, სხვანაირად რომ ვთქვათ, გენერალური ერთობლიობიდან ელემენტების შემთხვევითი ამოკრევით. ამონარჩევის შემთხვევითობა მიიღება ან მექანიკური მეთოდებით, რომლებიც ეფუძნება გენერალური ერთობლიობის ნაწილ-ნაწილ დაყოფას, ანდა მათემატიკური მეთოდებით, მაგალითად, მონტე-კარლოს მეთოდით.

იყენებს რა ალბათობის თეორიის მეთოდებს, მათემატიკური სტატისტიკის მიზანია, ამონარჩევის საშუალებით შეაფასოს გენერალური ერთობლიობის მახასიათებლები. ალბათობის თეორიიდან ვიცით, რომ შემთხვევით სიდიდეს გააჩნია გარკვეული სახის განაწილების ფუნქცია მისი შესაბამისი რიცხვითი მახასიათებლებით (მაგ. მათემატიკური ლოდინი ან სხვა სანყისი და ცენტრალური მომენტები), რომლებსაც შემდგომში თეორიულს ვუწოდებთ, ხოლო ამონარჩევის საშუალებით მიღებულ მნიშვნელობებს – შერჩევითს ანუ ემპირიულს.

მათემატიკურ სტატისტიკაში ძალიან ხშირად გამოიყენება „თავისუფლების ხარისხის“ ცნება. **თავისუფლების ხარისხი** არის გამონათქვამი, რომელიც ნასესხებია ფიზიკიდან, სადაც იგი ახასიათებს ობიექტის მოძრაობას. თუ ობიექტს აქვს შესაძლებლობა იმოძრაოს მხოლოდ სწორხაზოვნად, მაშინ მას აქვს ერთი თავისუფლების ხარისხი. ობიექტს, რომელსაც შეუძლია სიბრტყეზე მოძრაობა, გააჩნია ორი თავისუფლების ხარისხი და ა.შ. თუ გამოვალთ თავისუფლების ხარისხის ასეთი გეომეტრიული ინტერპრეტაციიდან, მაშინ იგი შეიძლება ასე განისაზღვროს:

$$v = n - k,$$

სადაც, v – თავისუფლების ხარისხია;

n – ამონარჩევის განზომილება, რომლითაც განისაზღვრება შესაფასებელი პარამეტრი;

k – დამატებითი პარამეტრების რაოდენობა, რომლებიც განისაზღვრებიან იგივე ამონარჩევით და შედიან შესაფასებელი პარამეტრის გამოსახულებაში.

მაგალითად, დისპერსიისთვის $\sigma_x^2 = \frac{1}{n-1} \sum_{i=1}^n (x_i - \bar{x})^2$, $k = 1$,

რადგან დისპერსიის ფორმულაში არის ერთი დამატებითი ცვლადი – საშუალო არითმეტიკული \bar{x} , რომელიც განისაზღვრება იგივე ამონარჩევით. აქედან გამომდინარე, $v = n - 1$.

6. მონაცემების კლასიფიკაცია და მათი წარმოდგენის მეთოდები. სიხშირული ანალიზი

ობიექტის შესწავლისას საჭიროა მისი მთელი რიგი მაჩვენებლების გაზომვა და დაფიქსირება. დაკვირვებისა თუ ფიზიკური გაზომვების შედეგად მიღებულ ინფორმაციას უწოდებენ შესასწავლი ობიექტის პირველად (ნედლ, დაუმუშავებელ) მონაცემებს, ხოლო ამ მონაცემთა ერთობლიობას – სტატისტიკურ ერთობლიობას. საზოგადოდ, მონაცემები არის ობიექტთა რაიმე სიმრავლის მახასიათებელთა დაკვირვებული მნიშვნელობების ერთობლიობა. ამ ერთობლიობის ყოველ წევრს დაკვირვება, მონაცემი ანუ მონაცემი წერტილი ეწოდება.

თავდაპირველი მონაცემების დალაგების, დაჯგუფებისა და ერთიანი მიმოხილვის მეთოდოლოგია შეადგენს **დესკრიფციულ** ანუ **აღწერით სტატისტიკას**. უნდა გვახსოვდეს, რომ მონაცემთა გაერთიანება შესაძლებელია მათი ერთგვაროვანი ნიშნების მიხედვით. ტერმინი „ნიშნის“ ქვეშ იგულისხმება თვისება, რითაც ერთი საგანი განსხვავდება მეორისგან. ნიშნებს გააჩნიათ ცვალებადობის თვისება და ამიტომ მათ **ვარიანტს** უწოდებენ. საზოგადოდ, ყველა მონაცემი (ნიშანი) ცვალებადია და ექვემდებარება უშუალოდ გაზომვას. ისინი პირობითად შეიძლება დავყოთ თვისებრივ და რაოდენობრივ მონაცემებად.

თვისებრივი მონაცემები ფიზიკურად არ იზომება, ისინი მხოლოდ აღირიცხება ამა თუ იმ ნიშნის ყოფნა-არყოფნის საშუა-

ლებით. მაგალითად, პაციენტში ამა თუ იმ სიმპტომის არსებობა ან არარსებობა, პაციენტთა კლასიფიკაცია ავადმყოფობის სიმძიმის მიხედვით და სხვ. თვისებრივი მონაცემის კონკრეტულ დონეს **ატრიბუტი** ეწოდება.

რაოდენობრივი მონაცემები ექვემდებარება უშუალოდ გაზომვებს. ისინი პირობითად შეიძლება დავყოთ ორ კლასად: დისკრეტულად და უწყვეტად. გავიხსენოთ, რომ უწყვეტია ისეთი სიდიდე, რომელსაც შეუძლია მიიღოს ყველა შესაძლო მნიშვნელობა რაღაც გარკვეული ინტერვალის ფარგლებში. მაგ. მაქსიმალური არტერიული წნევის მნიშვნელობა, პულსის სიხშირე, სისხლში ლეიკოციტების რაოდენობა და სხვ. თუ ხდება რაიმე მონაცემების თვლა, მაგ. ოპერაციების რიცხვი, მწყობრიდან გამოსული ხელსაწყოების რაოდენობა და სხვა, რომლებიც დისკრეტულად იცვლებიან, მაშინ მათ დისკრეტული მონაცემები ეწოდებათ.

ცვალებად თვისებრივ და რაოდენობრივ მონაცემებს მათემატიკურ სტატისტიკაში ეწოდება ცვლადი შემთხვევითი სიდიდეები, რომლებიც ლათინური ალფავიტის დიდი ასოებით აღინიშნება X, Y, Z, \dots , ხოლო მათი რიცხვითი მონაცემები შესაბამისად მცირე ლათინური ასოებით – x_1, x_2, \dots, x_n ; y_1, y_2, \dots, y_m და ა.შ.

მონაცემების ზემოთ მოყვანილი კლასიფიკაციის გათვალისწინებით არსებობს გაზომვის სამი სკალა: ნომინალური, რიგითი და რიცხვითი.

ნომინალური სკალა წარმოადგენს გაზომვის უმარტივეს დონეს და გამოიყენება თვისებრივი (კატეგორიული, ბინარული) მონაცემების გასაზომად. როდესაც კატეგორიებს შორის არსებობს გარკვეული მიმართება ან რიგი, მაშინ თვისებრივი მონაცემები იზომება **რიგის სკალის** გამოყენებით. დაკვირვებები აქაც კლასიფიცირებულია, მაგრამ მათ შორის არსებობს „მეტობის“ ან „ნაკლებობის“ მიმართება. მაგ. სახსრებით დაავადებული ავადმყოფები კლასიფიცირდებიან ავადმყოფობის სიმძიმის მიხედვით ოთხ კლასად, სადაც პირველი კლასი შეესაბამება ნორმალურ აქტივობას, ხოლო მეოთხე კლასი – ინვალიდის ეტლით მოძრაობას.

რაოდენობრივი მონაცემები იზომებიან **რიცხვით სკალაზე**, რომელიც იყოფა ინტერვალურ ანუ უწყვეტ სკალად და დისკრეტულ სკალად. უწყვეტ სკალაზე იზომება უწყვეტი შემთხვევითი სიდიდეები, ხოლო დისკრეტულ სკალაზე – დისკრეტული შემთხვევითი სიდიდეები.

თვისებრივი მაჩვენებლების წარმოდგენა. ვთქვათ, ამონარჩევის n ელემენტები ხასიათდება ერთი თვისებრივი მაჩვენებლით, რომლის მიმართ შესაძლებელია ორი მსჯელობა „ნიშანი არის“ (აღვნიშნოთ A -თი) და „ნიშანი არ არის“ (\bar{A}). მაშინ დაკვირვებები შეიძლება ასე წარმოვადგინოთ:

i	1	2	3	...	n
ნიშანი	A	\bar{A}	\bar{A}	...	A

უფრო უკეთესია, თუ მოვახდენთ მათ დაჯგუფებას. მაშინ გვექნება:

ვარიანტი	A	\bar{A}
სიხშირე	m_A	$m_{\bar{A}}$

ცხადია, რომ სიხშირეთა ჯამი $\sum_i m_i = n$. სიხშირე ესაა რიცხვი,

რომელიც გვიჩვენებს, თუ რამდენჯერ გვხვდება მოცემული ვარიანტი (მნიშვნელობა) ამონარჩევში. განვიხილოთ ერთობლიობა, რომლის ელემენტები ხასიათდება ორი ალტერნატიული ნიშნით A და B . მაშინ მონაცემების დაჯგუფების შედეგად ვიღებთ შემდეგ ცხრილს, რომელსაც ოთხუჯრედიანი ანუ (2×2) ტიპის ცხრილი ეწოდება.

$\begin{matrix} & B \\ A \end{matrix}$	B	\bar{B}	სულ
A	m_{AB}	$m_{A\bar{B}}$	m_A
\bar{A}	$m_{\bar{A}B}$	$m_{\bar{A}\bar{B}}$	$m_{\bar{A}}$
სულ	m_B	$m_{\bar{B}}$	n

რაოდენობრივი მაჩვენებლების წარმოდგენა. დავუშვათ, რომ უნდა შევისწავლოთ რაიმე დისკრეტული შემთხვევითი სიდიდე X , რომლის განაწილების კანონი უცნობია. განაწილების კანონის შეფასებისათვის ან სტატისტიკური მახასიათებლების გამოთვლისათვის ტარდება დამოუკიდებელი გაზომვების სერია x_1, x_2, \dots, x_n . გაზომვების შედეგად მიღებული მასალა წარმოვადგინოთ ცხრილის სახით:

i	1	2	...	n
x_i	x_1	x_2	...	x_n

ამ ცხრილს ეწოდება **სტატისტიკური მწკრივი**. იგი წარმოადგენს სტატისტიკური მასალის წარმოდგენის პირველად ფორმას.

დიდი რაოდენობით გაზომვის შემთხვევაში სტატისტიკური მასალის ცხრილის სახით წარმოდგენა მოუხერხებელია და მისი საშუალებით პრაქტიკულად შეუძლებელია გამოსაკვლევი X შემთხვევითი სიდიდის განაწილების კანონის დადგენა. ამიტომ, მიზანშეწონილია მონაცემების დაჯგუფება. კერძოდ, გაზომვის შედეგად მიღებული მნიშვნელობებით გამოვთვალოთ m_i სიხშირეები ან ფარდობითი სიხშირეები

$$p_i^* = \frac{m_i}{n}, \quad i = 1, 2, \dots, k.$$

ამის შემდეგ ვიღებთ **სიხშირულ ცხრილს**, რომელსაც აქვს შემდეგი სახე:

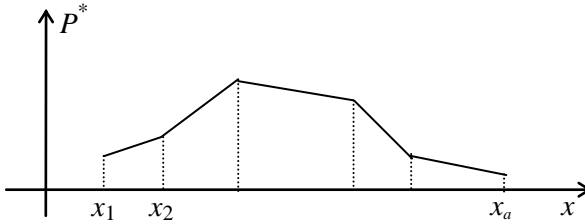
x_i	x_1	x_2	...	x_k
m_i	m_1	m_2	...	m_k
p_i^*	p_1^*	p_2^*	...	p_k^*

ცხრილი სწორადაა შედგენილი, თუ სრულდება შემდეგი პირობები:

$$\sum_{i=1}^k m_i = n \quad \text{და} \quad \sum_{i=1}^k p_i^* \approx 1.$$

ასეთი სიხშირული ცხრილით წარმოდგენილ სტატისტიკურ მწკრივს ეწოდება **ვარიაციული მწკრივი**.

თვალსაჩინოებისათვის სიხშირული ცხრილის მაგივრად იყენებენ მის გრაფიკულ გამოსახულებას, სადაც აბსცისათა ღერძზე გადაზომილია ვარიაციული მწკრივის მნიშვნელობები, ხოლო ორდინატთა ღერძზე – მათი შესაბამისი ფარდობითი სიხშირეები.



მიღებულ მრუდს უწოდებენ განაწილების მრავალკუთხედს ან განაწილების სიხშირის **პოლიგონს**.

თუ შეისწავლება უწყვეტი შემთხვევითი სიდიდე, მაშინ დაჯგუფება მდგომარეობს შემთხვევითი სიდიდის ცვალებადობის ინტერვალის თანაბარი სიგრძის k რაოდენობის კერძო ინტერვალებად დაყოფაში და ამ ინტერვალებისთვის m_i სიხშირეებისა და P_i^* ფარდობითი სიხშირეების გამოთვლაში.

ინტერვალების რაოდენობა აირჩევა ნებისმიერად, ჩვეულებრივ, 5-დან 15-მდე. ინტერვალის ოპტიმალური სიგრძის განსაზღვრისთვის იყენებენ სტერჯესის ფორმულას:

$$h = \frac{x_{\max} - x_{\min}}{1 + 3,32 \cdot \lg n},$$

სადაც, x_{\min} და x_{\max} ამონარჩევის მინიმალური და მაქსიმალური მნიშვნელობებია, n – ამონარჩევის განზომილება. მიღებული h სიდიდე, უმჯობესია, დავამრგვალოთ. მაგალითად, თუ $h < 1$, მაშინ დავამრგვალოთ მეთაქვამდე, სხვა შემთხვევაში – მთელამდე.

პირველი ინტერვალის საწყისად რეკომენდებულია მივიღოთ შემდეგი მნიშვნელობა:

$$x^{(0)} = x_{\min} - \frac{h}{2}.$$

მეორე ინტერვალის დასაწყისად, რომელიც ემთხვევა პირველი ინტერვალის ბოლოს – $x^{(1)} = x^{(0)} + h$ და ა.შ. მანამ, სანამ არ მივიღებთ ინტერვალს, რომელშიც მოხვდება x_{\max} მნიშვნელობა. ამ გზით მივიღებთ ინტერვალურ ვარიაციულ მწკრივს, რომელიც შეიძლება წარმოვადგინოთ შემდეგი სიხშირული ცხრილის სახით:

ინტერ- ვალეები	ინტერ- ვალეების საშ. მნი- შვნელობა x'_i	სიხ- ში- რეები m_i	ფარ- დობი- თი სიხშ. p_i^*	დაგროვილი სიხშირეები $\sum_i m_i$	დაგროვილი ფარდობითი სიხშირეები $\sum_i p_i^*$
$[x^{(0)}; x^{(1)}[$	x'_1	m_1	p_1^*	m_1	p_1^*
$[x^{(1)}; x^{(2)}[$	x'_2	m_2	p_2^*	$m_1 + m_2$	$p_1^* + p_2^*$
...
$[x^{(k-1)}; x^{(k)}]$	x'_k	m_k	p_k^*	$m_1 + \dots + m_k$	$p_1^* + \dots + p_k^*$

დაგროვილი სიხშირეები მიიღება სიხშირეების თანამიმდევრული აჯამებით, დანყებული პირველი ინტერვალისა.

თვალსაწინოებისათვის უმჯობესია, ინტერვალური ვარიაციული მსკრივი წარმოვადგინოთ გრაფიკულად ე.წ. **ჰისტოგრამის** საშუალებით. ჰისტოგრამა აიგება შემდეგნაირად: აბსცისათა ღერძზე გადაიზომება კერძო ინტერვალეები და თითოეულ მათგანზე აიგება ოთხკუთხედი, რომლის ფართობი უდრის მოცემული ინტერვალის სიხშირეს. ორდინატთა ღერძზე გადაზომილი მაჩვენებლის მიხედვით ჰისტოგრამა იყოფა ორ ტიპად: 1) ფარდობითი სიხშირეების ჰისტოგრამა, ანუ, როგორც მას ზოგჯერ უწოდებენ, ნორმირებული ჰისტოგრამა და 2) პროცენტებში გამოსახული სიხშირეების ჰისტოგრამა (პროცენტული ჰისტოგრამა). ამ შემთხვევაში, ორდინატთა ღერძზე გადაიზომება $P_i^* \cdot 100$ სიდიდეები.

ამ ორი ტიპის ჰისტოგრამის გამოყენების უპირატესობა ისაა, რომ ისინი გვაძლევენ ერთი და იგივე ინტერვალეებზე აგებული სხვადასხვა ჰისტოგრამების შედარების საშუალებას.

მაგალითი. მოცემულია 5 წლამდე ბავშვების სისხლის ნაკადის სიჩქარე (წმ-ში) 7 10 9 6 7 9 10 7 9 6 12 8 9 10 9 8 10 8 12 11 8 9 7 10 6 11 9 11 10 11 9 11 10 8 9 7 8 9 8 10 ($n = 40$). ავაგოთ ჰისტოგრამა, ემპირიული განაწილების სიმკვრივისა და განაწილების ფუნქციების გრაფიკები.

შევადგინოთ სიხშირული ცხრილი. ამისათვის განვსაზღვროთ კერძო ინტერვალის სიგრძე:

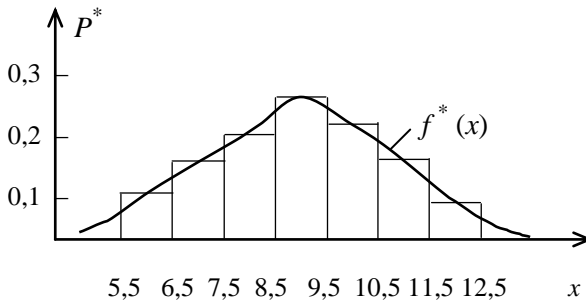
$$h = \frac{12 - 6}{1 + 3,321 \lg 40} = \frac{6}{1 + 5,312} = 0,95 \approx 1.$$

$x^{(0)} = 6 - \frac{1}{2} = 5,5$. სისშირულ ცხრილს ექნება შემდეგი სახე:

ინტერ- ვალები	ინტერ- ვალების საშ. მნი- შვნელობა x'_i	სიხ- ში- რეები m_i	ფარ- დობი- თი სიხშ. p_i^*	დაგროვილი სიხშირეები $\sum_i m_i$	დაგროვილი ფარდობითი სიხშირეები $\sum_i p_i^*$
[5,5;6,5[6	3	0,075	3	0,075
[6,5;7,5[7	5	0,125	8	0,200
[7,5;8,5[8	7	0,175	15	0,375
[8,5;9,5[9	10	0,25	25	0,625
[9,5;10,5[10	8	0,20	33	0,825
[10,5;11,5[11	5	0,125	38	0,95
[11,5;12,5]	12	2	0,05	40	1,00

ცხრილი სწორადაა შედგენილი, რადგან $\sum_{i=1}^7 m_i = 40$ და

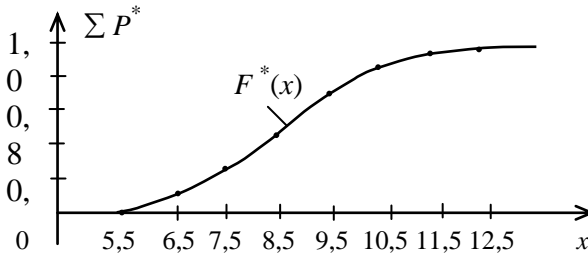
$\sum_{i=1}^7 p_i^* = 1$. ავაგოთ ნორმირებული ჰისტოგრამა.



თუ მრავალკუთხედის ზედა გვერდების შუა წერტილებს შევაერთებთ მრუდით, მაშინ ეს მრუდი მიახლოებით გვაძლევს წარმოდგენას ემპირიული განაწილების სიმკვრივის $f^*(x)$ ფუნქციის სახეზე. ჩვენი მაგალითის ემპირიული განაწილების სიმკვრივის

გრაფიკული გამოსახულებიდან ჩანს, რომ ბავშვების სისხლის ნაკადის სიჩქარის მნიშვნელობები დაახლოებით ნორმალურადაა განაწილებული.

თუ აბსცისათა ღერძზე გადაიზომება კერძო ინტერვალები, ხოლო ორდინატთა ღერძზე – დაგროვილი ფარდობითი სიხშირეები და მიღებულ წერტილებს შევაერთებთ, მაშინ მივიღებთ გრაფიკს, რომელსაც **კუმულატიური** გრაფიკი ეწოდება. პირველი ინტერვალის ქვედა ზღვრის მნიშვნელობას იღებენ ნულის ტოლად. ჩვენი მაგალითისათვის გვექნება:



ამ გზით მიღებული გრაფიკი წარმოადგენს შემთხვევითი სიდიდის ემპირიული განაწილების $F^*(x)$ ფუნქციის მიახლოებით მნიშვნელობას.

მონაცემების თვალსაჩინო წარმოდგენისათვის, ჰისტოგრამის გარდა, ზოგჯერ გამოიყენება **ფოთლებიანი ღეროების მსგავსი დიაგრამა**. ასეთი დიაგრამის ასაგებად, გავატაროთ ვერტიკალური წრფე და მის მარჯვნივ ჩავწეროთ ციფრების უმცირესი თანრიგები, რომლებსაც ფოთლები ეწოდებათ, ხოლო მარცხენა მხარეს – ციფრების ზედა თანრიგები, რომლებსაც ღეროები ეწოდებათ. დაკვირვებათა ღეროებად დაყოფა ხდება ციფრების თანრიგების რაოდენობის გათვალისწინებით. მაგალითად, თუ ღეროები წარმოდგენილია ერთთანრიგიანი ციფრებით, მაშინ დაყოფა ხდება ერთეულებით, ორთანრიგიანი ციფრების დროს – ათეულებით და ა.შ. აქვე უნდა შევნიშნოთ, რომ ასეთი დაყოფა პირობითია და შესაძლებელია დაყოფის სხვა კრიტერიუმის გამოყენება.

მაგალითი. ვთქვათ, მოცემულია 32 პაციენტის ასაკი (წლ.): 60, 72, 25, 81, 32, 80, 61, 73, 34, 65, 74, 45, 66, 63, 48, 62, 42, 65, 49, 54, 61, 52, 75, 57, 31, 58, 59, 69, 51, 82, 67, 72. ავაგოთ ფოთლებიანი ღეროების მსგავსი დიაგრამა.

ღეროები	ფოთლები
2	5
3	241
4	5892
5	427891
6	0156325197
7	23452
8	102

ფოთლებიანი ღეროების მსგავსი დიაგრამა ჰისტოგრამის მსგავსია, მაგრამ აქ შენარჩუნებულია ინდივიდუალური მნიშვნელობები. გარდა ამისა, მისი საშუალებით ადვილად დგინდება ამონარჩევის მინიმალური და მაქსიმალური მნიშვნელობები; რომელი მონაცემი გვხდება ყველაზე ხშირად; რომელი ყველაზე იშვიათად და სხვა. ნაკლი ის არის, რომ მისი გამოყენება დიდი მოცულობის ამონარჩევისათვის ნაკლებად მოსახერხებელია.

თუ მონაცემები მოცემულია პროცენტებში ან ფარდობითი სიხშირეების სახით, მაშინ მათი თვალსაჩინო წარმოდგენისთვის იყენებენ **წრიულ დიაგრამას**. წრიული დიაგრამის ასაგებად საჭიროა ვიპოვოთ ცენტრული კუთხეები შემდეგნაირად:

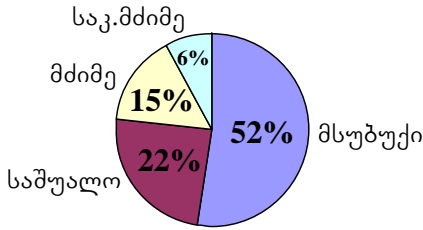
$$\varphi = 360 \cdot p^{\circ} \quad \text{ან} \quad \varphi = 360 \cdot \frac{a}{100},$$

სადაც p° – ფარდობითი სიხშირეა, a – პროცენტი. ზოგჯერ წრიულ დიაგრამებს მოცულობით სახეს აძლევენ.

მაგალითი. დიზენტერიით დაავადებულ პაციენტთა რაოდენობა 90-ია. აქედან, 47 პაციენტს (52%) აღენიშნება მსუბუქი, 22-ს (24%) საშუალო, 14-ს (16%) მძიმე და 7-ს (8%) საკმაოდ მძიმე ფორმა. მონაცემები წარმოვადგინოთ წრიული დიაგრამის სახით. ამისათვის გამოვთვალოთ ცენტრული კუთხეები.

$$\begin{aligned} \varphi_1 &= 360 \cdot 0,52 = 187^{\circ}, & \varphi_3 &= 360 \cdot 0,16 = 58^{\circ}, \\ \varphi_2 &= 360 \cdot 0,24 = 86^{\circ}, & \varphi_4 &= 360 \cdot 0,08 = 29^{\circ}. \end{aligned}$$

წრიულ დიაგრამას ექნება შემდეგი სახე:

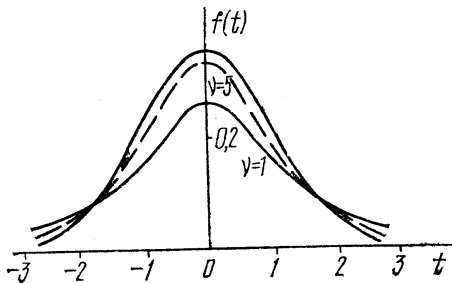


7. მათემატიკურ სტატისტიკაში გამოყენებული ძირითადი განაწილების კანონები

სტიუდენტის განაწილება. ვთქვათ, X შემთხვევითი სიდიდე განაწილებულია ნორმალურად. განაწილების სიმკვრივის ფორმულაში, როგორც ვიცით, არგუმენტად შედიან გენერალური ერთობლიობის მათემატიკური ლოდინი a და საშუალო კვადრატული გადახრა s , რომლებიც, როგორც წესი, უცნობები არიან. ამონარჩევის დროს ჩვენ ვსარგებლობთ მათი შეფასებებით, კერძოდ, საშუალო არითმეტიკულით \bar{x} და საშუალო კვადრატული გადახრით σ_x . აქედან გამომდინარე, ინგლისელმა მათემატიკოსმა კ. გოსეტმა იპოვა

$$t = \frac{\bar{x} - a}{\sigma_x} \sqrt{n}$$

სიდიდის განაწილების კანონი, სადაც გენერალური პარამეტრი a შეცვლილია მისი შეფასებით σ_x . აღმოჩნდა, რომ ამონარჩევის საშუალოსა და გენერალურ საშუალოს შორის სხვაობის განაწილების სიმკვრივის მრუდსა და ფუნქციას აქვს შემდეგი სახე:



$$f(t) = \frac{1}{\sqrt{\pi v}} \frac{\Gamma\left(\frac{v+1}{2}\right)}{\Gamma\left(\frac{v}{2}\right)} \left(1 + \frac{t^2}{v}\right)^{-\frac{v+1}{2}}, \quad -\infty < t < \infty,$$

სადაც, $\Gamma(\cdot)$ ნარმოადგენს გამა ფუნქციას. $f(t)$ ფუნქციას უნოდებენ სტიუდენტის განაწილებას (t -განაწილება) $v=n-1$ თავისუფლების ხარისხით.

სტიუდენტის განაწილების სიმკვრივე a და σ პარამეტრებზე არაა დამოკიდებული. ის განისაზღვრება მხოლოდ ამონარჩევის n მოცულობით. ეს თავისებურება ნარმოადგენს სტიუდენტის განაწილების ღირსებას. t -განაწილების მათემატიკური ლოდინი და დისპერსია ტოლია:

$$M(t) = 0, \quad D(t) = \frac{v}{v-2},$$

ხოლო მედიანა და მოდა ნულის ტოლია. თავისუფლების ხარისხის ზრდასთან ერთად, სტიუდენტის განაწილება სწრაფად უახლოვდება ნორმალურს. კერძოდ, როცა $v = 50$, სტიუდენტის განაწილება თითქმის არ განსხვავდება ნორმალურისგან. აქედან გამომდინარე, სტიუდენტის განაწილება ფართოდ გამოიყენება მცირე ამონარჩევების დროს (ცხრილი მოცემულია დანართში).

χ^2 (ხი-კვადრატი) განაწილება. განვიხილოთ n დამოუკიდებელი შემთხვევითი სიდიდეები X_1, X_2, \dots, X_n , რომლებიც ნორმალურად არიან განაწილებულნი a_1, a_2, \dots, a_n მათემატიკური ლოდინებითა და s_1, s_2, \dots, s_n საშუალო კვადრატული გადახრებით. ყოველი ამ შემთხვევითი სიდიდისთვის შემოვიტანოთ სტანდარტიზირებული შემთხვევითი სიდიდე

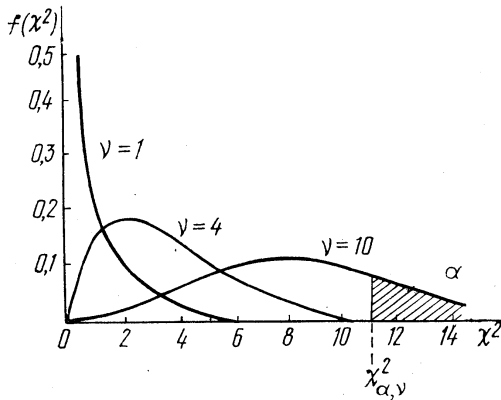
$$Z_i = \frac{x_i - a_i}{s_i}, \quad i = 1, 2, \dots, n.$$

სტანდარტიზირებული ცვლადების კვადრატების ჯამს

$$\chi^2 = Z_1^2 + Z_2^2 + \dots + Z_n^2$$

ენოდება χ^2 შემთხვევითი სიდიდე $v = n$ თავისუფლების ხარისხით. χ^2 შემთხვევითი სიდიდის განაწილების სიმკვრივის ფუნქციას აქვს შემდეგი სახე:

$$f(x^2) = \begin{cases} \frac{1}{2^{\frac{\nu}{2}} \Gamma\left(\frac{\nu}{2}\right)} (x^2)^{\frac{\nu}{2}-1} e^{-\frac{x^2}{2}}, & \text{როცა } x^2 \geq 0 \\ 0 & \text{, როცა } x^2 < 0 \end{cases}$$



როგორც განანილების სიმკვრივის ფუნქციიდან ჩანს, χ^2 განანილება არ არის დამოკიდებული არც X შემთხვევითი სიდიდის მათემატიკურ ლოდინზე და არც დისპერსიაზე, ის დამოკიდებულია მხოლოდ ამონარჩევის მოცულობაზე. ვიპოვოთ მათემატიკური ლოდინი და დისპერსია:

$$M(\chi^2) = M\left(\sum_{i=1}^n Z_i^2\right) = \sum_{i=1}^n M(Z_i^2) = n = \nu,$$

$$D(\chi^2) = 2n = 2\nu.$$

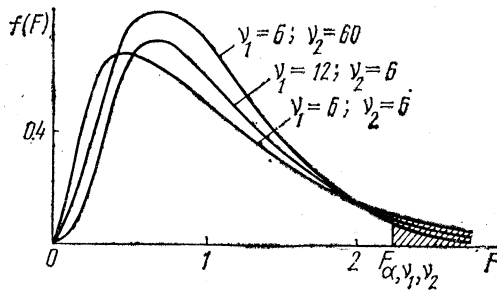
χ^2 განანილების სიმკვრივე რთული ფუნქციაა და მისი ინტეგრირება საკმაოდ რთულია, ამიტომ შედგენილია χ^2 განანილების სპეციალური ცხრილები (იხ. დანართი).

ფიშერის განანილება. ვთქვათ, მოცემულია დამოუკიდებელი შემთხვევითი სიდიდეები $X_1, X_2, \dots, X_n, Y_1, Y_2, \dots, Y_m$,

რომლებიც ნორმალურად არიან განაწილებული $N(0,1)$ პარამეტრებით. უგანზომილებო შემთხვევით სიდიდეს

$$F = \frac{m \sum_{i=1}^n X_i^2}{n \sum_{j=1}^m Y_j^2}$$

გააჩნია ფიშერის განაწილება (F -განაწილება) $\nu_1 = n$ მრიცხველისა და $\nu_2 = m$ მნიშვნელის თავისუფლების ხარისხით. განაწილების ფუნქციას აქვს შემდეგი სახე:



$$f(F) = \begin{cases} \frac{\Gamma\left(\frac{\nu_1 + \nu_2}{2}\right)}{\Gamma\left(\frac{\nu_1}{2}\right)\Gamma\left(\frac{\nu_2}{2}\right)} \left(\frac{\nu_1}{\nu_2}\right)^{\frac{\nu_1}{2}} \frac{F^{\frac{\nu_1}{2}-1}}{\left(1 + \frac{\nu_1}{\nu_2} F\right)^{\frac{\nu_1 + \nu_2}{2}}}, & \text{როცა } F > 0 \\ 0 & , \text{ როცა } F \leq 0 \end{cases}$$

როგორც ვხედავთ, F შემთხვევითი სიდიდის განაწილების სიმკვრივე რთული გამოსახულებაა და იგი დამოკიდებულია მხოლოდ ν_1 და ν_2 თავისუფლების ხარისხებზე. ეს გვაძლევს საშუალებას, რათა შევადგინოთ F -განაწილების სპეციალური ცხრილები (იხ. დანართი).

F -განაწილებას გააჩნია დადებითი ასიმეტრია და, ამონარჩევის მოცულობის ზრდასთან ერთად, უახლოვდება ნორმალურ განაწილებას.

8. ძირითადი სტატისტიკური მახასიათებლები

ცნობილია, რომ განაწილების კანონი შემთხვევით სიდიდეს სრულად ახასიათებს, მაგრამ ძალიან ხშირად განაწილების კანონი უცნობია. ამიტომ, ზოგჯერ გაცილებით უფრო მოსახერხებელია ზოგიერთი რაოდენობრივი მაჩვენებლების (სტატისტიკური მახასიათებლების) ცოდნა, რომლებიც მოგვანვდიან ინფორმაციას შემთხვევითი სიდიდის შესახებ. გარდა ამისა, შედარებითი ამოცანების გადაწყვეტისას, შესაძლებელია, რომ ორი შემთხვევითი სიდიდის განაწილების სიმკვრივის ფუნქციები იყოს ერთნაირი, მაგრამ ისინი შეიძლება განსხვავდებოდნენ ერთმანეთისგან სტატისტიკური მახასიათებლებით, რომლებიც პირობითად შეიძლება დავყოთ სამ ჯგუფად: განაწილების მდებარეობის, ცვალებადობისა და ფორმის მახასიათებლებად.

8.1. განაწილების მდებარეობის მახასიათებლები

განაწილების მდებარეობის ანუ ცენტრალური ტენდეციის საზომ მახასიათებლებს მიეკუთვნება საშუალო სიდიდეები, მოდა, მედიანა და კვანტილი. არსებობს საშუალო სიდიდეთა რამდენიმე სახე: საშუალო არითმეტიკული, საშუალო გეომეტრიული, საშუალო კვადრატული, საშუალო კუბური და სხვ. მათგან ყველაზე უფრო გავრცელებულია საშუალო არითმეტიკული.

ეთქვათ, მოცემულია X შემთხვევითი სიდიდის x_1, x_2, \dots, x_n ამონარჩევი. მაშინ საშუალო არითმეტიკული განისაზღვრება ფორმულით:

$$\bar{x} = \frac{x_1 + x_2 + \dots + x_n}{n} = \frac{1}{n} \sum_{i=1}^n x_i .$$

ამრიგად, საშუალო არითმეტიკული არის ამონარჩევის თითოეულ მონაცემზე მოსული ჯამურ დაკვირვებათა წილი, ანუ საშუალო რაოდენობა. თუ მონაცემები მოცემულია ვარიაციული

მნიშვნის სახით, მაშინ $\bar{x} = \frac{1}{n} \sum_{i=1}^k m_i x_i'$ და მას ანონილ (შენონილ)

საშუალო სიდიდეს უწოდებენ. აქ, m_i – ვარიანტების სიხშირეებია,

x'_i – ინტერვალების საშუალო მნიშვნელობა, k – ინტერვალების რაოდენობა. თუ გვაქვს რამდენიმე ჯგუფის საშუალო სიდიდეები და გვინდა გამოვთვალოთ მათი საერთო საშუალო, მაშინ გვექნება:

$$\bar{x} = \frac{n_1 \bar{x}_1 + n_2 \bar{x}_2 + \dots + n_k \bar{x}_k}{n_1 + n_2 + \dots + n_k} = \frac{\sum_{i=1}^k n_i \bar{x}_i}{\sum_{i=1}^k n_i}.$$

საშუალო არითმეტიკულის მიახლოებითი მნიშვნელობა შეიძლება ასე განისაზღვროს:

$$\bar{x} \approx \frac{x_{\min} + x_{\max}}{2},$$

სადაც, x_{\min} და x_{\max} ამონარჩევის მინიმალური და მაქსიმალური მნიშვნელობებია.

საშუალო არითმეტიკული გვიჩვენებს ერთობლიობის განლაგებას ჩვეულებრივ რიცხვით ღერძზე, ხოლო საშუალო გეომეტრიული – ლოგარითმულ ღერძზე. საშუალო გეომეტრიული გამოითვლება შემდეგი ფორმულით:

$$\bar{x}_g = \sqrt[n]{x_1 \cdot x_2 \cdot \dots \cdot x_n} \quad \text{ან} \quad \log \bar{x}_g = \frac{1}{n} \sum_{i=1}^n \log x_i.$$

საშუალო გეომეტრიულს იყენებენ იმ შემთხვევაში, როცა დაკვირვებათა რაოდენობა ხასიათდება გეომეტრიული პროპორციებით. როგორც საშუალო გეომეტრიულის ფორმულიდან ჩანს, მისი გამოყენება შესაძლებელია, თუ სრულდება პირობა $x_i > 0$, $i=1,2,\dots,n$.

საშუალო მნიშვნელობა არის ერთ-ერთი ძირითადი სტატისტიკური მახასიათებელი, რომელიც გამოითვლება ყველა მონაცემების მეშვეობით და შეიცავს მეტ ინფორმაციას, ვიდრე მდებარეობის სხვა მახასიათებლები. საშუალოს გააჩნია მთელი რიგი მნიშვნელოვანი თვისებები.

1. თუ ამონარჩევის მნიშვნელობებს შევამცირებთ ან გავზრდით რაიმე c მუდმივით, მაშინ საშუალო მნიშვნელობაც მცირდება ან იზრდება იგივე c სიდიდით. მართლაც,

$$\bar{x} = \frac{1}{n} \sum_{i=1}^n (x_i + c) = \frac{1}{n} \sum_{i=1}^n x_i + \frac{1}{n} \sum_{i=1}^n c = \bar{x} + \frac{1}{n} nc = \bar{x} + c.$$

2. თუ ამონარჩევის ყოველ მნიშვნელობას გავამრავლებთ ან გავყოფთ რაიმე c მუდმივზე, მაშინ საშუალო მნიშვნელობაც c -ჯერ იცვლება. მართლაც,

$$\bar{x} = \frac{1}{n} \sum_{i=1}^n cx_i = \frac{c}{n} \sum_{i=1}^n x_i = c\bar{x}.$$

$$\bar{x} = \frac{1}{n} \sum_{i=1}^n \frac{x_i}{c} = \frac{1}{cn} \sum_{i=1}^n x_i = \frac{\bar{x}}{c}.$$

3. ამონარჩევის თითოეული მნიშვნელობის საშუალო სიდიდესთან სხვაობის ჯამი ნულის ტოლია. მართლაც,

$$\begin{aligned} \sum_{i=1}^n (x_i - \bar{x}) &= \sum_{i=1}^n x_i - \sum_{i=1}^n \bar{x} = \sum_{i=1}^n x_i - n\bar{x} = \sum_{i=1}^n x_i - \frac{n}{n} \sum_{i=1}^n x_i = \\ &= \sum_{i=1}^n x_i - \sum_{i=1}^n x_i = 0 \end{aligned}$$

მიუხედავად იმისა, რომ საშუალო არითმეტიკული წარმოადგენს ერთ-ერთ უმნიშვნელოვანეს მახასიათებელს, მას გააჩნია ნაკლი. კერძოდ, იგი ძალზე მგრძობიარეა ისეთი მონაცემების გაზრდაზე ან შემცირებაზე, რომლებიც მკვეთრად განსხვავდებიან ძირითადი მონაცემებისაგან. ამრიგად, რანჟირებული ვარიაციული მწკრივის კიდურა მაჩვენებლები საშუალო არითმეტიკულზე ახდენენ ძლიერ ზეგავლენას, რაც მეტად არახელსაყრელია იმის გამო, რომ ეს მაჩვენებლები არ წარმოადგენს ტიპიურს ამ შემთხვევითი სიდიდისთვის. მაგალითად, თუ ამონარჩევის რომელიმე მნიშვნელობა იცვლება c ერთეულით, მაშინ საშუალო მნიშვნელობა იცვლება $\frac{c}{n}$ ერთეულით.

აქედან გამომდინარე, ბევრ შემთხვევაში ერთობლიობის განზოგადებულ მახასიათებლად მიზანშეწონილია **სტრუქტურული საშუალოების**, კერძოდ მედიანისა და მოდას გამოყენება.

მედიანა M_e არის საშუალო მნიშვნელობა, რომლის მიმართ სტატისტიკური მწკრივი იყოფა ორ თანაბარ ნაწილად და იგი დაკვირვებათა რანჟირებული (დალაგებული ზრდადობით ან კლებადობით) მწკრივის შუაში მდგარი ელემენტის მნიშვნელობის ტოლია. კერძოდ, თუ n კენტია, მაშინ მედიანა ტოლია რანჟირებული მწკრივის შუაში მყოფი ელემენტისა, თუ n ლუნია, მაშინ – რან-

ჟირებული მწკრივის შუაში მყოფი ორი მნიშვნელობის საშუალო არითმეტიკულისა. ე.ი. გვექნება:

$$M_e = \begin{cases} x_{\left(\frac{n+1}{2}\right)}, & \text{როცა } n \text{ კენტია} \\ \frac{1}{2} \left(x_{\left(\frac{n}{2}\right)} + x_{\left(\frac{n+1}{2}\right)} \right), & \text{როცა } n \text{ ლუნია} \end{cases}$$

ინტერვალური ანუ დაჯგუფებული მონაცემებისთვის მედიანა გამოითვლება შემდეგნაირად: პირველ როგში პოულობენ იმ ინტერვალს, სადაც მედიანა იმყოფება. ამისთვის საჭიროა გამოვთვალოთ ინტერვალების დაგროვილი სახშირეები იმ მნიშვნელობამდე, სანამ იგი $\frac{n}{2}$ -ს გაუტოლდება ან მეტი გახდება. ის

პირველი ინტერვალი, რომლის დაგროვილი სიხშირე $\frac{n}{2}$ -ზე მეტია, არის მედიანური ინტერვალი ანუ ინტერვალი, სადაც მედიანა იმყოფება. ამ შემთხვევაში მედიანა გამოითვლება შემდეგი ფორმულით:

$$M_e = a_0 + \frac{h}{m_e} \left(\frac{n}{2} - \sum_i m_i \right),$$

სადაც, a_0 – მედიანური ინტერვალის საწყისი (ქვედა ზღვრის) მნიშვნელობაა, h – ინტერვალის სიგრძე, m_e – მედიანური ინტერვალის სიხშირე, $\sum m_i$ – მედიანური ინტერვალის წინ (მარცხნივ) მდებარე ინტერვალის დაგროვილი სიხშირე.

მდებარეობის მაჩვენებლებს შორის მედიანა წარმოადგენს ერთ-ერთ მდგრად მახასიათებელს. მასზე არ მოქმედებს მონაცემების „დიდი“ და „მცირე“ მნიშვნელობები. მაგალითად, მედიანა არ შეიცვლება ამონარჩევის მაქსიმალური მნიშვნელობის გასამკვეცებითაც კი.

მოდა M_0 ეწოდება სტატისტიკური მწკრივის იმ მნიშვნელობას, რომელიც ყველაზე უფრო ხშირად გვხვდება ამონარჩევში. დაჯგუფებული (სიხშირული) მონაცემებისათვის მოდა გამოითვლება შემდეგნაირად: ჯერ საჭიროა განისაზღვროს მოდალური ინტერვალი. ინტერვალს, რომელსაც გააჩნია ყველაზე დიდი სიხშირე, ეწოდება მოდალური ინტერვალი ანუ ისეთი ინტერვალი, სადაც მოდაა მოთავსებული. მოდა გამოითვლება შემდეგი ფორმულით:

$$M_0 = a_0 + \frac{h(m_0 - m_1)}{2m_0 - m_1 - m_2},$$

სადაც, a_0 – მოდალური ინტერვალის საწყისი (ქვედა ზღვრის) მნიშვნელობაა, h – ინტერვალის სიგრძე, m_0 – მოდალური ინტერვალის სიხშირე, m_1 – მოდალური ინტერვალის წინ მდებარე (მარცხნივ მყოფი) ინტერვალის სიხშირე, m_2 – მოდალური ინტერვალის შემდგომი (მარჯვნივ მყოფი) ინტერვალის სიხშირე.

მცირე ამონარჩევებისთვის მოდა შეიძლება იყოს არასტაბილური. სამაგიეროდ, დიდი ამონარჩევის დროს იგი მდგრადი მახასიათებელია და რეაგირებს ამონარჩევის ელემენტების მხოლოდ ზოგიერთ ცვლილებაზე და არა ყოველგვარზე, როგორც საშუალო არითმეტიკული.

მდებარეობის სტრუქტურულ მაჩვენებლებს მიეკუთვნება **კვანტილი**. კვანტილი ზოგადი ცნებაა და საწყის მონაცემებს რიცხვითი ლერძის მიმართ ყოფს გარკვეული პროპორციებით. კვანტილის შემადგენელი ნაწილებია კვარტილი, დეცილი და პროცენტილი (პერცენტილი).

კვარტილი არის სამი Q_1 , Q_2 და Q_3 მნიშვნელობა, რომლებიც რანჟირებულ მწკრივს ყოფენ ოთხ თანაბარ ნაწილად (კვარტებად). **კვინტილი** (კვინტა-ხუთი) არის ოთხი K_1, K_2, K_3 და K_4 მნიშვნელობები, რომლებიც მოცემულ ერთობლიობას ყოფენ ხუთ თანაბარ ნაწილად. ცხრა D_1, D_2, \dots, D_9 **დეცილი** მწკრივს ყოფს ათ თანაბარ ნაწილად, ხოლო 99 **პერცენტილი** P_1, P_2, \dots, P_{99} მწკრივს ყოფენ 100 თანაბარ ნაწილად.

პრაქტიკაში, ძირითადად, გამოიყენება $P_3, P_{10}, P_{25}, P_{50}, P_{75}, P_{90}$ და P_{97} პერცენტილები. მაგალითად, P_{25} და P_{75} პერცენტილი შეესაბამება Q_1 და Q_3 კვარტილებს და მათ შორის მოთავსებულია ამონარჩევის წევრთა 50%. P_{50} პერცენტილი, რომელიც Q_2 კვარტილს შეესაბამება, მედიანას ტოლია, ე.ი. $P_{50} = M_e$. გარდა ამისა, პერცენტილი ძალზე მოსახერხებელია მონაცემების განზოგადებისთვის. მაგ. თუ $P_5 = 10,7$, ხოლო $P_{15} = 16,8$, შეიძლება ითქვას, რომ მონაცემების 5% ნაკლებია 10,7 სიდიდეზე, ხოლო 10% იმყოფება 10,7-სა და 16,8-ს შორის.

რადგან სხვადასხვა კვანტილებს შორის არსებობს გარკვეული ურთიერთკავშირები,

$$D_i = P_{10i}, \quad i = 1, 2, \dots, 9; \quad K_j = P_{20j}, \quad j = 1, 2, 3, 4; \quad Q_k = P_{25k}, \quad k = 1, 2, 3,$$

ამიტომ საკმარისია ვიცოდეთ პერცენტის მნიშვნელობა, რომელიც გამოითვლება შემდეგი ფორმულით:

$$P_j = x_{j3} + \frac{h}{m_p} \left(k - \sum_i m_i \right),$$

სადაც, x_{j3} – ინტერვალის ქვედა ზღვრული მნიშვნელობაა, რომელიც შეიცავს P_j პერცენტის; h – ინტერვალის სიგრძეა;

$k = \frac{L_j n}{100}$, n – დაკვირვებათა საერთო რაოდენობაა; L_j – პერცენტის რიგი (მაგალითად, P_{25} და P_{75} პერცენტების რიგი ტოლია შესაბამისად 25% და 75%); m_p – ინტერვალის სიხშირე, რომელიც

პერცენტის შეიცავს; $\sum_i m_i$ – დაგროვილი სიხშირე.

მაგალითი. ვთქვათ, მოცემულია სიხშირული ცხრილი

x_i	5	6	7	8	9	10	11	12
m_i	4	7	13	15	7	9	6	3
$\sum m_i$	4	11	24	39	46	55	61	64

საჭიროა ვიპოვოთ 50%-იანი პერცენტული P_{50} . გამოვთვალოთ k .

$k = \frac{50 \cdot 64}{100} = 32$. ეს მნიშვნელობა მეტია $\sum_i m_i = 24$ -ზე და

ნაკლებია $\sum_i m_i = 39$ -ზე. ამიტომ x_{j3} ვეძებთ მე-7 და მე-8 ინტერვალის (კლასის) მნიშვნელობების შუა, ე.ი.

$$x_{j3} = \frac{1}{2}(7 + 8) = 7,5; m_p = 15; P_{50} = 7,5 + \frac{32 - 24}{15} = 8,03.$$

ანალოგიურად გამოითვლება ინტერვალური მწკრივისთვისაც.

x	[2500;2600[[2600;2700[[2700;2800[[2800;3000[
m_i	2	5	13	20
$\sum m_i$	2	7	20	40

x	[3000;3100[[3100;3200[[3200;3300[[3300;3400[
m_i	16	17	4	3
$\sum m_i$	56	73	77	80

$$k = \frac{50 \cdot 80}{100} = 40, \text{ ეს შეესაბამება მე-4 ინტერვალს, ე.ი. } x_{j3} = 2800;$$

$$m_p = 20; h = 100; P_{50} = 2800 + 100 \frac{40 - 20}{20} = 2900.$$

ბოლო წლებში განსაკუთრებული პოპულარობით სარგებლობს მონაცემთა სიმრავლის გრაფიკული დახასიათება – **ბოქსპლოტი**, რომელიც მდგრადი საზომების Q_1 , Q_3 , IQR , მედიანას მეშვეობით საშუალებას იძლევა, გამოიკვეთოს მონაცემთა სიმრავლის ზოგიერთი თავისებურებანი, კერძოდ, განლაგების ცენტრი, გაფანტულობის გავრცობა და სიმეტრიულობიდან გადახრა. და ბოლოს, მოხდეს ამოვარდნილი დაკვირვებების იდენტიფიკაცია. IQR სიდიდე განისაზღვრება ფორმულით: $IQR = Q_3 - Q_1$ და მას კვარტილთაშორისო გაბნევის დიაპაზონი ეწოდება (IQR – *Interquartil range*).

8.2. განაწილების ცვალებადობის მახასიათებლები

დაკვირვებათა მწკრივის ერთ-ერთ ძირითად მაჩვენებელს წარმოადგენს მისი ცვალებადობის დონე. ერთი და იგივე საშუალოების მქონე შემთხვევითი სიდიდეები ერთმანეთისაგან შეიძლება განსხვავდებოდნენ ცვალებადობის ანუ ვარიაციის დონით. ამიტომ საშუალო მაჩვენებლებთან ერთად, სტატისტიკური მწკრივის სრული დახასიათებისათვის უნდა განისაზღვროს ვარიაციის მახასიათებლებიც.

სტატისტიკური მწკრივის ცვალებადობის მარტივ მაჩვენებელს წარმოადგენს **გაბნევის დიაპაზონი** R , რომელიც განისაზღვრება მოცემული x_1, x_2, \dots, x_n ამონარჩევის მაქსიმალური და მინიმალური მნიშვნელობების სხვაობით, ე.ი. $R = x_{\max} - x_{\min}$. იგი გვიჩვენებს, თუ რა დიაპაზონში იცვლება შემთხვევითი სიდიდე X . გაბნევის დიაპაზონის სიდიდეზე დიდ გავლენას ახდენს არტიფაქტური მნიშვნელობები. გარდა ამისა, ის არ შეიცავს ინფორმაციას, თუ როგორაა გაბნეული დანარჩენი მონაცემები მაქსიმალურ და მინიმალურ მნიშვნელობებს შორის.

სტატისტიკური მწკრივის ცვალებადობის მნიშვნელოვან მახასიათებელს წარმოადგენს **დისპერსია** (ლათინური სიტყვა

dispersio-გაფანტვა) ან ვარიანსა (ინგლისური სიტყვა-*variance*-ვარიაცია). დისპერსია აღინიშნება σ_x^2 ან $\text{Var}(X)$ სიმბოლოებით და წარმოადგენს საშუალოდან გადახრების კვადრატების საშუალო არითმეტიკულს

$$\sigma_x^2 = \frac{1}{n-1} \sum_{i=1}^n (x_i - \bar{x})^2,$$

ან

$$\sigma_x^2 = \frac{1}{n-1} \left[\sum_{i=1}^n x_i^2 - \frac{1}{n} \left(\sum_{i=1}^n x_i \right)^2 \right] = \frac{1}{n-1} \left[\sum_{i=1}^n x_i^2 - n \bar{x}^2 \right].$$

ვარიაციული მსკრივისთვის გვექნება:

$$\sigma_x^2 = \frac{1}{n-1} \sum_{i=1}^k m_i (x'_i - \bar{x})^2.$$

იმისათვის, რომ დისპერსიის შეფასება იყოს გადაუადგილებადი, იგი იყოფა არა n -ზე, არამედ $(n-1)$ სიდიდეზე (იხ. § 9.1).

დისპერსია შემთხვევითი სიდიდეა, რომელიც თავისივე საშუალო მნიშვნელობის ირგვლივ გაფანტვის მახასიათებელს წარმოადგენს და მისი განზომილება შემთხვევითი სიდიდის განზომილების კვადრატის ტოლია. ისევე როგორც საშუალო არითმეტიკულს, დისპერსიასაც გააჩნია შემდეგი ძირითადი თვისებები:

1. თუ ამონარჩევის ყოველ მნიშვნელობას შევამცირებთ ან გავზრდით რაიმე c მუდმივით, ამით დისპერსიის მნიშვნელობა არ იცვლება. მართლაც,

$$\begin{aligned} \frac{1}{n-1} \sum_{i=1}^n [(x_i + c) - (\bar{x} + c)]^2 &= \frac{1}{n-1} \sum_{i=1}^n \left[(x_i + c) - \left(\frac{1}{n} \sum_{i=1}^n (x_i + c) \right) \right]^2 = \\ &= \frac{1}{n-1} \sum_{i=1}^n \left[x_i + c - \left(\frac{1}{n} \sum_{i=1}^n x_i \right) - \frac{nc}{n} \right]^2 = \frac{1}{n-1} \sum_{i=1}^n (x_i + c - \bar{x} - c)^2 = \\ &= \frac{1}{n-1} \sum_{i=1}^n (x_i - \bar{x})^2. \end{aligned}$$

ანალოგიურად მტკიცდება სხვაობის დროს.

2. თუ ამონარჩევის თითოეულ მნიშვნელობას გავყოფთ ან გავამრავლებთ ერთი და იგივე c მუდმივზე, მაშინ დისპერსია შესაბამისად მცირდება ან იზრდება c^2 სიდიდით. მართლაც,

$$\sigma_x^2 = \frac{1}{n-1} \sum_{i=1}^n \left(\frac{x_i}{c} - \frac{\bar{x}}{c} \right)^2 = \frac{1}{c^2(n-1)} \sum_{i=1}^n (x_i - \bar{x})^2 = \frac{\sigma_x^2}{c^2};$$

$$\sigma_x^2 = \frac{1}{n-1} \sum_{i=1}^n (x_i c - \bar{x} c)^2 = \frac{1}{n-1} \sum_{i=1}^n (x_i - \bar{x})^2 c^2 = \sigma_x^2 c^2.$$

დისპერსიის გარდა ვარიაციის მნიშვნელოვან მახასიათებელს წარმოადგენს **საშუალო კვადრატული გადახრა**

$\sigma_x = \sqrt{\sigma_x^2}$. ეს მახასიათებელი, რომელსაც ზოგჯერ **სტანდარტულ გადახრას** უწოდებენ, გამოისახება იმავე ფიზიკურ ერთეულში, რომლითაც ამონარჩევის ელემენტებია გამოსახული. ამიტომ ის უფრო მოსახერხებელია, ვიდრე დისპერსია. რაც უფრო მეტად განიცდის პარამეტრი ცვლილებას, მით უფრო დიდია სტანდარტული გადახრა და პირიქით, რაც უფრო სუსტია ეს ცვლილება, მით უფრო მცირეა საშუალო კვადრატული გადახრა.

საშუალო კვადრატული გადახრის მიახლოებითი მნიშვნელობა განისაზღვრება შემდეგი ფორმულით:

$$\sigma_x \approx \frac{x_{\max} - x_{\min}}{k},$$

სადაც, k კოეფიციენტია, რომლის მნიშვნელობა დამოკიდებულია ამონარჩევის n მოცულობაზე ისე, როგორც ეს შემდეგ ცხრილშია წარმოდგენილი.

n	2-5	6-15	16-49	50-200	201-1000	>1000
k	2	3	4	5	6	7

დისპერსია და საშუალო კვადრატული გადახრა აბსოლუტური სიდიდეებია და იზომება იგივე ფიზიკური ერთეულით, რითაც ამონარჩევის ელემენტებია გაზომილი. ამიტომ, როდესაც საჭიროა სხვადასხვა ერთეულებში გამოსახული ცვლადების შედარება, უმჯობესია გამოვიყენოთ ვარიაციის ფარდობითი მაჩვენებელი. ერთ-ერთი ასეთი მაჩვენებელია **ვარიაციის კოეფი-**

ციენტი, რომელიც პირსონმა შემოიტანა და განისაზღვრება შემდეგი ფორმულით:

$$V = \frac{\sigma_x}{\bar{x}} 100\%,$$

სადაც, σ_x – საშუალო კვადრატული გადახრაა, \bar{x} – საშუალო არითმეტიკული. ვარიაციის კოეფიციენტი უგანზომილებო სიდიდეა, რომელიც გამოსახულია პროცენტებში. მისი გამოყენების დროს საჭიროა მხედველობაში გვექონდეს შემდეგი გარემოებანი: ჯერ ერთი, იგი მკვეთრად იზრდება ასიმეტრიული განაწილების დროს; ამის გარდა, თუ მოცემულ პარამეტრს აქვს მცირე ან უბრალოდ უარყოფითი მნიშვნელობები, მაშინ ვარიაციის კოეფიციენტმა შეიძლება 100%-ს გადააჭარბოს. ეს ყველაფერი რამდენადმე ამცირებს ამ მაჩვენებლის ფასს და გარკვეულად ზღუდავს მის გამოყენებას პრაქტიკაში.

8.3. განაწილების ფორმის მახასიათებლები

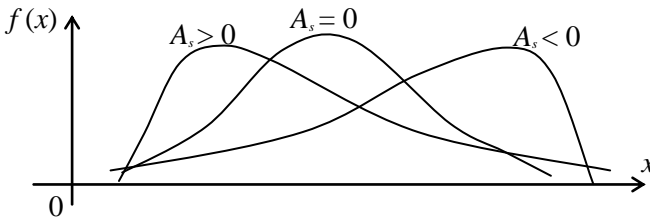
განაწილების ფორმის მახასიათებლებიდან განვიხილოთ ექსცესა და ასიმეტრია. ვიზუალურად X შემთხვევითი სიდიდის განაწილების ასიმეტრია შეიძლება დადგინდეს განაწილების სიმკვრივის ფუნქციის ან ჰისტოგრამის საშუალებით. ყველა განაწილების წირს გააჩნია ასიმეტრიის სხვადასხვა ხარისხი. თუ შემთხვევითი სიდიდე თავისი მათემატიკური ლოდინის მიმართ სიმეტრიულადაა განაწილებული, მაშინ ასიმეტრია ნულის ტოლია. აქედან გამომდინარე, ნორმალური განაწილების მრუდს ასიმეტრია არ გააჩნია. ზოგადად, ასიმეტრიის დახასიათებისთვის მესამე რიგის ცენტრალური მომენტის საშუალებით გამოითვლება

ასიმეტრიის კოეფიციენტი:
$$A_s = \frac{\frac{1}{n} \sum_{i=1}^n (x_i - \bar{x})^3}{\sigma_x^3},$$

სადაც, σ_x^3 – საშუალო კვადრატული გადახრაა აყვანილი მესამე ხარისხში. თუ მოცემულია ვარიაციული მწკრივი, მაშინ

$$A_s = \frac{\frac{1}{n} \sum_{i=1}^k m_i (x'_i - \bar{x})^3}{\sigma_x^3},$$

როცა $A_s = 0$, განაწილების წირს არ გააჩნია ასიმეტრია. განაწილებას გააჩნია დადებითი ასიმეტრია, როცა $A_s > 0$ და უარყოფითი, როცა $A_s < 0$.



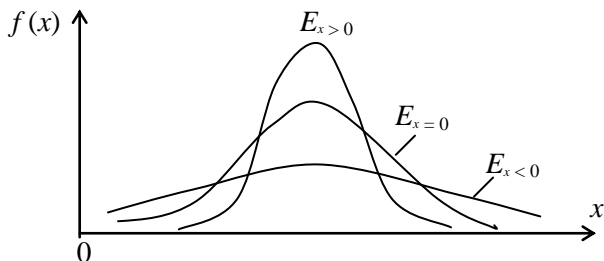
განაწილების მრუდის წამახვილების ხარისხის დასადგენად გამოიყენება **ექსცესის კოეფიციენტი** E_x , რომელიც გამოითვლება მეოთხე რიგის ცენტრალური მომენტის საშუალებით

შემდეგი ფორმულით:
$$E_x = \frac{\frac{1}{n} \sum_{i=1}^n (x_i - \bar{x})^4}{\sigma_x^4} - 3.$$

ვარიაციული მწკრივისათვის გვექნება:
$$E_x = \frac{\frac{1}{n} \sum_{i=1}^k m_i (x'_i - \bar{x})^4}{\sigma_x^4} - 3,$$

სადაც, σ_x^4 – საშუალო კვადრატული გადახრაა აყვანილი მეოთხე ხარისხში.

ნორმალური განაწილების მრუდისთვის ექსცესის კოეფიციენტი $E_s = 0$. იგი მიღებულია ეტალონად, რომელთანაც შედარდებიან სხვა განაწილების მრუდეები. მრუდს, რომელსაც უფრო მაღალი წვერო აქვს ვიდრე ნორმალურს, ე.ი. უფრო მახვილწვერიანია, შეესაბამება დადებითი ექსცესა, ხოლო მრუდს, რომელსაც უფრო დაბალი და ბრტყელი წვერო აქვს – უარყოფითი ექსცესა.



მაგალითი 1. მოცემულია კარდიოგენურ შოკში მყოფი პაციენტების სისტოლური წნევის მნიშვნელობები (x). გამოვთვალოთ სტატისტიკური მახასიათებლები. შევადგინოთ შემდეგი ცხრილი:

	x_i	$x_i - \bar{x}$	$(x_i - \bar{x})^2$	$(x_i - \bar{x})^3$	$(x_i - \bar{x})^4$
1	98	3,40	11,56	39,30	133,63
2	92	-2,60	6,76	-17,58	45,70
3	103	8,40	70,56	592,70	4978,71
4	85	-9,60	92,16	-884,74	8493,47
5	88	-6,60	43,56	-287,50	1897,47
6	94	-0,60	0,36	-0,22	0,13
7	96	1,40	1,96	2,74	3,84
8	105	10,40	108,16	1124,86	11698,59
9	95	0,40	0,16	0,064	0,026
10	90	-4,60	21,16	-97,34	447,75
Σ	946	-	356,40	472,28	27699,32

საშუალო არითმეტიკული $\bar{x} = \frac{1}{n} \sum_{i=1}^n x_i = \frac{946}{10} = 94,60$.

მედიანას განსაზღვრისათვის მოვახდინოთ მონაცემების რანჟირება: 85 88 90 92 94 95 96 98 103 105. რადგან $n = 10$, ამიტომ

$$M_e = \frac{1}{2} \left(x_{\left(\frac{n}{2}\right)} + x_{\left(\frac{n}{2}+1\right)} \right) = \frac{1}{2} (x_{(5)} + x_{(6)}) = \frac{1}{2} (94 + 95) = 94,5.$$

მოცემულ ამონარჩევს მოდა არ გააჩნია, რადგან რან-
ჟირებულ მწკრივში ერთნაირი სიდიდის მონაცემები არ გვხვდები-
ან. გაბნევის დიაპაზონი $R = x_{\max} - x_{\min} = 105 - 85 = 20$;

$$\text{დისპერსია } \sigma_x^2 = \frac{1}{n-1} \sum_{i=1}^n (x_i - \bar{x})^2 = \frac{356,40}{9} = 39,60;$$

$$\text{საშუალო კვადრატული გადახრა } \sigma_x = \sqrt{\sigma_x^2} = \sqrt{39,6} = 6,29;$$

$$\text{ვარიაციის კოეფიციენტი } V = \frac{\sigma_x}{\bar{x}} 100\% = \frac{6,29 \cdot 100}{94,6} = 6,65\%;$$

$$\text{ასიმეტრიის კოეფიციენტი } A_s = \frac{\frac{1}{n} \sum_{i=1}^n (x_i - \bar{x})^3}{\sigma_x^3} = \frac{472,28}{(6,29)^3} = 0,19;$$

ექსცესის კოეფიციენტი

$$E_x = \frac{\frac{1}{n} \sum_{i=1}^n (x_i - \bar{x})^4}{\sigma_x^4} - 3 = \frac{27699,32}{(6,29)^4} - 3 = -1,23.$$

მაგალითი 2. მე-6 თავში მოყვანილი მაგალითისთვის, სა-
დაც მონაცემები წარმოდგენილია სიხშირული ცხრილის ანუ ვარი-
აციული მწკრივის სახით, გამოვთვალოთ ძირითადი სტატისტიკუ-
რი მახასიათებლები. ამისათვის შევადგინოთ შემდეგი ცხრილი:

ინტერვ.	x'_i	m_i	$x'_i - \bar{x}$	$m_i(x'_i - \bar{x})^2$	$m_i(x'_i - \bar{x})^3$	$m_i(x'_i - \bar{x})^4$
[5,5;6,5[6	3	-2,95	26,11	-77,02	227,20
[6,5;7,5[7	5	-1,95	19,01	-37,07	72,30
[7,5;8,5[8	7	-0,95	6,32	-6,00	5,70
[8,5;9,5[9	10	0,05	0,03	0,00	0,00
[9,5;10,5[10	8	1,05	8,82	9,26	9,72
[10,5;11,5[11	5	2,05	21,01	43,08	88,31
[11,5;12,5]	12	2	3,05	18,61	56,75	173,07
Σ				99,91	11,00	576,30

$n = 40, k = 7$. საშუალო არითმეტიკული $\bar{x} = \frac{1}{n} \sum_{i=1}^k m_i x'_i =$
 $= \frac{1}{40} (3 \cdot 6 + 5 \cdot 7 + \dots + 2 \cdot 12) = 8,95$. მედიანური ინტერვალი არის

მე-4 ინტერვალი, რადგან $\sum m_i = 25 > \frac{n}{2} = 20$. აქედან გამომდინარე, მედიანა ტოლია:

$$M_e = a_0 + \frac{h}{m_e} \left(\frac{n}{2} - \sum m_i \right) = 8,5 + \frac{1(20-15)}{10} = 9,0.$$

მე-4 ინტერვალი მოდალური ინტერვალაცაა, რადგან მას გააჩნია უდიდესი სიხშირე ($m = 10$), ამიტომ

$$M_0 = a_0 + \frac{h(m_0 - m_1)}{2m_0 - m_1 - m_2} = 8,5 + \frac{1(10-7)}{2 \cdot 10 - 7 - 8} = 9,10;$$

$$\text{დისპერსია: } \sigma_x^2 = \frac{1}{n-1} \sum_{i=1}^k m_i (x'_i - \bar{x})^2 = \frac{99,91}{39} = 2,56;$$

$$\text{საშუალო კვადრატული გადახრა: } \sigma_x = \sqrt{\sigma_x^2} = \sqrt{2,56} = 1,60;$$

$$\text{გაბნევის დიაპაზონი: } R = x'_7 - x'_1 = 12 - 6 = 6;$$

$$\text{ვარიაციის კოეფიციენტი: } V = \frac{\sigma_x}{\bar{x}} 100\% = \frac{1,60 \cdot 100}{8,95} = 17,88\% ;$$

$$\text{ასიმეტრიის კოეფიციენტი: } A_s = \frac{\frac{1}{n} \sum_{i=1}^n (x_i - \bar{x})^3}{\sigma_x^3} = \frac{11,0}{(1,6)^3} = 0,067;$$

ექსცესის კოეფიციენტი:

$$E_x = \frac{\frac{1}{n} \sum_{i=1}^k m_i (x'_i - \bar{x})^4}{\sigma_x^4} - 3 = \frac{576,3}{(1,6)^4} - 3 = -0,80.$$

8.4. თვისებრივი მაჩვენებლების სტატისტიკური მახასიათებლები

თვისებრივი მონაცემების დროს საშუალო მნიშვნელობების გამოთვლა უაზრობაა. ამ შემთხვევაში, მის მაგიერ იყენებენ ფარდობით სიხშირეს ან პროცენტს. თუ ამონარჩევის მოცულობაა n , ხოლო რომელიმე p ნიშნის არსებობის რაოდენობა – m , მაშინ მისი ფარდობითი სიხშირე ტოლია:

$$p^* = \frac{m}{n}.$$

ალტერნატიული მონაცემების დროს, როდესაც ერთი m -განზომილებიანი p მაჩვენებელი დაპირისპირებულია მეორე q მაჩვენებელთან, მაშინ ფარდობითი სიხშირეები ტოლია:

$$p^* = \frac{m}{n}; \quad q^* = \frac{n-m}{n} = 1 - p^*.$$

ხოლო პროცენტებში:

$$p^* = \left(\frac{m}{n}\right)100\%; \quad q^* = 100 - p^* \%.$$

ცხადია, რომ $p^* + q^* = 1$ ან $p^* \% + q^* \% = 100$.

ცვალებადობის მაჩვენებლებიდან შეიძლება განისაზღვროს საშუალო კვადრატული გადახრა

$$\sigma_p = \sqrt{p^*(1-p^*)} = \sqrt{p^*q^*} \quad \text{ან} \quad \sigma_p = \sqrt{p^*\% (100 - p^*\%)}$$

ეს მაჩვენებლები ცვალებადობის თვალსაზრისით ერთნაირად ახასიათებს ორივე ალტერნატიულ მაჩვენებელს. თუ ალტერნატიული ჯგუფები წარმოდგენილია აბსოლუტური რიცხვებით, მაშინ

$$\sigma_p = \sqrt{np^*q^*}.$$

მიღებული სტანდარტული გადახრა σ_p დამოკიდებულია p^* -ზე და აღწევს მაქსიმალურ მნიშვნელობას, როცა $p^* = 0,5$, ხოლო ნულის ტოლია, როცა $p^* = 0$ ან $p^* = 1$. თუ ამონარჩევის მოცულობა საკმაოდ დიდია, მაშინ ცენტრალური ზღვართი თეორემებიდან გამომდინარე, p^* შეფასებას გააჩნია ნორმალური განაწილება, რასაც ვერ ვიტყვით მცირე ამონარჩევის დროს. მა-

თემატიკურ სტატისტიკაში მტკიცდება, რომ თუ $np^* > 5$ და $n(1-p^*) > 5$, მაშინ ფარდობითი სიხშირე p^* დაახლოებით ნორმალურადაა განაწილებული.

9. უცნობი პარამეტრების სტატისტიკური შეფასება

9.1. პარამეტრების შეფასების ცნება

ამონარჩევის მახასიათებლები – საშუალო არითმეტიკული, საშუალო კვადრატული გადახრა, დისპერსია და სხვა – შემთხვევითი სიდიდეებია, რომლებიც იცვლებიან თავიანთი გენერალური პარამეტრების – გენერალური საშუალოს, სტანდარტული გადახრის ან დისპერსიის გარშემო. ამონარჩევის მახასიათებლები განიხილებიან, როგორც მიახლოებითი მნიშვნელობები, რის გამოც აუცილებელია მათი შეფასება.

ვთქვათ, θ^* არის უცნობი θ პარამეტრის სტატისტიკური შეფასება. დავუშვათ, რომ n -განზომილებიანი ამონარჩევიდან ნაპოვნია θ_1^* შეფასება. გავიმეოროთ ცდა, ე.ი. გენერალური ერთობლიობიდან ამოვიღოთ იგივე განზომილების სხვა ამონარჩევი და მისი საშუალებით მივიღოთ θ_2^* შეფასება. თუ ცდას მრავალჯერ გავიმეორებთ, მივიღებთ $\theta_1^*, \theta_2^*, \dots, \theta_k^*$ შეფასებებს, რომლებიც ერთმანეთისაგან განსხვავდებიან. ამრიგად, θ^* შეფასება შეიძლება განვიხილოთ, როგორც შემთხვევითი სიდიდე, ხოლო $\theta_1^*, \theta_2^*, \dots, \theta_k^*$ რიცხვები, როგორც მისი შესაძლო მნიშვნელობები.

ჩამოვაცალიბოთ ის ძირითადი თვისებები, რომლებიც უნდა გააჩნდეს უცნობი პარამეტრის „კარგ“ შეფასებას. შეფასების სიზუსტისა და იმედიანობის თვალსაზრისით, სასურველია, რომ უცნობი პარამეტრის შეფასება θ^* შეძლებისდაგვარად მჭიდროდ იყოს კონცენტრირებული შესაფასებელი θ პარამეტრის ირგვლივ. ანუ, სხვა სიტყვებით რომ ვთქვათ, θ -ს გარშემო θ^* გაბნევა (გაფანტვა) უნდა იყოს უმცირესი. აქედან გამომდინარე, შეფასება

უნდა აკმაყოფილებდეს გადაუადგილებადობის, საფუძვლიანობისა და ეფექტურობის მოთხოვნებს.

გადაუადგილებადი ეწოდება ისეთ სტატისტიკურ θ^* შეფასებას, რომლის მათემატიკური ლოდინი უცნობი პარამეტრის ტოლია ამონარჩევის ნებისმიერი განზომილების დროს, ე.ი.

$$M(\theta^*) = \theta.$$

ხშირად, გადაუადგილებად შეფასებასთან ერთად, გამოიყენება ასიმპტოტიურად გადაუადგილებადი, ე.ი. ისეთი შეფასება, რომლისთვისაც ამონარჩევის მოცულობის გაზრდისას $M(\theta^*) \rightarrow \theta$.

გადაადგილებადი ეწოდება ისეთ შეფასებას, რომლის მათემატიკური ლოდინი შესაფასებელი პარამეტრის ტოლი არ არის, ე.ი. $M(\theta^*) \neq \theta$.

ეფექტური ეწოდება ისეთ სტატისტიკურ შეფასებას, როდესაც მოცემულ ამონარჩევს გააჩნია უმცირესი შესაძლო დისპერსია.

საფუძვლიანი ეწოდება შეფასებას, თუ ამონარჩევის განზომილების გაზრდისას, შეფასება რაღაც ალბათობით უახლოვდება შესაფასებელ პარამეტრს, ე.ი. $\lim_{n \rightarrow \infty} P(|\theta^* - \theta| < \varepsilon) = 1$.

მაგალითად, თუ გადაუადგილებადი შეფასების დისპერსია, როცა $n \rightarrow \infty$, მიისწრაფვის ნულისაკენ, მაშინ ასეთი შეფასება საფუძვლიანია.

ამრიგად, როდესაც ვეძებთ უცნობი პარამეტრის შეფასებას, უნდა გავითვალისწინოთ შეფასების ზემოთ მოყვანილი მოთხოვნები. განვიხილოთ მათემატიკური ლოდინისა და დისპერსიის შეფასებები.

მათემატიკური ლოდინისა და დისპერსიის შეფასებისათვის განვიხილოთ X გენერალური ერთობლიობიდან x_1, x_2, \dots, x_n ამონარჩევი, სადაც, x_i წარმოადგენენ დამოუკიდებელ და ერთნაირად განაწილებულ შემთხვევით სიდიდეებს. აღვნიშნოთ მათემატიკური ლოდინი $M[x_i] = a$ და დისპერსია $D[x_i] = s^2$. განვიხილოთ ამონარჩევიდან მიღებული საშუალო არითმეტიკულის \bar{x} მათემატიკური ლოდინი და დისპერსია:

$$M(\bar{x}) = M\left[\frac{1}{n} \sum_{i=1}^n x_i\right] = \frac{1}{n} \sum_{i=1}^n M(x_i) = \frac{na}{n} = a.$$

$$D(\bar{x}) = D\left[\frac{1}{n} \sum_{i=1}^n x_i\right] = \frac{1}{n^2} \sum_{i=1}^n D(x_i) = \frac{1}{n^2} ns^2 = \frac{s^2}{n}.$$

ე.ი. ამონარჩევის საშუალო არითმეტიკული არის გენერალური ერთობლიობის მათემატიკური ლოდინის გადაუადგილებადი შეფასება, ხოლო საშუალო არითმეტიკულის დისპერსია n -ჯერ მცირეა, ვიდრე შემთხვევითი X სიდიდის დისპერსია.

ამრიგად, თუ შემთხვევითი სიდიდე ნორმალურადაა განაწილებული $N(a, s)$ პარამეტრებით, მაშინ მათემატიკური ლოდინის გადაუადგილებად შეფასებას გააჩნია მინიმალური დისპერსია. აქედან გამომდინარე, საშუალო არითმეტიკული წარმოადგენს გენერალური ერთობლიობის მათემატიკური ლოდინის ეფექტურ შეფასებას.

განვიხილოთ დისპერსიის შეფასება. კერძოდ, დავადგინოთ, არის თუ არა ამონარჩევით მიღებული დისპერსია σ_x^2 გენერალური s^2 დისპერსიის გადაუადგილებადი შეფასება. ამისათვის განვსაზღვროთ დისპერსიის მათემატიკული ლოდინი:

$$M(\sigma_x^2) = M\left(\frac{1}{n} \sum_{i=1}^n x_i^2 - \bar{x}^2\right) = M\left(\frac{1}{n} \sum_{i=1}^n x_i^2\right) - M^2(\bar{x}) = M(X^2) - M^2(\bar{x}).$$

თუ გამოვიყენებთ დისპერსიის განსაზღვრის ფორმულას

$$D(X) = M(X^2) - M^2(\bar{x}), \quad \text{მაშინ გვექნება:}$$

$$D(X^2) = D(X) + M^2(\bar{x}) = s^2 + a^2, \quad M(\bar{x}^2) = D(\bar{x}) + M^2(\bar{x}) = \frac{s^2}{n} + a^2$$

$$\text{ე.ი. მივიღეთ: } M(\sigma_x^2) = s^2 + a^2 - \frac{s^2}{n} - a^2 = s^2 \frac{n-1}{n}.$$

აქედან გამომდინარეობს, რომ შერჩევითი დისპერსია σ_x^2 არის გენერალური s^2 დისპერსიის გადაუადგილებადი შეფასება.

გადაუადგილებად შეფასებას მივიღებთ, თუ განვიხილავთ σ_x^2 შესწორებულ მნიშვნელობას. ამისათვის, იგი უნდა გავამრავლოთ

$\left(\frac{n}{n-1}\right)$ სიდიდებზე, რომელსაც **ბესელის** შესწორებას უწოდებენ,

მაშინ გვექნება:

$$\hat{\sigma}_x^2 = \frac{n}{n-1} \sigma_x^2 = \frac{1}{n-1} \sum_{i=1}^n (x_i - \bar{x})^2.$$

მცირე ამონარჩევის დროს ბესელის შესწორება საგრძნობლად განსხვავდება ერთისაგან, ხოლო როცა ამონარჩევის n მოცულობა იზრდება, იგი სწრაფად მიისწრაფვის ერთისკენ და როცა $n > 50$ პრაქტიკულად განსხვავება $\hat{\sigma}_x^2$ და σ_x^2 შორის უმნიშვნელოა.

არსებობს პარამეტრთა შეფასების ორი მეთოდი: წერტილოვანი და ინტერვალური. პარამეტრის წერტილოვანი შეფასება განისაზღვრება ერთი მნიშვნელობით, ხოლო ინტერვალური – ორი რიცხვით. განვიხილოთ ეს შეფასებები.

9.2. პარამეტრთა შეფასების წერტილოვანი მეთოდები

განვიხილოთ უცნობი პარამეტრების წერტილოვანი შეფასების ძირითადი მეთოდები. ესენია მომენტთა მეთოდი და მაქსიმალური დასაჯერობის მეთოდი. აქედან ყველაზე მარტივი შეფასების მეთოდია – მომენტთა მეთოდი, რომელიც შემოიტანა ინგლისელმა სტატისტიკოსმა პირსონმა. ამ მეთოდის თეორიული საფუძველი ემყარება დიდ რიცხვთა კანონს, რომლის თანახმად დიდი მოცულობის მქონე ამონარჩევისთვის ამონარჩევის მომენტები გენერალური ერთობლიობის მომენტების ტოლია.

მომენტთა მეთოდის არსი შემდეგში მდგომარეობს: ხდება განაწილების თეორიული μ და ემპირიული m მომენტების გატოლება. თეორიულს ვუნოდებთ გენერალური ერთობლიობის მომენტებს, ხოლო ემპირიულს – ამონარჩევის საფუძველზე მიღებულ მომენტებს. ვთქვათ, საშუალოებით უნდა შეფასდეს არა ერთი, არამედ რამდენიმე k უცნობი პარამეტრი $\theta_1^*, \theta_2^*, \dots, \theta_k^*$. საჭიროა ვიპოვოთ გენერალური ერთობლიობის პირველი, მეორე და ა.შ. k -ური რიგის თეორიული მომენტები:

$$\mu_1(\theta_1^*, \theta_2^*, \dots, \theta_k^*), \mu_2(\theta_1^*, \theta_2^*, \dots, \theta_k^*), \dots, \mu_k(\theta_1^*, \theta_2^*, \dots, \theta_k^*)$$

და შემდეგ შესაბამისი ემპირიული მომენტები:

$$m_1(X_1, X_2, \dots, X_k), m_2(X_1, X_2, \dots, X_k), \dots, m_k(X_1, X_2, \dots, X_k).$$

თუ ამ მომენტებს ერთმანეთს გავუტოლებთ,

$$\mu_j(\theta_1^*, \theta_2^*, \dots, \theta_k^*) = m_j(X_1, X_2, \dots, X_k), \quad j = 1, 2, \dots, k,$$

მივიღებთ სისტემას, რომლის ამონახსნი $\theta_1^*, \theta_2^*, \dots, \theta_k^*$ იქნება უცნობი პარამეტრების სტატისტიკური შეფასება.

მომენტთა მეთოდი გამოიყენება სიმარტივეთ, მაგრამ მისი საშუალებით მიღებული შეფასებები ხშირად გადაადგილებადია და ნაკლებად ეფექტური. გამონაკლისს წარმოადგენს მხოლოდ ნორმალური განაწილების შემთხვევა, რომლის დროსაც მომენტთა მეთოდი იძლევა ეფექტურ და საფუძვლიან შეფასებებს.

მომენტთა მეთოდით მიღებული შეფასებების თვისებების შესწავლისას, ინგლისელმა მათემატიკოსმა ფიშერმა შემოგვთავაზა პარამეტრთა შეფასების უფრო საიმედო მეთოდი – **მაქსიმალური დასაჯერობის მეთოდი**. ამ მეთოდის ძირითადი არსი მდგომარეობს შემდეგში: ვთქვათ, მოცემულია X შემთხვევითი სიდიდე და მისი განაწილების სიმკვრივე $f(X, \theta)$, რომელიც დამოკიდებულია შესაფასებელ θ უცნობ პარამეტრზე.

თუ X_1, X_2, \dots, X_n დამოუკიდებელი შემთხვევითი სიდიდეებია, მაშინ **დასაჯერობის ფუნქცია** ეწოდება შემდეგ გამოსახულებას:

$$L = f(X_1, \theta) \cdot f(X_2, \theta) \cdot \dots \cdot f(X_n, \theta). \quad (9.1)$$

უცნობი θ პარამეტრის შესაფასებლად შევარჩიოთ ისეთი θ^* მნიშვნელობა, რომლის (9.1) ტოლობაში θ -ს მაგივრად ჩასმისას მივიღებთ L ფუნქციის მაქსიმალურ მნიშვნელობას. ასეთ θ^* შეფასებას უწოდებენ მაქსიმალური დასაჯერობის შეფასებას.

L ფუნქციის მაქსიმიზაციისას იგულისხმება, რომ X_1, X_2, \dots, X_n მნიშვნელობები დაფიქსირებულია, ხოლო ცვლად სიდიდეს წარმოადგენს θ პარამეტრი. თუ L ფუნქცია დიფერენცირებადია θ პარამეტრის მიმართ, მაშინ მისი მაქსიმუმის მოსაძებნად საჭიროა, რომ

$$\frac{\partial L}{\partial \theta} = 0. \quad (9.2)$$

თუ L ფუნქცია არადიფერენცირებადია θ -ს მიმართ, მაშინ მაქსიმუმის მოსაძებნად გამოიყენება მათემატიკის სხვა მეთოდები, რაც ხშირად დიდ გამოთვლებთანაა დაკავშირებული. ხშირად (9.2) ტოლობის ნაცვლად გამოიყენება შემდეგი ტოლობა:

$$\frac{\partial \ln L}{\partial \theta} = 0, \quad (9.3)$$

რადგან ცნობილია, რომ L და $\ln L$ ფუნქციების ექსტრემუმები ერთმანეთს ემთხვევა, ე.ი. ერთი და იგივე მნიშვნელობის წერტილებში მიიღებიან.

მაქსიმალური დასაჯერობის მეთოდის საშუალებით შევაფასოთ ნორმალურად განაწილებული X შემთხვევითი სიდიდის a და s პარამეტრები. ვთქვათ, გვაქვს ამონარჩევი x_1, x_2, \dots, x_n . როგორც ვიცით, ნორმალურად განაწილებული შემთხვევითი სიდიდის განაწილების სიმკვრივის ფუნქციას აქვს შემდეგი სახე:

$$f(x, a, s) = \frac{1}{s\sqrt{2\pi}} e^{-\frac{1}{2s^2}(x-a)^2}.$$

შესაბამისად, მაქსიმალური დასაჯერობის ფუნქციას ექნება შემდეგი სახე:

$$\begin{aligned} L &= \frac{1}{(s\sqrt{2\pi})^n} \exp\left\{-\frac{1}{2s^2} \sum_{i=1}^n (x_i - a)^2\right\} = \\ &= s^{-n} (2\pi)^{-\frac{n}{2}} \exp\left\{-\frac{1}{2s^2} \sum_{i=1}^n (x_i - a)^2\right\} \end{aligned}$$

ჩვენერთ მაქსიმალური დასაჯერობის ლოგარითმული ფუნქცია

$$\ln L = -n \ln s - \frac{n}{2} \ln(2\pi) - \frac{1}{2s^2} \sum_{i=1}^n (x_i - a)^2.$$

გავანარმოთ ეს ფუნქცია a და s პარამეტრებით, მაშინ გვექნება:

$$\begin{aligned} \frac{\partial \ln L}{\partial a} &= \frac{1}{s^2} \sum_{i=1}^n (x_i - a), \\ \frac{\partial \ln L}{\partial s} &= -\frac{n}{s} + \frac{1}{s^3} \sum_{i=1}^n (x_i - a)^2. \end{aligned}$$

ამრიგად, მივიღეთ განტოლებათა შემდეგი სისტემა:

$$\left. \begin{aligned} \frac{1}{s^2} \sum_{i=1}^n (x_i - a) &= 0 \\ -\frac{n}{s} + \frac{1}{s^3} \sum_{i=1}^n (x_i - a)^2 &= 0 \end{aligned} \right\}. \quad (9.4)$$

რადგან $\frac{1}{s^2} \neq 0$, ამიტომ (9.4) სისტემის პირველი განტოლებიდან გვექნება:

$$\sum_{i=1}^n (x_i - a) = 0, \text{ აქედან } a = \frac{1}{n} \sum_{i=1}^n x_i = \bar{x}.$$

ე.ი. a პარამეტრი ყოფილა საშუალო არითმეტიკულის შეფასება. (9.4) სისტემის მეორე განტოლება გადავწეროთ შემდეგნაირად:

$$\frac{1}{s} \left[-n + \frac{1}{s^2} \sum_{i=1}^n (x_i - a)^2 \right] = 0.$$

რადგან $\frac{1}{s} \neq 0$, ამიტომ $\left[-n + \frac{1}{s^2} \sum_{i=1}^n (x_i - a)^2 \right] = 0$, აქედან

$$s^2 = \frac{1}{n} \sum_{i=1}^n (x_i - \bar{x})^2 = \sigma^2.$$

ე.ი. მივიღეთ დისპერსიის გადაადგილებადი შეფასება. გადაუადგილებადი შეფასებისათვის, როგორც ვიცით, საჭიროა დისპერსიის ფორმულაში კორექტივის შეტანა. მაშინ მივიღებთ:

$$\sigma^2 = \frac{1}{n-1} \sum_{i=1}^n (x_i - \bar{x})^2.$$

მომენტთა მეთოდთან შედარებით, მაქსიმალური დასაჯერობის მეთოდს გააჩნია მთელი რიგი უპირატესობანი, კერძოდ:

1. ეს მეთოდი იძლევა საფუძვლიან შეფასებას;
2. თუ არსებობს ეფექტური შეფასება, მაშინ მას იძლევა მაქსიმალური დასაჯერობის მეთოდი;

3. შეფასებები ასიმპტოტურად ეფექტურია.

მეთოდის ნაკლი მდგომარეობს იმაში, რომ ზოგჯერ ამ მეთოდით მიღებული შეფასება გადაადგილებადია, რაც მოითხოვს ფორმულაში შესწორების შეტანას. აქვე უნდა აღვნიშნოთ, რომ ამონარჩევის მოცულობის ზრდისას, მაქსიმალური დასაჯერობის მეთოდი იძლევა ასიმპტოტურად გადაუადგილებად შეფასებას.

სტატისტიკური შეცდომები. ჩვენ ვნახეთ, რომ ამონარჩევის საშუალებით მიღებული სტატისტიკური ანუ შერჩევითი მახასიათებლები, როგორც წესი, თავისი აბსოლუტური მნი-

შენელობით არ ემთხვევიან შესაბამისი გენერალური ერთობლიობის მახასიათებლებს. სტატისტიკური მახასიათებლის გადახრას მის შესაბამის გენერალურ მახასიათებლებთან უწოდებენ სტატისტიკურ შეცდომას ან რეპრეზენტატიულობის შეცდომას.

სტატისტიკური შეცდომები გამოწვეულია ამონარჩევის სასრულობით. რაც უფრო დიდია ამონარჩევის მოცულობა, მით უფრო მცირეა სტატისტიკური შეცდომა.

რეპრეზენტატიულობის შეცდომის გამოსათვლელად იყენებენ ამონარჩევის საშუალებით მიღებულ დისპერსიას ან საშუალო კვადრატულ გადახრას, რომელსაც ზოგჯერ სტატისტიკურ კვადრატულ შეცდომას უწოდებენ.

თუ შემთხვევითი სიდიდის X განაწილების კანონი არც ისე ძლიერ განსხვავდება ნორმალურისგან და ამონარჩევის მოცულობა არც ისე მცირეა ($n \geq 30$), მაშინ საშუალო არითმეტიკულის სტატისტიკური შეცდომა გამოითვლება ფორმულით:

$$\varepsilon_{\bar{x}} = \frac{\sigma_x}{\sqrt{n}}$$

და საშუალო არითმეტიკული ასე ჩაინერება $\bar{x} \pm \varepsilon_{\bar{x}}$.

ხშირად საინტერესოა ვიცოდეთ, თუ რა სიზუსტით არის გამოთვლილი ესა თუ ის საშუალო არითმეტიკული (განსაკუთრებით მაშინ, როდესაც საშუალო სიდიდეები სხვადასხვა ფიზიკურ ერთეულებში არიან წარმოდგენილი), მაშინ სიზუსტე შეიძლება განისაზღვროს ფორმულით:

$$Cs = \frac{\varepsilon_{\bar{x}}}{\bar{x}} 100\% \quad \text{ან} \quad Cs = \frac{V}{\sqrt{n}},$$

სადაც, V – ვარიაციის კოეფიციენტია პროცენტებში, n – ამონარჩევის განზომილება. სიზუსტის ამ მაჩვენებელს გააჩნია თავისი ცდომილება, რომელიც განისაზღვრება შემდეგი ფორმულით:

$$\varepsilon_{Cs} = Cs \sqrt{\frac{1}{2n} + \left(\frac{Cs}{100}\right)^2}.$$

დანარჩენი სტატისტიკური მახასიათებლების შეცდომები ასე განისაზღვრება:

$$- \text{დისპერსიის: } \varepsilon_{\sigma^2} = \frac{\sigma_x^2}{\sqrt{2n}};$$

– საშუალო კვადრატული გადახრის: $\varepsilon_{\sigma} = \frac{\sigma_x}{\sqrt{2n}}$;

– ვარიაციის კოეფიციენტის:

$$\varepsilon_V = \frac{V}{\sqrt{n-1}} \sqrt{\frac{1}{2} + \left(\frac{V}{100}\right)^2} \approx \frac{V}{\sqrt{2n}};$$

– მედიანის: $\varepsilon_{Me} = \varepsilon_x \sqrt{\frac{\pi}{2}} = 1,2533 \frac{\sigma_x}{\sqrt{n}}$;

– ასიმეტრიის კოეფიციენტის:

$$\varepsilon_{A_s} = \sqrt{\frac{6(n-1)}{(n+1)(n+3)}};$$

– ექსცესის კოეფიციენტის:

$$\varepsilon_{E_x} = \sqrt{\frac{24n(n-2)(n-3)}{(n-1)^2(n+3)(n+5)}};$$

სადაც, n — ამონარჩევის განზომილება.

– ფარდობითი სიხშირის:

$$\varepsilon_p = \sqrt{\frac{p^*(1-p^*)}{n}} = \sqrt{\frac{p^*q^*}{n}}.$$

თუ ფარდობითი სიხშირე პროცენტებშია გამოსახული, მაშინ

$$\varepsilon_{p\%} = \sqrt{\frac{p^*\%(100-p^*\%)}{n}}.$$

მაგალითი. §8.3-ში მოყვანილ პირველ მაგალითში მიღებული სტატისტიკური მახასიათებლებისათვის განვსაზღვროთ სტატისტიკური შეცდომები.

საშუალო არითმეტიკულის: $\varepsilon_{\bar{x}} = \frac{\sigma_x}{\sqrt{n}} = \frac{6,29}{\sqrt{10}} = 1,99$;

საშუალო არითმეტიკულის გამოთვლის სიზუსტე:

$$Cs = \frac{\varepsilon_{\bar{x}}}{\bar{x}} 100 = \frac{1,99}{94,6} 100 = 2,10\%;$$

დისპერსიის: $\varepsilon_{\sigma_x^2} = \frac{\sigma_x^2}{\sqrt{2n}} = \frac{39,6}{\sqrt{20}} = 8,86;$

საშუალო კვადრატული გადახრის: $\varepsilon_{\sigma} = \frac{\sigma_x}{\sqrt{2n}} = \frac{6,29}{\sqrt{20}} = 1,41;$

ვარიაციის კოეფიციენტის: $\varepsilon_V = \frac{V}{\sqrt{2n}} = \frac{6,65}{\sqrt{20}} = 1,49;$

ასიმეტრიის კოეფიციენტის:

$$\varepsilon_{A_s} = \sqrt{\frac{6(n-1)}{(n+1)(n+3)}} = \sqrt{\frac{6(10-1)}{(10+1)(10+3)}} = 0,61;$$

ექსცესის კოეფიციენტის:

$$\varepsilon_{E_x} = \sqrt{\frac{24 \cdot 10(10-2)(10-3)}{(10-1)^2(10+3)(10+5)}} = 0,92.$$

ამრიგად, მივიღეთ:

საშუალო არითმეტიკული: $\bar{x} = 94,60 \pm 1,99;$

დისპერსია: $\sigma_x^2 = 39,60 \pm 8,86;$

საშუალო კვადრატული გადახრა: $\sigma_x = 6,29 \pm 1,41;$

ვარიაციის კოეფიციენტი: $V = 6,65 \pm 1,49;$

ასიმეტრიის კოეფიციენტი: $A_s = 0,19 \pm 0,61;$

ექსცესის კოეფიციენტი: $E_x = -1,23 \pm 0,92.$

9.3. პარამეტრთა შეფასების ინტერვალური მეთოდი

მცირე მოცულობის ამონარჩევების დროს პარამეტრთა შეფასების წერტილოვანი მეთოდები იძლევა მნიშვნელოვან ცდომილებებს. ამიტომ მათი გამოყენება პრაქტიკაში მიზანშეწონილი არ არის. ასეთ დროს ფართოდ იყენებენ პარამეტრთა შეფასების ინტერვალურ მეთოდს.

ვთქვათ, მოცემული ამონარჩევით მიღებული სტატისტიკური მახასიათებელი θ^* წარმოადგენს θ პარამეტრის შეფასებას. ცხადია, რომ რაც უფრო ნაკლებია $|\theta - \theta^*|$ სხვაობა, მით უფრო უკეთესია შეფასების ხარისხი, ანუ მით უფრო ზუსტია შეფასება.

სხვა სიტყვებით რომ ვთქვათ, თუ $\varepsilon > 0$ და $|\theta - \theta^*| < \varepsilon$ (9.5), მაშინ რაც უფრო ნაკლები იქნება ε , მით უფრო ზუსტი იქნება შეფასება. პრაქტიკულად, სტატისტიკური მეთოდები არ გვაძლევს იმის საშუალებას, რომ კატეგორიულად დავამტკიცოთ, რომ θ^* შეფასება აკმაყოფილებს (9.5) უტოლობას. ჩვენ შეგვიძლია მხოლოდ იმის თქმა, რომ (9.5) უტოლობა სრულდება რაღაც $\gamma = 1 - \alpha$ ალბათობით (მნიშვნელოვნების დონე α იხ. §10.1).

θ პარამეტრის შეფასების **ნდობის ალბათობა** ეწოდება იმ $\gamma = 1 - \alpha$ ალბათობას, როდესაც სრულდება (9.5) უტოლობა. ჩვეულებრივ, ნდობის ალბათობა წინასწარ არის ხოლმე მოცემული და საზოგადოდ, მას იღებენ ერთთან ახლოს, კერძოდ, 0,95; 0,99 ან 0,999 სიდიდეების ტოლად.

ვთქვათ, ალბათობა იმისა, რომ (9.5) უტოლობა სრულდება, არის γ

$$P(|\theta - \theta^*| < \varepsilon) = \gamma = 1 - \alpha.$$

თუ შევცვლით (9.5) უტოლობას მისი ტოლფასოვანი უტოლობით $-\varepsilon < \theta - \theta^* < \varepsilon$ ან $\theta^* - \varepsilon < \theta < \theta^* + \varepsilon$, მაშინ მივიღებთ:

$$P(\theta^* - \varepsilon < \theta < \theta^* + \varepsilon) = \gamma. \quad (9.6)$$

ამ გამოსახულების ინტერპრეტაცია შეიძლება შემდეგნაირად: ალბათობა იმისა, რომ $|\theta^* - \varepsilon; \theta^* + \varepsilon|$ ინტერვალში მოქცეულია უცნობი პარამეტრი θ , ტოლია γ . θ პარამეტრის **ნდობის ინტერვალი** $|\theta^* - \varepsilon; \theta^* + \varepsilon|$ ეწოდება ისეთ ინტერვალს, რომლის მიმართ წინასწარ მოცემული $\gamma = 1 - \alpha$ ნდობის ალბათობით შეიძლება იმის მტკიცება, რომ ის შეიცავს θ უცნობ პარამეტრს. როგორც (9.6) ფორმულიდან ჩანს, ნდობის ინტერვალის სიგრძე დამოკიდებულია ორ სიდიდეზე: ნდობის γ ალბათობაზე და ამონარჩევის n მოცულობაზე.

ნორმალურად განაწილებული შემთხვევითი სიდიდისათვის ნდობის ინტერვალის განსაზღვრის ზოგადი სქემა შემდეგია:

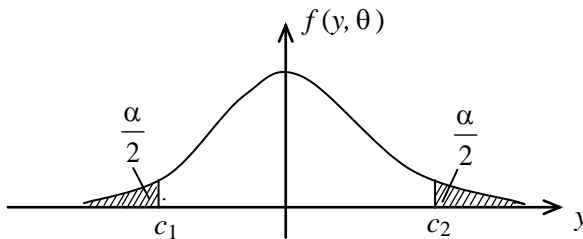
1. $F(x, \theta)$ განაწილების ფუნქციის გენერალური ერთობლიობიდან ამოღებული ამონარჩევით ახდენენ θ პარამეტრის წერტილოვან შეფასებას.

2. განიხილავენ შემთხვევით სიდიდეს, მაგალითად, $Y(\theta)$, რომელიც დამოკიდებულია θ -ზე და ცნობილია მისი განაწილების სიმკვრივის ფუნქცია $f(y, \theta)$.

3. დაუშვებენ ნდობის ალბათობის $\gamma = 1 - \alpha$ მნიშვნელობას.

4. გამოიყენებენ რა $f(y, \theta)$ განაწილების სიმკვრივის ფუნქციას, პოულობენ c_1 და c_2 რიცხვით მნიშვნელობებს ისე, რომ სრულდებოდეს შემდეგი პირობა:

$$P(c_1 < y(\theta) < c_2) = \int_{c_1}^{c_2} f(y, \theta) dy = 1 - \alpha.$$



როგორც წესი, c_1 და c_2 მნიშვნელობებს პოულობენ შემდეგი პირობებიდან:

$$P(y(\theta) < c_1) = \frac{\alpha}{2} \text{ და } P(y(\theta) > c_2) = \frac{\alpha}{2}$$

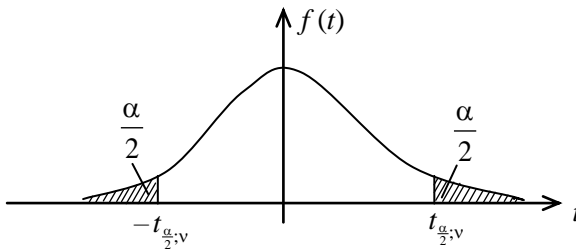
ე.ი. ნახაზზე დაშტრიხული ფართობების ჯამი უნდა უდრიდეს α -ს. ამ სქემის გამოყენებით განვიხილოთ ნორმალური განაწილების კანონის a და s პარამეტრების ნდობის ინტერვალები.

გენერალური საშუალო არითმეტიკულის ნდობის ინტერვალი. ვთქვათ, მოცემულია ნორმალურად განაწილებული შემთხვევითი სიდიდე $X \rightarrow N(a, s)$, სადაც, a და s უცნობია. გენერალური ერთობლიობიდან ამოღებული ამონარჩევის საშუალებით მოვახდინოთ \bar{x} საშუალო არითმეტიკულისა და σ საშუალო კვადრატული გადახრის წერტილოვანი შეფასებები. განვიხილოთ შემთხვევითი სიდიდე

$$t = \frac{\bar{x} - a}{\frac{\sigma}{\sqrt{n}}} = \frac{\bar{x} - a}{\sigma} \sqrt{n},$$

რომელსაც გააჩნია სტიუდენტის განაწილება $v = n-1$ თავისუფლების ხარისხით. ალბათობა იმისა, რომ შემთხვევითი სიდიდე t მოხვდება $]-t_{\frac{\alpha}{2};v}; t_{\frac{\alpha}{2};v}[$ ინტერვალში, გამოითვლება შემდეგი ფორმულით:

$$P\left(-t_{\frac{\alpha}{2};v} < \frac{\bar{x} - a}{\frac{\sigma}{\sqrt{n}}} < t_{\frac{\alpha}{2};v}\right) = 2 \int_0^{t_{\frac{\alpha}{2};v}} f(t) dt \quad (9.7)$$



თუ დავუშვებთ, რომ ეს ალბათობა $1-\alpha$ სიდიდის ტოლია, მაშინ

$$2 \int_0^{t_{\frac{\alpha}{2};v}} f(t) dt = 1 - \alpha$$

და სტიუდენტის განაწილების ცხრილიდან α მნიშვნელოვნების დონისა და v თავისუფლების ხარისხის მიხედვით შეირჩევა $t_{\frac{\alpha}{2};v}$

კრიტიკული წერტილი. ამრიგად, გვექნება:

$$P\left(-t_{\frac{\alpha}{2};v} \frac{\sigma}{\sqrt{n}} < \bar{x} - a < t_{\frac{\alpha}{2};v} \frac{\sigma}{\sqrt{n}}\right) = 1 - \alpha,$$

ან

$$P\left(\bar{x} - t_{\frac{\alpha}{2};v} \frac{\sigma}{\sqrt{n}} < a < \bar{x} + t_{\frac{\alpha}{2};v} \frac{\sigma}{\sqrt{n}}\right) = 1 - \alpha. \quad (9.8)$$

ამრიგად, ნდობის ინტერვალი $(1 - \alpha)$ ალბათობით შეიცავს გენერალური ერთობლიობის საშუალო არითმეტიკულის მნიშვნელობას.

ვნელობას. საშუალო არითმეტიკულის შეფასების სიზუსტე ტოლია:

$$\varepsilon = t_{\frac{\alpha}{2};v} \frac{\sigma}{\sqrt{n}}. \quad (9.9)$$

(9.8) ფორმულიდან გამომდინარეობს, რომ თუ ამონარჩევის n მოცულობას გავზრდით, მაშინ ნდობის ინტერვალის სიგრძე და, შესაბამისად, საშუალო არითმეტიკულის შეფასების ცდომილებაც მცირდება. თუ ნდობის ალბათობას გავზრდით, მაშინ გაიზრდება ნდობის ინტერვალის სიგრძე და, შესაბამისად, ε სიზუსტე შემცირდება. თუ ε და γ მნიშვნელობებს წინასწარ დავუშვებთ, მაშინ შეგვიძლია ვიპოვოთ ამონარჩევის ის მინიმალური მნიშვნელობა, რომელიც უზრუნველყოფს შეფასების საჭირო სიზუსტეს. ამისათვის (9.9) ფორმულიდან გვაქვს:

$$n = \frac{t_{\frac{\alpha}{2};v}^2 \sigma^2}{\varepsilon^2}.$$

მაგალითი. §8.3-ში განხილულ პირველ მაგალითში გამოთვლილი საშუალო არითმეტიკულითა $\bar{x} = 94,6$ და სტანდარტული გადახრით $\sigma = 6,29$ განვსაზღვროთ გენერალური ერთობლიობის საშუალო არითმეტიკულის a ნდობის ინტერვალი. ნდობის ალბათობა ავიღოთ $\gamma = 0,95$ ტოლი. მაშინ $\alpha = 1 - \gamma = 0,05$. სტიუდენტის განაწილების ცხრილიდან α და $v = n - 1$ სიდიდეებით შევარჩიოთ $t_{\frac{\alpha}{2};v} = 2,26$, მაშინ საშუალო არითმეტიკული შეფასების სი-

ზუსტე ტოლია:

$$\varepsilon = t_{\frac{\alpha}{2};v} \frac{\sigma}{\sqrt{n}} = \frac{2,26 \cdot 6,29}{\sqrt{10}} = 4,50.$$

გენერალური საშუალო არითმეტიკულის ნდობის ინტერვალი იქნება:

$$94,6 - 4,5 < a < 94,6 + 4,5,$$

ანუ $90,1 < a < 99,1$ ე.ი. გენერალური საშუალო არითმეტიკული $0,95$ ალბათობით მოთავსებულია $]90,1; 99,1[$ ინტერვალში. იმისათვის, რომ გავზარდოთ საშუალო არითმეტიკულის შეფასების

სიზუსტე, რომელიც არ აღემატება, მაგალითად, $\varepsilon = 2$ სიდიდეს, საჭიროა ამონარჩევის შემდეგი მინიმალური რაოდენობა:

$$n = \frac{2,26^2 \cdot 39,6}{4} \approx 50.$$

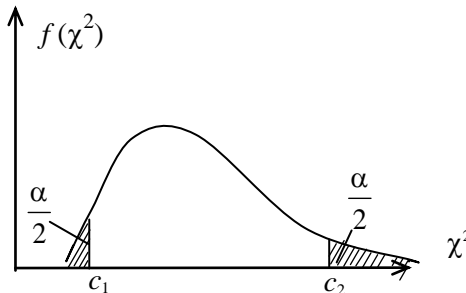
გენერალური საშუალო კვადრატული გადახრის ნდობის ინტერვალი. ვთქვათ, მოცემულია ნორმალურად განაწილებული შემთხვევითი სიდიდე $X \rightarrow N(a, s)$, რომლის a და s პარამეტრები უცნობია. n -განზომილებიანი ამონარჩევით მოვახდინოთ a მათემატიკური ლოდინისა და s საშუალო კვადრატული გადახრის წერტილოვანი შეფასებები, ე.ი.

$$a = \bar{x} \text{ და } s = \sigma.$$

საშუალო კვადრატული გადახრის ნდობის ინტერვალის ასაგებად განვიხილოთ სიდიდე

$$\chi^2 = \frac{(n-1)\sigma^2}{s^2},$$

რომელსაც გააჩნია χ^2 განაწილება $v=n-1$ თავისუფლების ხარისხით.



აღბათობა იმისა, რომ χ^2 შემთხვევითი სიდიდე მოხვდება $[c_1; c_2]$ ინტერვალში ტოლია:

$$P(c_1 < \chi^2 < c_2) = \int_{c_1}^{c_2} f(\chi^2) d(\chi^2). \quad (9.10)$$

დავუშვათ, რომ ეს აღბათობები $(1 - \alpha)$ სიდიდის ტოლია, მაშინ c_1 და c_2 წერტილები განისაზღვრებიან შემდეგი პირობების გათვალისწინებით:

$$P(\chi^2 \geq c_2) = \frac{\alpha}{2}; \quad P(\chi^2 \leq c_1) = \frac{\alpha}{2}.$$

თუ ამ განტოლებებს ამოვხსნით, მივიღებთ c_1 და c_2 კვანტილების მნიშვნელობებს

$$c_1 = \chi_{(1-\frac{\alpha}{2})v}^2; \quad c_2 = \chi_{\frac{\alpha}{2}v}^2; \quad v = n - 1.$$

(9.10) წარმოვადგინოთ შემდეგნაირად:

$$P\left(\chi_{(1-\frac{\alpha}{2})v}^2 < \frac{(n-1)\sigma^2}{s^2} < \chi_{\frac{\alpha}{2}v}^2\right) = 1 - \alpha, \quad \text{ან}$$

$$P\left(\frac{1}{\chi_{\frac{\alpha}{2}v}^2} < \frac{s^2}{(n-1)\sigma^2} < \frac{1}{\chi_{(1-\frac{\alpha}{2})v}^2}\right) = P\left(\sigma \sqrt{\frac{n-1}{\chi_{\frac{\alpha}{2}v}^2}} < s < \sigma \sqrt{\frac{n-1}{\chi_{(1-\frac{\alpha}{2})v}^2}}\right) = 1 - \alpha.$$

ამრიგად, $P(\sigma\gamma_1 < s < \sigma\gamma_2) = 1 - \alpha$,

სადაც,

$$\gamma_1 = \sqrt{\frac{n-1}{\chi_{\frac{\alpha}{2}v}^2}}, \quad \gamma_2 = \sqrt{\frac{n-1}{\chi_{(1-\frac{\alpha}{2})v}^2}}.$$

ზოგიერთ მათემატიკური სტატისტიკის სახელმძღვანელოში მოცემულია γ_1 და γ_2 კოეფიციენტების მნიშვნელობები α და v სიდიდეების გათვალისწინებით (იხ. დანართი).

მაგალითი. §8.3-ში განხილულ პირველ მაგალითში გამოთვლილი საშუალო კვადრატული გადახრით $\sigma_x = 6,29$ განვსაზღვროთ გენერალური ერთობლიობის საშუალო კვადრატული გადახრის s ნდობის ინტერვალი:

$$\sigma_x \gamma_1 < s < \sigma_x \gamma_2.$$

s ნდობის ინტერვალის ცხრილებიდან (იხ. დანართი) $\alpha = 0,05$ და $v = n - 1 = 9$ მნიშვნელობებით ვიღებთ $\gamma_1 = 0,69$ და $\gamma_2 = 1,83$ ე.ი.

$$6,29 \cdot 0,69 < s < 6,29 \cdot 1,83 \quad \text{ანუ} \quad 4,34 < s < 11,51.$$

ამრიგად, გენერალური საშუალო კვადრატული გადახრა $0,95$ ალბათობით მოთავსებულია $]4,34; 11,51[$ ინტერვალში.

10. ჰიპოთეზის სტატისტიკური შემოწმების პარამეტრული მეთოდები

10.1. სტატისტიკური ჰიპოთეზის ცნება

პარამეტრების წერტილოვანი და ინტერვალური შეფასებები წარმოადგენენ სტატისტიკური კვლევის ერთ-ერთ საწყის ეტაპს. კვლევის საბოლოო მიზანს შეიძლება წარმოადგენდეს სხვადასხვა ტექნოლოგიური პროცესების შედარებითი შეფასება, გამზომი ხელსაწყოების მახასიათებლების შედარება და ა.შ. ასეთი ტიპის ამოცანებს შედარებით ამოცანებს უწოდებენ.

მათემატიკურად, შედარების ამოცანებს უწოდებენ ჰიპოთეზების სტატისტიკურ შემოწმებას. სტატისტიკური ჰიპოთეზა ეს არის ნებისმიერი გამონათქვამი გენერალურ ერთობლიობაზე, რომლის შემოწმება შეიძლება ამონარჩევის მიხედვით.

თუ მოცემული ამონარჩევები ნორმალურადაა განაწილებული, მაშინ უნდა გამოვიყენოთ ჰიპოთეზების შემოწმების პარამეტრული მეთოდები, ხოლო როცა განაწილების კანონი უცნობია ან განსხვავდება ნორმალურისგან – არაპარამეტრული მეთოდები.

სტატისტიკური ჰიპოთეზები იყოფა ჰიპოთეზებად, რომლებიც ეხება განაწილების კანონებს და ჰიპოთეზებად, რომლებიც ეხება განაწილების პარამეტრებს. მაგალითად, თუ განაწილების კანონი უცნობია, მაშინ არის იმის საფუძველი, რომ მას გააჩნია გარკვეული სახე, მაგალითად, A . ჩამოყალიბდება ჰიპოთეზა: გენერალური ერთობლიობა განაწილებულია A კანონის სახით.

შესაძლებელია ისეთი შემთხვევა, როცა განაწილების კანონი ცნობილია, ხოლო მისი პარამეტრები – უცნობი. თუ გვაქვს იმის საფუძველი, რომ მპარამეტრი ტოლია რაიმე θ_0 ნიშვნელობისა, მაშინ დაუშვებენ ჰიპოთეზას: $\theta = \theta_0$. შესაძლებელია სხვა ჰიპოთეზებიც, მაგ. ორი და რამდენიმე პარამეტრების ტოლობის შესახებ, ამონარჩევის დამოუკიდებლობის შესახებ და სხვ.

ჰიპოთეზის ჭეშმარიტების დასადგენად, ყველაზე ზუსტი და უშეცდომო მსჯელობა შეიძლება ჩატარდეს მხოლოდ მთელი გენერალური ერთობლიობის გამოკვლევის შედეგად. მაგრამ, პრაქტიკულად, სხვადასხვა მიზეზების გამო, ასეთი კვლევის ჩატარება შეუძლებელია. ამრიგად, მსჯელობა გენერალური ერთობლიობის განაწილების ფუნქციის სახეზე ან განაწილების ფუნქციის პარამეტრების შესახებ სტატისტიკური ჰიპოთეზის ჭეშმა-

რიტებასა ან მცდარობაზე, ტარდება მხოლოდ ამონარჩევის ერთობლიობისთვის.

ამონარჩევის გამოყენებას სტატისტიკური ჰიპოთეზების შესამოწმებლად, ეწოდება დაშვებული ჰიპოთეზის ჭეშმარიტების ან მცდარობის სტატისტიკური დამტკიცება.

ზოგადად, დაშვებულ ჰიპოთეზასთან ერთად, განიხილება ალტერნატიული (კონკურირებული) ჰიპოთეზა. თუ დაშვებული ჰიპოთეზა უარყოფილი იქნება, მაშინ მის ადგილს დაიკავებს ალტერნატიული ჰიპოთეზა. ამ თვალსაზრით, სტატისტიკური ჰიპოთეზები იყოფა ნულოვან და ალტერნატიულ ჰიპოთეზებად.

ნულოვანი ეწოდება დაშვებულ (H_0) ჰიპოთეზას და აღინიშნება H_0 : სიმბოლოთი. საზოგადოდ, ნულოვანი ჰიპოთეზით მტკიცდება, რომ შესაძარებელ სიდიდეებს შორის განსხვავება არ არის, ხოლო არსებული გადახრები აიხსნება მხოლოდ და მხოლოდ ამონარჩევის შემთხვევითი ცვალებადობით.

ალტერნატიული ეწოდება H_1 : ჰიპოთეზას, რომელიც ნულოვანი ჰიპოთეზის კონკურენტია იმ თვალსაზრისით, რომ თუ ნულოვანი ჰიპოთეზა უარყოფილი იქნება, მაშინ მიიღება ალტერნატიული.

დაშვებული ჰიპოთეზა შეიძლება იყოს ჭეშმარიტი ან მცდარი, ამიტომ საჭიროა მისი შემოწმება. ამონარჩევის მონაცემებით სტატისტიკური ჰიპოთეზის შემოწმებისას, ყოველთვის არსებობს მცდარი გადაწყვეტილების მიღების რისკი. ეს აიხსნება ამონარჩევის მოცულობის სასრულობით, რის გამოც ძნელია ზუსტად დადგინდეს განანიღების ფუნქციის სახე ან მისი პარამეტრების მნიშვნელობები. აქედან გამომდინარე, სტატისტიკური ჰიპოთეზის შემოწმებისას, შესაძლებელია, დაშვებულ იქნეს ორი ტიპის შეცდომა, ანუ, როგორც მათ უწოდებენ, პირველი და მეორე გვარის შეცდომა.

პირველი გვარის შეცდომა ეწოდება ისეთ შეცდომას, როდესაც ხდება ჭეშმარიტი ნულოვანი ჰიპოთეზების უარყოფა. პირველი გვარის შეცდომის მოხდენის ალბათობას, რომელიც აღინიშნება α სიმბოლოთი, ეწოდება **მნიშვნელოვნების დონე**. უმეტეს შემთხვევაში, α მნიშვნელობას იღებენ 0,05 ან 0,01-ის ტოლად. მაგალითად, თუ $\alpha = 0,05$, ეს იმას ნიშნავს, რომ 100-დან 5 შემთხვევაში შესაძლებელია დავეუშვათ პირველი გვარის შეცდომა, ე.ი. ვუარყოთ სწორი ჰიპოთეზა.

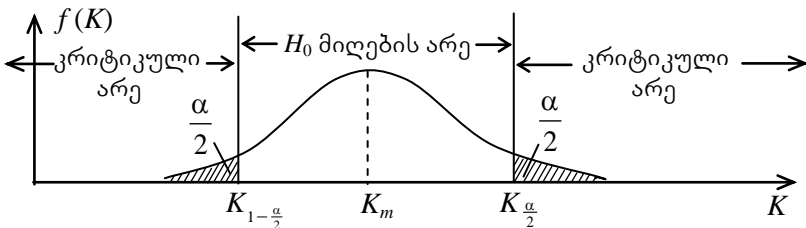
მეორე გვარის შეცდომა ეწოდება ისეთ შეცდომას, როდესაც ხდება მცდარი ნულოვანი ჰიპოთეზის მიღება. მეორე გვარის შეცდომის მოხდენის ალბათობა აღინიშნება β სიმბოლოთი.

სტატისტიკური კრიტერიუმების α მნიშვნელოვნების დონის შერჩევა დამოკიდებულია შედეგის სიმძიმეზე, რომელსაც ინვესს პირველი და მეორე გვარის შეცდომები. მაგალითად, თუ პირველი გვარის შეცდომის მოხდენა ინვესს მეტ დანაკარგებს, ვიდრე მეორე გვარის შეცდომის მოხდენა, მაშინ α -ს მნიშვნელობა უნდა მივიღოთ რაც შეიძლება ნაკლები. რალა თქმა უნდა, არ შეიძლება, რომ $\alpha = 0$, რადგან ამ შემთხვევაში მიღებული იქნება ყველა ნულოვანი ჰიპოთეზა, მათ შორის მცდარიც. უნდა გვახსოვდეს, რომ რაც უფრო მცირეა α -ს მნიშვნელობა, მით უფრო ნაკლებად ხდება ნულოვანი ჰიპოთეზის უარყოფა.

სტატისტიკური ჰიპოთეზების შემონიშნება ხდება მოცემული ამონარჩევით მიღებული მონაცემების საფუძველზე სტატისტიკური კრიტერიუმების გამოყენებით.

სტატისტიკური კრიტერიუმი (ტესტი, სტატისტიკა) ეწოდება რაიმე K შემთხვევით სიდიდეს, რომლის საშუალებითაც ხდება დაშვებული ნულოვანი ჰიპოთეზის მიღება ან უარყოფა. ნულოვანი ჰიპოთეზის შესამოწმებლად, ამონარჩევის მონაცემებით გამოითვლება კრიტერიუმში შემავალი სიდიდეების კერძო მნიშვნელობები და მიიღება თვით სტატისტიკური კრიტერიუმის კერძო მნიშვნელობა.

ვთქვათ, რაიმე ნულოვანი ჰიპოთეზის შესამოწმებლად, რომელიც ეხება განაწილების პარამეტრებს, გვაქვს K სტატისტიკა, რომლის განაწილების სიმკვრივის ფუნქცია $f(K)$ ცნობილია.

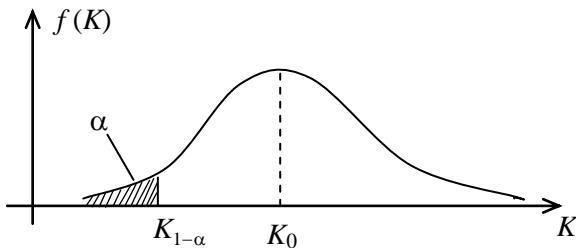


აქ, K_m არის K -ს მათემატიკური ლოდინი. $[K_{(1-\frac{\alpha}{2})}; K_{\frac{\alpha}{2}}]$ ინტერვალს ეწოდება K შემთხვევითი სიდიდის დასაშვებ მნიშვნელობათა არე,

რომლისთვისაც ნულოვანი ჰიპოთეზა მიიღება. $]-\infty; K_{(1-\frac{\alpha}{2})}[$ და $]K_{\frac{\alpha}{2}}; \infty[$ ინტერვალებს ეწოდებათ ნულოვანი ჰიპოთეზის გადახრის არეები ანუ K კრიტერიუმის კრიტიკული არეები, სადაც ხდება ნულოვანი ჰიპოთეზის უარყოფა. თუ კრიტიკული არეები მოთავსებულია მათემატიკური ლოდინის მარჯვნივ და მარცხნივ, მაშინ კრიტიკულ არეებს ეწოდებათ ორმხრივი, ხოლო K კრიტერიუმის მნიშვნელობას – ორმხრივი კრიტერიუმი. წინააღმდეგ შემთხვევაში, საქმე გვაქვს ცალმხრივ, კერძოდ მარცხენა (როცა $K < K_m$), ან მარჯვენა (როცა $K > K_m$) კრიტიკული არეებთან. კრიტიკული წერტილები (საზღვრები) $K_{(1-\frac{\alpha}{2})}; K_{\frac{\alpha}{2}}$ ეწოდებათ იმ წერტილებს, რომლებიც გამოყოფენ კრიტიკულ არეებს დასაშვებ მნიშვნელობათა არისგან.

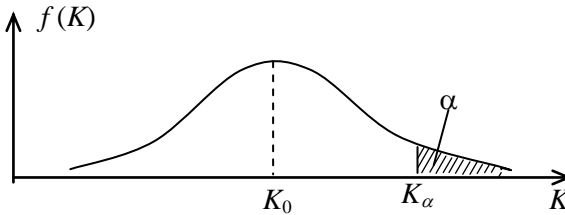
ნულოვანი ჰიპოთეზის შემონმება სტატისტიკური K კრიტერიუმით ხდება შემდეგნაირად: თუ გვაქვს ორმხრივი კრიტიკული არე, მაშინ რაიმე ნულოვანი ჰიპოთეზის $H_0: K_1 = K_2$ ალტერნატიული ჰიპოთეზა იქნება: $H_1: K_1 \neq K_2$. ამ შემთხვევაში, კრიტიკული (სასაზღვრო) წერტილი $K_{\frac{\alpha}{2};v}$ მოიძებნება $\frac{\alpha}{2}$ მნიშვნელოვნების დონისა და v თავისუფლების ხარისხის მიხედვით. თუ ჩვენს მიერ გამოთვლილი სტატისტიკა $K < K_{\frac{\alpha}{2};v}$, მაშინ არა გვაქვს საფუძველი ნულოვანი ჰიპოთეზის უარსაყოფად, წინააღმდეგ შემთხვევაში, როცა $K \geq K_{\frac{\alpha}{2};v}$, ნულოვანი ჰიპოთეზა უარყოფილი იქნება ალტერნატიულის სასარგებლოდ.

თუ გვაქვს მხოლოდ მარცხენა კრიტიკული არე, მაშინ ნულოვანი და ალტერნატიული ჰიპოთეზები ასე ჩამოყალიბდება: $H_0: K_1 = K_2$, $H_1: K_1 < K_2$.



კრიტიკული წერტილი $K_{\alpha;v}$ შეირჩევა α მნიშვნელოვნების დონისა და v თავისუფლების ხარისხის საშუალებით. თუ $K > K_{(1-\alpha);v}$, მაშინ ნულოვანი ჰიპოთეზა მიიღება, წინააღმდეგ შემთხვევაში, როცა $K \leq K_{(1-\alpha);v}$, ნულოვანი ჰიპოთეზა უარყოფილი იქნება ალტერნატიული ჰიპოთეზის სასარგებლოდ.

ანალოგიურად, როცა გვაქვს მარჯვენა კრიტიკული არე, მაშინ ნულოვანი და ალტერნატიული ჰიპოთეზები ასე ჩამოყალიბდება: $H_0: K_1 = K_2$, $H_1: K_1 > K_2$.



კრიტიკული წერტილი $K_{\alpha;v}$ შეირჩევა α მნიშვნელოვნების დონითა და v თავისუფლების ხარისხით. თუ $K \geq K_{\alpha;v}$, მაშინ ნულოვანი ჰიპოთეზა უარყოფილი იქნება H_1 ჰიპოთეზის სასარგებლოდ, წინააღმდეგ შემთხვევაში, როცა $K < K_{\alpha;v}$, არა გვაქვს საფუძველი ნულოვანი ჰიპოთეზის უარსაყოფად.

საზოგადოდ, რაც უფრო დიდია K კრიტერიუმის მნიშვნელობა, მით უფრო დიდია განსხვავება შესაძარებელ სიდიდეებს შორის.

კრიტერიუმის სიმძლავრე დამოკიდებულია მეორე გვარის შეცდომის მოხდენის β ალბათობაზე. იგი აღინიშნება M სიმბოლოთი $M = 1 - \beta$. თუ სიმძლავრე იზრდება, მაშინ შესაბამისად მცირდება β ალბათობა. თუ α სიდიდეს შევამცირებთ, მაშინ შემცირდება პირველი გვარის შეცდომა და იზრდება მეორე გვარის შეცდომა, რაც, თავის მხრივ, იწვევს კრიტერიუმის სიმძლავრის შემცირებას.

გენერალური ერთობლიობისათვის $\alpha = 0$ და $\beta = 0$. აქედან გამომდინარე, პირველი და მეორე გვარის შეცდომების მოხდენის ალბათობების ერთდროული შემცირება შესაძლებელია ამონარჩევის განზომილების გაზრდით.

ამრიგად, კრიტერიუმის სიმძლავრის გაზრდა შეიძლება ამონარჩევის მოცულობის გაზრდით. მაგრამ, თუ ამონარჩევის მოცულობა მცირეა, მაშინ α მნიშვნელობა უნდა ავიღოთ არც ისე მცირე, რადგან მცირე n და α იწვევს კრიტერიუმის სიმძლავრის შემცირებას.

ცალმხრივი და ორმხრივი კრიტიკული არეების შედარებისას უნდა აღინიშნოს, რომ ცალმხრივ კრიტერიუმს გააჩნია უფრო დიდი სიმძლავრე, ვიდრე ორმხრივს, რადგან ორმხრივი კრიტიკული არის დროს ვიღებთ $\frac{\alpha}{2}$ მნიშვნელობას, რაც იწვევს

როგორც კრიტიკული წერტილის მნიშვნელობის გაზრდას, ასევე β სიდიდის გაზრდასაც და საბოლოოდ, სიმძლავრის შემცირებას. ამიტომ, უმჯობესია ცალმხრივი კრიტიკული არის გამოყენება, თუ ნულოვანი ჰიპოთეზა ამის საშუალებას იძლევა. ასე მაგალითად, თუ გვინდა დავადგინოთ, რამდენად ეფექტურია ახალი პრეპარატით მკურნალობა ძველთან შედარებით, უნდა ავიღოთ ცალმხრივი კრიტერიუმი. მაგრამ, თუ გვინდა ერთმანეთს შევადაროთ ორი ახალი მეთოდი, მაშინ საჭიროა გამოვიყენოთ მხოლოდ ორმხრივი კრიტერიუმი, რადგან ცალმხრივი კრიტერიუმი თითქმის არ არის მგრძობიარე მეორე მეთოდთან მიმართებაში.

აქვე უნდა შევნიშნოთ, რომ არაპარამეტრულ კრიტერიუმებს, პარამეტრულ კრიტერიუმებთან შედარებით, გააჩნიათ მცირე სიმძლავრე. მაგრამ, თუ გვაქვს მცირე განზომილების ამონარჩევი, მაშინ არაპარამეტრული კრიტერიუმის გამოყენება უფრო ეფექტურია პარამეტრულთან შედარებით.

P-მნიშვნელობა. ზოგჯერ, უფრო მოსახერხებელია, ჰიპოთეზების შესამოწმებლად გამოვიყენოთ სხვა პროცედურა, რომელიც წარმოადგენს ზემოთ განხილული მეთოდის შებრუნებულ მეთოდს. კერძოდ, იმის მაგივრად, რომ ვიმუშავოდ ფიქსირებული α მნიშვნელოვნების დონით, შეგვიძლია ვიპოვოთ α -ს ზუსტი მნიშვნელობა, რომლის დროსაც უარყოფილია (მიღებულია) ნულოვანი ჰიპოთეზა. α -ს ასეთ მნიშვნელობას აღნიშნავენ P სიმბოლოთი და მას P -მნიშვნელობა ეწოდება. ამრიგად, P -მნიშვნელობა არის ის მინიმალური მნიშვნელოვნების დონე, როდესაც ხდება ნულოვანი ჰიპოთეზის უარყოფა.

თუ Z არის რაიმე K სტატისტიკის გამოთვლის შედეგად მიღებული კრიტიკული წერტილი, მაშინ $P = P(K \leq Z)$. რადგან K კრიტერიუმს გააჩნია t , F , χ^2 ან $N(0,1)$ განაწილების კანონებიდან

ერთ-ერთი მათგანი, ამიტომ P -მნიშვნელობა განისაზღვრება შემდეგნაირად:

P -მნიშ. = $1 - F(z)$ ცალმხრივი კრიტიკული არის დროს,

P -მნიშ. = $2[1 - F(z)]$ ორმხრივი კრიტიკული არის დროს.

$F(z)$ წარმოადგენს ნორმალიზებული ნორმალური განაწილების ფუნქციის მნიშვნელობას, რომელიც აღებულია შესაბამისი ცხრილიდან (იხ. დანართი) ან განისაზღვრება ლაპლასის ფუნქციის საშუალებით (3.2) ფორმულით.

თუ $P \leq \alpha$, მაშინ ნულოვანი ჰიპოთეზა უარყოფილია α მნიშვნელოვნების დონით. თუ $P > \alpha$, მაშინ ნულოვანი ჰიპოთეზის უარყოფის საფუძველი არ გაგვაჩნია.

10.2. შემთხვევითი სიდიდის საშუალოსა და დისპერსიაზე სტატისტიკური ჰიპოთეზის შემოწმება

ვთქვათ, მოცემულია X გენერალური ერთობლიობა, რომელსაც გააჩნია ნორმალური განაწილება $N(a, s)$, თანაც პარამეტრები a და s უცნობია. ამონარჩევის საშუალებით შეიძლება მივიღოთ მათი წერტილოვანი შეფასებები \bar{x} და σ_x . საჭიროა შევამოწმოთ ნულოვანი ჰიპოთეზა $H_0: \bar{x} = a_0$ ალტერნატიულის $H_1: \bar{x} \neq a_0$ საწინააღმდეგოდ. a_0 წარმოადგენს რაიმე წინასწარ ცნობილ ჰიპოთეტიურ საშუალო მნიშვნელობას.

ნულოვანი ჰიპოთეზის შესამოწმებლად საჭიროა გამოვთვალოთ სტატისტიკა:

$$t = \frac{\bar{x} - a_0}{\sigma_x} \sqrt{n-1}.$$

თუ ნულოვანი ჰიპოთეზა სამართლიანია, მაშინ t გამოსახულებას გააჩნია სტიუდენტის განაწილება $v = n-1$ თავისუფლების ხარისხით. სტიუდენტის განაწილების ცხრილიდან α მნიშვნელოვნების დონითა და v თავისუფლების ხარისხით მოიძებნება კრიტიკული წერტილი $t_{\alpha, v}$. თუ $t \geq t_{\alpha, v}$, მაშინ ნულოვანი ჰიპოთეზა უარყოფილია ალტერნატიულის სასარგებლოდ, წინააღმდეგ შემ-

თხვევაში, როცა $t < t_{\alpha;v}$, არა გვაქვს საფუძველი ნულოვანი ჰიპოთეზის უარსაყოფად.

თუ ალტერნატიულ ჰიპოთეზას აქვს $H_1: a > a_0$ სახე, მაშინ გამოიყენება t კრიტერიუმი მარჯვენა კრიტიკული არით და ამიტომ კრიტიკული წერტილი იქნება $t_{\alpha;v}$, ხოლო როცა $H_1: a < a_0$, მაშინ გვაქვს მარცხენა კრიტიკული არე და შესაბამისად, გვექნება $-t_{\alpha;v}$.

$H_0: \sigma^2 = \sigma_0^2$ ნულოვანი ჰიპოთეზის შესამოწმებლად, $H_1: \sigma^2 \neq \sigma_0^2$ ალტერნატიულის საწინააღმდეგოდ, საჭიროა გამოვთვალოთ სტატისტიკა:

$$\chi^2 = \frac{n\sigma^2}{\sigma_0^2}$$

და თუ ნულოვანი ჰიპოთეზა სამართლიანია, მაშინ მას გააჩნია χ^2 განაწილება $v = n - 1$ თავისუფლების ხარისხით. χ^2 განაწილების ცხრილიდან α და v პარამეტრების საშუალებით მოიძებნება კრიტიკული წერტილი $\chi_{\alpha;v}^2$. თუ $\chi^2 \geq \chi_{\alpha;v}^2$, მაშინ H_0 უარყოფილია H_1 -ის სასარგებლოდ, ხოლო როცა $\chi^2 < \chi_{\alpha;v}^2$, მაშინ არა გვაქვს საფუძველი ნულოვანი ჰიპოთეზის უარსაყოფად.

თუ ალტერნატიულ ჰიპოთეზას აქვს $H_1: \sigma^2 < \sigma_0^2$ სახე, მაშინ გამოიყენება χ^2 კრიტერიუმი მარცხენა კრიტიკული არით და კრიტიკული წერტილი იქნება $\chi_{(1-\alpha);v}^2$, ხოლო როცა $H_1: \sigma^2 > \sigma_0^2$, მაშინ საქმე გვაქვს χ^2 კრიტერიუმთან მარჯვენა კრიტიკული არით და სათანადოდ, გვექნება $\chi_{\alpha;v}^2$.

მაგალითი. ავტომატური ჩარხის სიზუსტე მონმდება დამუშავებული დეტალების ზომის დისპერსიის საშუალებით, რომლის სიდიდე არ უნდა აღემატებოდეს $\sigma_0^2 = 0,04$ სიდიდეს. ჩარხის საშუალებით დამუშავებულ იქნა 10 დეტალი და მიღებულ იქნა $\sigma^2 = 0,09$. საჭიროა შევამოწმოთ, აკმაყოფილებს თუ არა ჩარხი მოთხოვნილ სიზუსტეს. მნიშვნელოვნების დონე ავიღოთ $\alpha = 0,05$ ტოლად.

როგორც მოცემული პირობიდან ჩანს, უნდა შევამოწმოთ ჰიპოთეზა $H_0: \sigma^2 = 0,04$. ამ შემთხვევაში ალტერნატიული ჰიპოთეზა იქნება $H_1: \sigma^2 > 0,04$. გამოვითვალოთ სტატისტიკა:

$$\chi^2 = \frac{10 \cdot 0,09}{0,04} = 22,5.$$

χ^2 განაწილების ცხრილიდან $\chi_{0,05;9}^2 = 16,92$. რადგან $22,5 > 16,92$, ამიტომ ნულოვანი ჰიპოთეზა უარყოფილია ალტერნატიულის სასარგებლოდ, ე.ი. ჩარხი ვერ უზრუნველყოფს დეტალის დამუშავების მოთხოვნილ სიზუსტეს.

10.3. დისპერსიების ტოლობის ჰიპოთეზის შემოწმება. ფიშერის კრიტერიუმი

ვთქვათ, მოცემულია ორი X და Y ნორმალურად განაწილებული $N(a_x, s_x)$, $N(a_y, s_y)$ შემთხვევითი სიდიდე, რომელთა პარამეტრები უცნობია. განვიხილოთ ამონარჩევები x_1, x_2, \dots, x_n და y_1, y_2, \dots, y_m , რომელთა საშუალებითაც შევაფასოდ \bar{x} , \bar{y} საშუალო არითმეტიკულები და σ_x^2, σ_y^2 დისპერსიები. დისპერსიების ტოლობის $H_0: \sigma_x^2 = \sigma_y^2$ ნულოვანი ჰიპოთეზის შესამოწმებლად, უნდა გამოვიყენოთ ფიშერის კრიტერიუმი. ამისათვის განვიხილოთ სტატისტიკა:

$$F = \frac{\max(\sigma_x^2, \sigma_y^2)}{\min(\sigma_x^2, \sigma_y^2)},$$

რომელსაც გააჩნია ფიშერის განაწილება $v_1 = n - 1$ (მრიცხველისა) და $v_2 = m - 1$ (მნიშვნელის) თავისუფლების ხარისხებით.

თუ მოცემულია α მნიშვნელოვნების დონე, მაშინ F განაწილების ცხრილიდან v_1 და v_2 თავისუფლების ხარისხების საშუალებით მოიძებნება კრიტიკული წერტილი $F_{\alpha;v_1,v_2}$. თუ აღმოჩნდება, რომ $F \geq F_{\alpha;v_1,v_2}$, მაშინ ნულოვანი ჰიპოთეზა უარყოფილია

ალტერნატიულის სასარგებლოდ, ე.ი. დისპერსიები სტატისტიკურად განსხვავდება ერთმანეთისგან. თუ $F < F_{\alpha;v_1,v_2}$, მაშინ ითვლება, რომ არა გვაქვს საფუძველი ნულოვანი ჰიპოთეზის უარსაყოფად, ე.ი. დისპერსიები არ განსხვავდება ერთმანეთისგან.

მაგალითი. მოცემულია ბავშვებში სისხლის ნაკადის სიჩქარე (წმ-ში), გაზომილი ორი სხვადასხვა მეთოდით

x : 9 5 6 12 8 7 5 9 11 8 11 5 6

y : 11 4 11 9 13 8 4 12 14 9 10 7 9

უნდა შევამოწმოთ, გააჩნია თუ არა ორივე მეთოდს ერთნაირი გაზომვის სიზუსტე. მნიშვნელოვნების დონე ავიღოთ $\alpha = 0,05$ ტოლად. ჩავთვალოთ, რომ ამონარჩევები ნორმალურად არის განაწილებული.

შევამოწმოთ დისპერსიების ტოლობის ნულოვანი ჰიპოთეზა $H_0: \sigma_x^2 = \sigma_y^2$. მოცემული ამონარჩევებისთვის გვაქვს $\sigma_x^2 = 5,97$ და $\sigma_y^2 = 9,40$.

$$F = \frac{9,4}{5,97} = 1,57, \quad F_{0,05;12,12} = 2,69.$$

რადგან $1,57 < 2,69$, ამიტომ ნულოვანი ჰიპოთეზა მიიღება, ე.ი. ორივე მეთოდი ერთნაირი სიზუსტით განსაზღვრავს სისხლის ნაკადის სიჩქარეს.

იგივე ნულოვანი ჰიპოთეზა შევამოწმოთ P -მნიშვნელობით. სტანდარტიზირებული ნორმალური განაწილების ფუნქციის ცხრილიდან $F(1,57) = 0,9418$. $P = 1 - F(1,57) = 1 - 0,9418 = 0,0582 \approx 0,058$.

რადგან $0,058 > 0,05$, ამიტომ ნულოვანი ჰიპოთეზა მიიღება, ე.ი. დისპერსიები არ განსხვავდება ერთმანეთისგან.

10.4. საშუალოების ტოლობის ჰიპოთეზის შემოწმება. სტიუდენტის კრიტერიუმში

ვთქვათ, მოცემულია ორი X და Y ნორმალურად განაწილებული $N(a_x, s_x)$, $N(a_y, s_y)$ შემთხვევითი სიდიდე, სადაც, a_x , s_x , a_y , s_y პარამეტრები უცნობია. განვიხილოთ ამონარჩევები

x_1, x_2, \dots, x_n და y_1, y_2, \dots, y_m და განვსაზღვროთ \bar{x}, \bar{y} საშუალოებისა და σ_x^2, σ_y^2 დისპერსიების წერტილოვანი შეფასებები.

საშუალოების ტოლობის $H_0: \bar{x} = \bar{y}$ ნულოვანი ჰიპოთეზის შესამოწმებლად იყენებენ სტიუდენტის კრიტერიუმს, რომელსაც ზოგადად აქვს შემდეგი სახე:

$$t = \frac{\text{საშუალოების სხვაობა}}{\text{საშუალოების სხვაობის სტანდარტული შეცდომა}}$$

განვიხილოთ ორი შემთხვევა:

1. როცა $\sigma_x^2 = \sigma_y^2$. შეიძლება გვექონდეს ორი შემთხვევა:

ა) $n = m = n$. განვიხილოთ სტატისტიკა

$$t = \frac{|\bar{x} - \bar{y}|}{\sqrt{\frac{\sigma_x^2 + \sigma_y^2}{n}}}$$

რომელსაც გააჩნია სტიუდენტის განაწილება $v = 2n - 2$ თავისუფლების ხარისხით.

ბ) თუ $n \neq m$, მაშინ გვექნება:

$$t = \frac{|\bar{x} - \bar{y}|}{\sqrt{\frac{(n-1)\sigma_x^2 + (m-1)\sigma_y^2}{n+m-2} \left(\frac{n+m}{nm}\right)}}, \quad v = n + m - 2.$$

2. როცა $\sigma_x^2 \neq \sigma_y^2$. აქაც განვიხილოთ ორი შემთხვევა:

ა) თუ $n = m = n$, მაშინ:

$$t = \frac{|\bar{x} - \bar{y}|}{\sqrt{\frac{\sigma_x^2 + \sigma_y^2}{n}}}, \quad v = n - 1 + \frac{2n - 2}{\frac{\sigma_x^2}{\sigma_y^2} + \frac{\sigma_y^2}{\sigma_x^2}};$$

ბ) თუ $n \neq m$, მაშინ:

$$t = \frac{|\bar{x} - \bar{y}|}{\sqrt{\frac{\sigma_x^2}{n} + \frac{\sigma_y^2}{m}}}, \quad v = \frac{\left(\frac{\sigma_x^2}{n} + \frac{\sigma_y^2}{m}\right)^2}{\frac{\left(\frac{\sigma_x^2}{n}\right)^2}{n+1} + \frac{\left(\frac{\sigma_y^2}{m}\right)^2}{m+1}} - 2.$$

მოცემული α მნიშვნელოვნების დონითა და v თავისუფლების ხარისხით სტიუდენტის განაწილების ცხრილიდან ვპოულობთ $t_{\frac{\alpha}{2};v}$ კრიტიკულ მნიშვნელობას ორმხრივი კრიტიკული არის დროს.

თუ აღმოჩნდება, რომ $t \geq t_{\frac{\alpha}{2};v}$, მაშინ ნულოვანი ჰიპოთეზა უარყოფილია ალტერნატიულის $H_1: \bar{x} \neq \bar{y}$ სასარგებლოდ. თუ $t < t_{\frac{\alpha}{2};v}$, მაშინ ნულოვანი ჰიპოთეზა მიიღება, ე.ი. განსხვავება ორ საშუალოს შორის შემთხვევითია.

მაგალითი. წინა პარაგრაფში მოყვანილი მაგალითისთვის შევამოწმოთ საშუალოების ტოლობის ჰიპოთეზა $H_0: \bar{x} = \bar{y}$. მნიშვნელოვნების დონე ავიღოთ $\alpha = 0,05$ ტოლად.

მოცემული ამონარჩევებისთვის გვაქვს: $\bar{x} = 7,85$; $\bar{y} = 9,31$. რადგან $\sigma_x^2 = \sigma_y^2$ და ამონარჩევების განზომილება ერთნაირია, ამიტომ

$$t = \frac{|7,85 - 9,31|}{\sqrt{\frac{5,97 + 9,4}{13}}} = 1,34, \quad v = 2 \cdot 13 - 2 = 24.$$

სტიუდენტის განაწილების ცხრილიდან $t_{0,05;24} = 2,06$. რადგან $1,34 < 2,06$, ამიტომ ნულოვანი ჰიპოთეზა მიიღება, ე.ი. განსხვავება საშუალო სიდიდეებს შორის არ შეიმჩნევა.

შევამოწმოთ ნულოვანი ჰიპოთეზა P -მნიშვნელობით. $F(1,34) = 0,9099$. $P = 2(1 - 0,9099) = 0,1802$ რადგან $0,18 > 0,05$, ამიტომ ნულოვანი ჰიპოთეზა მიიღება, ე.ი. საშუალო სიდიდეები არ განსხვავდება ერთმანეთისგან.

უნდა გვახსოვდეს, რომ სტიუდენტის კრიტერიუმის გამოყენება ითვალისწინებს შემთხვევითი სიდიდეების ნორმალურ განაწილებასა და გენერალური დისპერსიების ტოლობას. თუ ეს პირობები არ სრულდება, მაშინ t -კრიტერიუმის გამოყენება მიზანშეწონილი არ არის. ამ შემთხვევაში, უფრო ეფექტურია არაპარამეტრული კრიტერიუმების გამოყენება, მაგალითად U -კრიტერიუმის.

10.5. საშუალოების მრავლობითი შედარება

თუ დისპერსიული ანალიზის გამოყენების შემდეგ აღმოჩნდება, რომ საშუალო სიდიდეები განსხვავდებიან ერთმანეთისგან, მაშინ შეუძლებელია დავადგინოთ, კერძოდ რომელი ამონარჩევები განსხვავდებიან. ასეთ შემთხვევაში, საჭიროა საშუალოები წყვილ-წყვილად შევადაროთ ერთმანეთს. უნდა გვახსოვდეს, რომ სტიუდენტის კრიტერიუმი გამოიყენება მხოლოდ ორი ამონარჩევის შედარებისთვის. თუ მას გამოვიყენებთ მრავლობითი შედარებისთვის, მაშინ ადგილი ექნება მრავლობითი შედარების ეფექტს, რომელიც იწვევს პირველი გვარის შეცდომის მოხდენის ალბათობის ზრდას, რადგან იზრდება მნიშვნელოვნების დონის სიდიდე, რომელიც განისაზღვრება შემდეგი ფორმულით:

$$\alpha' = 1 - (1 - \alpha)^k,$$

სადაც, k – შედარებათა რაოდენობაა. თუ k სიდიდე არც ისე დიდია, მაშინ შეიძლება გამოვიყენოთ მიახლოებითი ფორმულა $\alpha' \approx \alpha k$. მაგალითად, თუ $k = 3$ და ავიღებთ 5% მნიშვნელოვნების დონეს, ე.ი. $\alpha = 0,05$, მაშინ $\alpha' = 0,15$ და პირველი გვარის შეცდომის მოხდენის ალბათობა შეიძლება 15%-მდე გაიზარდოს. როცა $k = 6$, მაშინ იგი 30%-ის ტოლია და ა.შ.

ამ ეფექტის შესუსტება შეგვიძლია **ბონფერონის შესწორებით**, რომლის თანახმად, თითოეული შედარების მნიშვნელოვნების დონედ უნდა ავიღოთ $\frac{\alpha'}{k}$ სიდიდე. მაგალითად, სამჯერადი შედარებისას, მნიშვნელოვნების დონე უნდა იყოს $\frac{0,05}{3} \approx 1,7\%$. ბონფერონის შესწორება შედარებით კარგად მუშაობს მცირე რაოდენობის შედარებისას ($k < 8$). უფრო დიდი რაოდენობის შედარებისათვის, უმჯობესია, გამოვიყენოთ **ნიუმენ-კეილის კრიტერიუმი**, რომელსაც აქვს შემდეგი სახე:

$$q = \frac{|\bar{x}_j - \bar{x}_i|}{\sqrt{\frac{\bar{\sigma}^2}{2} \left(\frac{1}{n_j} + \frac{1}{n_i} \right)}}$$

სადაც, $\bar{\sigma}^2$ – შესადარებელი ამონარჩევების საშუალო დისპერსიაა, n_j, n_i – ამონარჩევების განზომილება. მიღებული q მნიშვნელობა უნდა შევადაროთ ცხრილიდან აღებულ კრიტიკულ მნიშვნელობას α მნიშვნელოვნების დონით, $v = N - m$ თავისუფლების ხარისხითა

და l შედარების ინტერვალით (იხ. დანართი). აქ, $N = \sum_{i=1}^m n_i$, m –

ამონარჩევების რაოდენობაა. თუ აღმოჩნდება, რომ $q < q_{\alpha;v,l}$, მაშინ ნულოვანი ჰიპოთეზა მიიღება, ე.ი. საშუალოები არ განსხვავდებიან ერთმანეთისგან, წინააღმდეგ შემთხვევაში, როცა $q \geq q_{\alpha;v,l}$ – საშუალოები განსხვავდებიან.

შედარების ინტერვალი l განისაზღვრება შემდეგნაირად: საშუალო სიდიდეები \bar{x}_i , $i = 1, 2, \dots, m$ უნდა დავალაგოთ ზრდადობის მიხედვით. მაგალითად, თუ ვადარებთ $\bar{x}_{(j)}$ და $\bar{x}_{(i)}$ საშუალოებს, რომლებსაც რანჟირებულ მსკრივში უკავიათ j -ური და i -ური ადგილი, მაშინ $l = j - i + 1$. მაგალითად, თუ ვადარებთ $\bar{x}_{(4)}$ და $\bar{x}_{(1)}$, მაშინ $l = 4 - 1 + 1 = 4$; $\bar{x}_{(2)}$ და $\bar{x}_{(1)}$ შედარებისას: $l = 2 - 1 + 1 = 2$ და ა.შ.

ნიუმენ-კეილის კრიტერიუმის გამოყენების შედეგი დამოკიდებულია შედარების გარკვეულ რიგზე. კერძოდ, თუ საშუალო სიდიდეებს დავალაგებთ ზრდადობით $1, 2, \dots, m$, მაშინ ჯერ უნდა შევადაროთ მწკრივის კიდურა (მაქსიმალური და მინიმალური) სიდიდეები, ე.ი. m -ური და 1 -ლი საშუალო სიდიდეები, შემდეგ m -ური და მე-2 და ა.შ. მაგალითად ოთხი საშუალო სიდიდისათვის გვექნება შედარების ასეთი თანმიმდევრობა: $4 - 1$; $4 - 2$; $4 - 3$; $3 - 1$; $3 - 2$ და $2 - 1$. აქვე უნდა შევნიშნოთ, რომ შედარება ყველა წყვილისთვის არაა საჭირო. იმ შემთხვევაში, როცა რომელიმე საშუალოების წყვილი არ განსხვავდება ერთმანეთისგან, მაშინ შედარების ამოცანა წყდება, რადგან დანარჩენები მით უფრო არ იქნება განსხვავებული. მაგალითად, თუ აღმოჩნდება, რომ $3-1$ წყვილი არ განსხვავდება ერთმანეთისგან, მაშინ აღარაა საჭირო $3 - 2$ და $2 - 1$ წყვილების შედარება.

მაგალითი. §10.4-ში მოყვანილი მაგალითისთვის, სადაც მივიღეთ, რომ საშუალოები განსხვავდებიან, ჩავატაროთ წყვილ-წყვილად შედარება. დავალაგოთ საშუალოები ზრდადობით. $\bar{x}_{(1)} = 9,1$; $\bar{x}_{(2)} = 10,1$; $\bar{x}_{(3)} = 11,5$. შევამოწმოთ $H_0 : \bar{x}_{(3)} = \bar{x}_{(1)}$ ნულოვანი ჰიპოთეზა

$$q = \frac{11,5 - 9,1}{\sqrt{\frac{3,95}{2} \left(\frac{1}{26} + \frac{1}{26} \right)}} = 6,16.$$

ამ შემთხვევაში, $l = 3 - 1 + 1 = 3$. q განაწილების ცხრილიდან $\alpha = 0,05$; $v = 3 \cdot 26 - 3 = 75$; ვლებულობთ $q_{0,05;75;3} = 3,39$. რადგან $6,16 > 3,39$, ამიტომ ნულოვანი ჰიპოთეზა უარყოფილია, ე.ი. განსხვავება სარწმუნოა. შევამოწმოთ $H_0 : \bar{x}_{(3)} = \bar{x}_{(2)}$ ნულოვანი ჰიპოთეზა

$$q = \frac{11,5 - 10,1}{\sqrt{\frac{3,95}{2} \left(\frac{1}{26} + \frac{1}{26} \right)}} = 3,59.$$

აქ, $l = 3 - 2 + 1 = 2$, $q_{0,05;75;2} = 2,82$. რადგან $3,59 > 2,82$, ამიტომ ნულოვანი ჰიპოთეზა უარყოფილია, ე.ი. განსხვავება სარწმუნოა.

$H_0 : \bar{x}_2 = \bar{x}_1$ ჰიპოთეზის შემოწმებისთვის გვექნება:

$$q = \frac{10,1 - 9,1}{\sqrt{\frac{3,95}{2} \left(\frac{1}{26} + \frac{1}{26} \right)}} = 2,57.$$

აქ, $l = 2 - 1 + 1 = 2$, $q_{0,05;75;2} = 2,82$. რადგან $2,57 < 2,82$, ამიტომ ნულოვანი ჰიპოთეზა მიიღება, ე.ი. საშუალოები არ განსხვავდება ერთმანეთისგან.

10.6. ორი დამოკიდებული ამონარჩევის შედარება

პრაქტიკაში ხშირად გვხვდება ამონარჩევები, რომლებიც არ შეიძლება განვიხილოთ როგორც დამოუკიდებელი, ისინი ერთმანეთის მიმართ დამოკიდებული არიან (მაგალითად, მონაცემები მკურნალობამდე და მკურნალობის შემდეგ). ვთქვათ, მოცემულია ორი ასეთი დამოკიდებული ამონარჩევი x_1, x_2, \dots, x_n და y_1, y_2, \dots, y_n . განვიხილოთ მათ შორის სხვაობები $d_i = x_i - y_i$, $i = 1, 2, \dots, n$, რომლებიც ნორმალურად არიან განაწილებული. განვსაზღვროთ სხვაობების საშუალო არითმეტიკული

$$\bar{d} = \frac{1}{n} \sum_{i=1}^n d_i.$$

რაც უფრო მცირეა \bar{d} სიდიდე, მით უფრო ნაკლებია განსხვავება ამ ორ ამონარჩევს შორის. აქედან გამომდინარე, საჭიროა შემოწმდეს $H_0: \bar{d} = 0$ ნულოვანი ჰიპოთეზა. ამისათვის განვიხილოთ სტატის-

ტიკა $t = \frac{|\bar{d}|}{\varepsilon_{\bar{d}}}$, სადაც, $\varepsilon_{\bar{d}}$ სხვაობათა საშუალო არითმეტიკულის შეცდომაა, რომელიც განისაზღვრება შემდეგნაირად:

$$\varepsilon_{\bar{d}} = \sqrt{\frac{1}{n(n-1)} \sum_{i=1}^n (d_i - \bar{d})^2}, \text{ ან } \varepsilon_{\bar{d}} = \sqrt{\frac{1}{n(n-1)} \left[\sum_{i=1}^n d_i^2 - \frac{1}{n} \left(\sum_{i=1}^n d_i \right)^2 \right]}.$$

t სტატისტიკას გააჩნია სტიუდენტის განაწილება $v = n - 1$ თავისუფლების ხარისხით. თუ $t < t_{\alpha;v}$, მაშინ ნულოვანი ჰიპოთეზა მიიღება, ე.ი. ამონარჩევები არ განსხვავდება ერთმანეთისგან. როცა $t \geq t_{\alpha;v}$, მაშინ ნულოვანი ჰიპოთეზა უარყოფილია და ამონარჩევები განსხვავდებიან ერთმანეთისგან.

მაგალითი. ერთი და იგივე პაციენტზე გამოცადეს ორი ძილის წამალი (x , y). შედეგები წარმოდგენილია ცხრილში დამატებითი ძილის ხანგრძლივობით საათებში. გვინტერესებს, არის თუ არა განსხვავება ამ წამლებს შორის? შევადგინოთ შემდეგი ცხრილი:

	x	y	d	d^2
1	4,0	3,0	1,0	1,00
2	3,5	3,0	0,5	0,25
3	4,1	3,8	0,3	0,09
4	5,5	2,1	3,4	11,56
5	4,6	4,9	-0,3	0,09
6	6,0	5,3	0,7	0,49
7	5,1	3,1	2,0	4,00
8	4,3	2,7	1,6	2,56
Σ			9,2	20,04

$$t = \frac{\frac{9,2}{8}}{\sqrt{\frac{20,04 - \frac{9,2^2}{8}}{8(8-1)}}} = 2,80; \alpha = 0,05; v = 7; t_{0,05;7} = 2,365.$$

რადგან $2,80 > 2,365$, ამიტომ ნულოვანი ჰიპოთეზა უარყოფილია, ე.ი. განსხვავება ძილის წამლებს შორის სარწმუნოა. შევამოწმოთ ნულოვანი ჰიპოთეზა P -მნიშვნელობით. $F(2,80) = 0,9974$. $P = 2(1 - 0,9974) = 0,0052$. რადგან $0,0052 < 0,05$, ამიტომ ნულოვანი ჰიპოთეზა უარყოფილია.

თუ გვინდა ორი დამოკიდებული ამონარჩევის დისპერსიების ტოლობის $H_0: \sigma_x^2 = \sigma_y^2$ ნულოვანი ჰიპოთეზის შემონმება, მაშინ უნდა განისაზღვროს შემდეგი სტატისტიკა:

$$t = \frac{|(Q_x - Q_y)\sqrt{n-2}|}{2\sqrt{Q_x Q_y - (Q_{xy})^2}},$$

სადაც,

$$Q_x = \sum_{i=1}^n x_i^2 - \frac{1}{n} \left(\sum_{i=1}^n x_i \right)^2, \text{ ან } Q_x = (n-1)\sigma_x^2,$$

$$Q_y = \sum_{i=1}^n y_i^2 - \frac{1}{n} \left(\sum_{i=1}^n y_i \right)^2, \text{ ან } Q_y = (n-1)\sigma_y^2,$$

$$Q_{xy} = \sum_{i=1}^n x_i y_i - \frac{1}{n} \left(\sum_{i=1}^n x_i \sum_{i=1}^n y_i \right).$$

t სიდიდეს გააჩნია სტიუდენტის განაწილება $v = n - 2$ თავისუფლების ხარისხით. α მნიშვნელოვნების დონითა და v სიდიდით სტიუდენტის განაწილების ცხრილიდან მოიძებნება $t_{\alpha;v}$ კრიტიკული მნიშვნელობა. თუ $t < t_{\alpha;v}$, მაშინ ნულოვანი ჰიპოთეზა მიიღება, ე.ი. დისპერსიები არ განსხვავდება ერთმანეთისგან. როცა $t \geq t_{\alpha;v}$, მაშინ ნულოვანი ჰიპოთეზა უარყოფილია და დისპერსიები განსხვავდება ერთმანეთისგან.

10.7. ამონარჩევი უხევი შეცდომების გამოვლენის მეთოდები

გაზომვების ჩატარების დროს დაცული უნდა იყოს გარკვეული პირობები. კერძოდ, გაზომვები უნდა ჩატარდეს ერთნაირ პირობებში, ერთნაირი კლასის ხელსაწყოებით. მიუხედავად ამისა, გაზომვის შედეგად მიღებულ ერთობლიობაში ზოგჯერ გვხვდება უხევი შეცდომები, ამიტომ ის მონაცემი, რომელიც მკვეთრად განსხვავდება სხვებისგან, უნდა იქნეს შეფასებული და თუ ის აღმოჩნდება უხევი შეცდომა (არტეფაქტი), მაშინ იგი უნდა გამოირიცხოს

ერთობლიობიდან. უხეში შეცდომის თავიდან აცილების ერთ-ერთ უმარტივეს მეთოდს წარმოადგენს ე.წ. „სამი სიგმას“ წესი, რომლის არსი შემდეგში მდგომარეობს: ალბათობის თეორიიდან ცნობილია, რომ თუ X შემთხვევით სიდიდეს გააჩნია $f(x)$ განაწილების სიმკვრივე, მაშინ რაიმე $[\alpha; \beta]$ ინტერვალში მისი მოხვედრის ალბათობა გამოითვლება ფორმულით:

$$P(\alpha < X < \beta) = \int_{\alpha}^{\beta} f(x) dx.$$

თუ შემთხვევითი სიდიდე ნორმალურადაა განაწილებული, მაშინ

$$P(\alpha < X < \beta) = \frac{1}{\sigma\sqrt{2\pi}} \int_{\alpha}^{\beta} e^{-\frac{(x-a)^2}{2\sigma^2}} dx.$$

შემოვიღოთ აღნიშვნა $z = \frac{x-a}{\sigma}$. აქედან,

$$x-a = z\sigma, \quad x = a + z\sigma, \quad dx = \sigma dz.$$

ვიპოვოთ ინტეგრირების ახალი ზღვრები. თუ $x = \alpha$, მაშინ

$$z = \frac{\alpha-a}{\sigma}. \text{ თუ } x = \beta, \text{ მაშინ } z = \frac{\beta-a}{\sigma}. \text{ ე.ი. გვექნება:}$$

$$\begin{aligned} P(\alpha < X < \beta) &= \frac{1}{\sigma\sqrt{2\pi}} \int_{\frac{\alpha-a}{\sigma}}^{\frac{\beta-a}{\sigma}} e^{-\frac{z^2}{2}} \sigma dz = \frac{1}{\sqrt{2\pi}} \int_{\frac{\alpha-a}{\sigma}}^{\frac{\beta-a}{\sigma}} e^{-\frac{z^2}{2}} dz + \frac{1}{\sqrt{2\pi}} \int_0^{\frac{\beta-a}{\sigma}} e^{-\frac{z^2}{2}} dz = \\ &= \frac{1}{\sqrt{2\pi}} \int_0^{\frac{\beta-a}{\sigma}} e^{-\frac{z^2}{2}} dz - \frac{1}{\sqrt{2\pi}} \int_0^{\frac{\alpha-a}{\sigma}} e^{-\frac{z^2}{2}} dz. \end{aligned}$$

თუ გამოვიყენებთ ლაპლასის ფუნქციას $\Phi(x) = \frac{2}{\sqrt{2\pi}} \int_0^x e^{-\frac{t^2}{2}} dt$, მაშინ

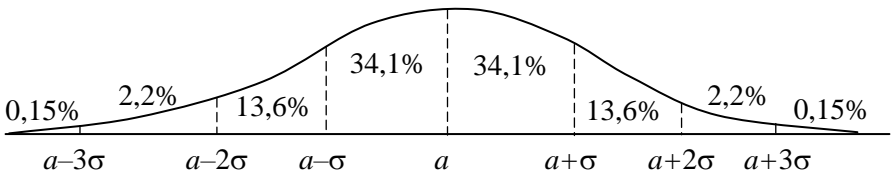
მივიღებთ:

$$P(\alpha < X < \beta) = \frac{1}{2} \left[\Phi\left(\frac{\beta-a}{\sigma}\right) - \Phi\left(\frac{\alpha-a}{\sigma}\right) \right].$$

ლაპლასის ფუნქციის ანუ ალბათობის ინტეგრალის მნიშვნელობები მოცემულია სპეციალურ ცხრილებში (იხ. დანართი). x -ის

უარყოფითი მნიშვნელობებისთვის $\Phi(-x)$ ფუნქციის მნიშვნელობათა მოსაძებნად უნდა ვისარგებლოთ ამ ფუნქციის კენტობით: $\Phi(-x) = -\Phi(x)$.

ნორმალურად განაწილებული მრუდის მათემატიკური ლოდინიდან როგორც მარცხნივ, ისე მარჯვნივ, გადავზომოთ სტანდარტული გადახრა σ , 2σ და 3σ . რადგან ნორმალური განაწილება ალბათური განაწილებაა, ამიტომ ფართობი $f(x)$ მრუდის ქვეშ ერთის ტოლია და σ -ს საშუალებით მიღებული მონაკვეთების ფართობები პროცენტებში სათანადოდ ნაჩვენებია შემდეგ ნახაზზე:



გამოვთვალოთ $]a - 3\sigma; a + 3\sigma[$ ინტერვალში X შემთხვევითი სიდიდის მოხვედრის ალბათობა

$$P[a - 3\sigma < X < a + 3\sigma] = \frac{1}{2} \left[\Phi\left(\frac{a + 3\sigma - a}{\sigma}\right) - \Phi\left(\frac{a - 3\sigma - a}{\sigma}\right) \right] =$$

$$= \frac{1}{2} [\Phi(3) - \Phi(-3)] = \Phi(3) = 0,9973$$

ამრიგად, **0,9973** ალბათობით შეგვიძლია ჩავთვალოთ, რომ ნორმალურად განაწილებული შემთხვევითი სიდიდე მიიღებს მნიშვნელობებს $[a - 3\sigma; a + 3\sigma]$ შუალედიდან, ხოლო ალბათობა იმისა, რომ შემთხვევითი სიდიდე მოხვდეს ამ შუალედის გარეთ, ძალიან მცირეა და იგი **0,0027**-ის ტოლია, ე.ი. ეს ხდომილობა შეგვიძლია ჩავთვალოთ პრაქტიკულად შეუძლებელ ხდომილობად. ამ ფაქტს ემყარება ე.წ. **სამი სიგმას წესი**, რომელიც შეიძლება ასე ჩამოვაყალიბოთ: თუ შემთხვევითი სიდიდე ნორმალურადაა განაწილებული, მაშინ მისი გადახრა მათემატიკური ლოდინიდან (საშუალო არითმეტიკულიდან) აბსოლუტური სიდიდით პრაქტიკულად არ აღემატება გასამკვეცებულ საშუალო კვადრატულ გადახრას.

აქედან გამომდინარე, თუ ამონარჩევში გვხვდება ისეთი, ვთქვათ, x_j სიდიდე, რომელიც აკმაყოფილებს $|\bar{x} \pm x_j| > 3\sigma$ უტოლობას, მაშინ ეს სიდიდე შეიძლება ჩაითვალოს არტეფაქტად და გამოირიცხოს ამონარჩევიდან.

არსებობენ სხვა, უფრო ზუსტი ტესტები უხეში გაზომვების გამოსავლენად. განვიხილოთ ერთ-ერთი მათგანი, რომელსაც **ტომპსონის წესს** უწოდებენ. ამ შემთხვევაში, ნულოვანი ჰიპოთეზა ჩამოყალიბდება შემდეგნაირად: „ამონარჩევში არ არის უხეში გაზომვები“. ამ ჰიპოთეზის შესამოწმებლად, საჭიროა გამოითვალოს შემდეგი სტატისტიკა:

$$t_i = \frac{x_i - \bar{x}}{\sigma}, \quad i = 1, 2, \dots, n,$$

სადაც, \bar{x} – საშუალო არითმეტიკულის, ხოლო σ – საშუალო კვადრატული გადახრის შეფასებებია.

ტომპსონის წესით, ამონარჩევის ყველა ის x_j მნიშვნელობა, რომლისთვისაც სრულდება უტოლობა $|t_i| \geq z_{\alpha;v}$, უნდა ჩაითვალოს არტეფაქტად და გამოირიცხოს ამონარჩევიდან. $z_{\alpha;v}$ კრიტიკული მნიშვნელობა α მნიშვნელოვნების დონითა და $v = n - 2$ თავისუფლების ხარისხით მოიძებნება სპეციალურ ცხრილში (იხ. დანართი).

მაგალითი. მოცემულ ამონარჩევში x : 23 40 9 25 38 32 37 26 28 $x_3 = 9$ მნიშვნელობა ძლიერ განსხვავდება სხვა მნიშვნელობებისგან. შევამოწმოთ ნულოვანი ჰიპოთეზა H_0 : „ x_3 მნიშვნელობა ეკუთვნის მოცემულ ამონარჩევს“. ამისათვის გამოვთვალოთ:

$$\bar{x} = 28,67 \text{ და } \sigma_x = 9,59. \quad |t| = \frac{9 - 28,67}{9,59} = 2,07; \quad z_{0,05;7} = 1,896.$$

რადგან $2,07 > 1,896$, ამიტომ H_0 ჰიპოთეზა უარყოფილია, ე.ი. $x_3 = 9$ წარმოადგენს არტეფაქტს და იგი უნდა გამოირიცხოს ამონარჩევიდან.

როცა ამონარჩევის მოცულობა დიდია და $n \rightarrow \infty$, მაშინ შეიძლება ვისარგებლოთ შემდეგი აპროქსიმაციით:

$$z_{\alpha;v} \approx \lambda_q \left(1 + \frac{3 - \lambda_q^2}{4v} + \frac{3 - 32\lambda_q^2 + 5\lambda_q^4}{96v^2} \right),$$

სადაც, λ_q არის ნორმალური განაწილების კვანტილი. ამასთან, $q = 1 - \frac{\alpha}{2}$. მაგალითად, როცა $v = 100$, $\alpha = 0,05$ ე.ი. $q = 0,975$, მაშინ $z_{\alpha;v} = 1,956$, რაც ზუსტად ემთხვევა ცხრილის მნიშვნელობას.

10.8. ორზე მეტი ამონარჩევის ერთდროული შედარება

დისპერსიების ტოლობის ჰიპოთეზა. თუ მოცემულია k რაოდენობის ნორმალურად განაწილებული ამონარჩევი $x_{ij}, i = 1, 2, \dots, n_j, j = 1, 2, \dots, k$ და საჭიროა $H_0 : \sigma_1^2 = \sigma_2^2 = \dots = \sigma_k^2$ ნულოვანი ჰიპოთეზის შემოწმება, მაშინ შეგვიძლია გამოვიყენოთ **ბარტლეტის კრიტერიუმი**:

$$\chi^2 = \frac{2,303}{C} \left[(N - k) \lg \bar{\sigma}^2 - \sum_{i=1}^k (n_i - 1) \lg \sigma_i^2 \right],$$

სადაც,

$$C = \frac{1}{3(k-1)} \left[\sum_{i=1}^k \frac{1}{n_i - 1} - \frac{1}{N - k} \right] + 1,$$

$$N = \sum_{i=1}^k n_i; \quad \sigma_j^2 = \frac{1}{n_j - 1} \sum_{i=1}^{n_j} (x_{ij} - \bar{x}_j)^2$$

$$\bar{x}_j = \frac{1}{n_j} \sum_{i=1}^{n_j} x_{ij}, \quad j = 1, 2, \dots, k.$$

$\bar{\sigma}^2$ — გაერთიანებული ანუ საშუალო დისპერსიაა, რომელიც ასე განისაზღვრება:

$$\bar{\sigma}^2 = \frac{1}{N - k} \sum_{j=1}^k (n_j - 1) \sigma_j^2 = \frac{1}{N - k} \sum_{j=1}^k \sum_{i=1}^{n_j} (x_{ij} - \bar{x}_j)^2.$$

თუ ამონარჩევები ერთნაირი განზომილებისაა, ე.ი. $n_1 = n_2 = \dots = n_k = n_0$, მაშინ χ^2 სტატისტიკის გამოსახულება გამარტივდება და მიიღებს შემდეგ სახეს:

$$\chi^2 = \frac{2,303}{C} \left[k(n_0 - 1) \left\{ \lg \bar{\sigma}^2 - \frac{1}{k} \sum_{i=1}^k \lg \sigma_i^2 \right\} \right],$$

სადაც,

$$C = \frac{k+1}{3k(n_0-1)} + 1; \quad \bar{\sigma}^2 = \frac{1}{k} \sum_{i=1}^k \sigma_i^2.$$

χ^2 სტატისტიკას გააჩნია χ^2 განაწილება $v = k - 1$ თავისუფლების ხარისხით. α მნიშვნელოვნების დონითა და v სიდიდით χ^2 განაწილების ცხრილიდან მოიძებნება $\chi_{\alpha;v}^2$ კრიტიკული წერტილი. თუ $\chi^2 \geq \chi_{\alpha;v}^2$, მაშინ ნულოვანი ჰიპოთეზა უარყოფილია ალტერნატიულის სასარგებლოდ. როცა $\chi^2 < \chi_{\alpha;v}^2$, მაშინ ნულოვანი ჰიპოთეზა მიიღება. ე.ი. დისპერსიები არ განსხვავდება ერთმანეთისგან.

მაგალითი. მოცემულია $n_1 = 9$, $n_2 = 6$ და $n_3 = 5$ განზომილებიანი ამონარჩევები $\sigma_1^2 = 8,00$, $\sigma_2^2 = 4,67$, $\sigma_3^2 = 4,00$ დისპერსიებით. შევამოწმოთ დისპერსიების ტოლობის ნულოვანი ჰიპოთეზა. ავიღოთ $\alpha = 0,05$.

$$\bar{\sigma}^2 = \frac{1}{20-3} (8 \cdot 8 + 5 \cdot 4,67 + 4 \cdot 4) = 6,079; \quad \lg \bar{\sigma}^2 = 0,7838;$$

$$C = 1 + \frac{1}{6} \left[\left(\frac{1}{8} + \frac{1}{5} + \frac{1}{4} \right) - \frac{1}{17} \right] = 1,086;$$

$$\chi^2 = \frac{2,303}{1,086} [17 \cdot 0,7838 - (8 \cdot 0,9031 + 5 \cdot 0,6693 + 4 \cdot 0,6021)] = 0,731;$$

$$\chi_{0,05;2}^2 = 5,99.$$

რადგან $0,731 < 5,99$, ამიტომ ნულოვანი ჰიპოთეზა მიიღება, ე.ი. დისპერსიები არ განსხვავდება ერთმანეთისაგან.

საშუალოების ტოლობის ჰიპოთეზა. თუ ამონარჩევების რაოდენობა $k > 2$, მაშინ საშუალოების ტოლობის ჰიპოთეზა ეფუძ-

ნება დისპერსიულ ანალიზს, კერძოდ, ერთფაქტორიანი დისპერსიული ანალიზის შედეგებს. იდეა მდგომარეობს საერთო დისპერსიის ორ დამოუკიდებელ მდგენელად: ფაქტორულ (ჯგუფთაშორისო) და ნარჩენ (შიგაჯგუფური) დისპერსიებად წარმოდგენაში. ამრიგად, $\sigma^2 = \sigma_{\text{ფ}}^2 + \sigma_{\text{ნარ}}^2$. ფიშერის კრიტერიუმის

$F = \frac{\sigma_{\text{ფ}}^2}{\sigma_{\text{ნარ}}^2}$ გამოყენებით მიდიან დასკვნამდე, არის თუ არა საშუალოებს შორის განსხვავება.

ვთქვათ, მოცემულია k რაოდენობის n_i განზომილებიანი ამონარჩევები x_{ij} , $i=1,2,\dots,n_j$, $j=1,2,\dots,k$, რომლებიც ნორმალურად არიან განაწილებული $N(a_i, s_i)$, სადაც, a_i და s_i პარამეტრები უცნობია, მაგრამ გულისხმობენ, რომ $s_1^2 = s_2^2 = \dots = s_k^2$. ამ ტოლობის ჰიპოთეზა შეიძლება შევამოწმოთ ბარტლეტის კრიტერიუმით.

საშუალოების ტოლობის $H_0: \bar{x}_1 = \bar{x}_2 = \dots = \bar{x}_k$ ნულოვანი ჰიპოთეზის შესამოწმებლად განვიხილოთ სტატისტიკა:

$$F = \frac{\frac{1}{k-1} \sum_{i=1}^k n_i (\bar{x}_i - \bar{\bar{x}})^2}{\frac{1}{N-k} \sum_{j=1}^k \sum_{i=1}^{n_j} (x_{ij} - \bar{x}_j)^2},$$

ან

$$F = \frac{\frac{1}{k-1} \sum_{i=1}^k n_i (\bar{x}_i - \bar{\bar{x}})^2}{\frac{1}{N-k} \sum_{j=1}^k (n_j - 1) \sigma_j^2},$$

სადაც,

$$\bar{x}_j = \frac{1}{n_j} \sum_{i=1}^{n_j} x_{ij}, \quad j=1,2,\dots,k$$

$$\bar{\bar{x}} = \frac{1}{N} \sum_{i=1}^k n_i \bar{x}_i; \quad N = \sum_{i=1}^k n_i.$$

თუ ამონარჩევების განზომილებები ერთმანეთის ტოლია, ე.ი. $n_1 = n_2 = \dots = n_k = n_0$, მაშინ გვექნება:

$$F = \frac{n_0 \sum_{i=1}^k (\bar{x}_i - \bar{\bar{x}})^2}{\bar{\sigma}^2},$$

სადაც,

$$\bar{\bar{x}} = \frac{1}{k} \sum_{i=1}^k \bar{x}_i, \quad \bar{\sigma}^2 = \frac{1}{k} \sum_{i=1}^k \sigma_i^2.$$

F სიდიდეს გააჩნია ფიშერის განაწილება $\nu_1 = k - 1$ და $\nu_2 = N - k$ თავისუფლების ხარისხებით. α, ν_1 და ν_2 სიდიდეებით ფიშერის განაწილების ცხრილიდან შეირჩევა $F_{\alpha; \nu_1; \nu_2}$ კრიტიკული ნერტილი. თუ $F \geq F_{\alpha; \nu_1; \nu_2}$, მაშინ ნულოვანი ჰიპოთეზა უარყოფილია, ხოლო როცა $F < F_{\alpha; \nu_1; \nu_2}$, მაშინ ნულოვანი ჰიპოთეზა მიიღება.

მაგალითი. მოცემულია $n_0 = 26$ განზომილებიანი სამი ამონარჩევი $\sigma_1^2 = 1,69$; $\sigma_2^2 = 4,41$ და $\sigma_3^2 = 5,76$ დისპერსიებისა და $\bar{x}_1 = 11,5$; $\bar{x}_2 = 10,1$ და $\bar{x}_3 = 9,1$ საშუალო არითმეტიკულების შეფასებებით.

უნდა შევამოწმოთ $H_0 : \bar{x}_1 = \bar{x}_2 = \bar{x}_3$ ნულოვანი ჰიპოთეზა. ავიღოთ $\alpha = 0,05$.

$$\bar{\bar{x}} = \frac{1}{3}(11,5 + 10,1 + 9,1) = 10,2;$$

$$\bar{\sigma}^2 = \frac{1}{3}(1,69 + 4,41 + 5,76) = 3,95;$$

$$F = \frac{\frac{16}{3-1} [(11,5 - 10,2)^2 + (10,1 - 10,2)^2 + (9,1 - 10,2)^2]}{3,95} = 9,58;$$

$$\nu_1 = 3 - 1 = 2; \quad \nu_2 = 3 \cdot 26 - 3 = 75; \quad F_{0,05; 2; 75} = 3,15.$$

რადგან $9,58 > 3,15$, ამიტომ ნულოვანი ჰიპოთეზა უარყოფილია, ე.ი. საშუალო სიდიდეები განსხვავდება ერთმანეთისგან.

შევამოწმოთ ნულოვანი ჰიპოთეზა P -მნიშვნელობით. $F(3,15) = 0,9992$ $P = 1 - 0,9992 = 0,0008$ რადგან $0,0008 < 0,05$, ამიტომ ნულოვანი ჰიპოთეზა უარყოფილია.

უნდა აღვნიშნოთ, რომ თუ შესადარებელი ამონარჩევების რაოდენობა $k = 2$, მაშინ ადვილად მტკიცდება, რომ სტიუდენტის კრიტერიუმი წარმოადგენს დისპერსიული ანალიზის კერძო შემთხვევას და სამართლიანია $F = t^2$ ტოლობა. მართლაც, თუ განვიხილავთ ერთნაირი განზომილების ორ ამონარჩევს \bar{x}_1, \bar{x}_2 საშუალოებითა და σ_1^2, σ_2^2 დისპერსიებით, მაშინ

$$F = \frac{n \left((\bar{x}_1 - \bar{\bar{x}})^2 + (\bar{x}_2 - \bar{\bar{x}})^2 \right)}{\frac{1}{2} (\sigma_1^2 + \sigma_2^2)},$$

სადაც, $\bar{\bar{x}} = \frac{1}{2} (\bar{x}_1 + \bar{x}_2)$. F -ის გამოსახულებიდან გამოვრიცხოთ $\bar{\bar{x}}$.

$$\begin{aligned} (\bar{x}_1 - \bar{\bar{x}})^2 + (\bar{x}_2 - \bar{\bar{x}})^2 &= \left[\bar{x}_1 - \frac{1}{2} (\bar{x}_1 + \bar{x}_2) \right]^2 + \left[\bar{x}_2 - \frac{1}{2} (\bar{x}_1 + \bar{x}_2) \right]^2 = \\ &= \left(\frac{1}{2} \bar{x}_1 - \frac{1}{2} \bar{x}_2 \right)^2 + \left(\frac{1}{2} \bar{x}_2 - \frac{1}{2} \bar{x}_1 \right)^2 = \frac{1}{2} (\bar{x}_1 - \bar{x}_2)^2, \end{aligned}$$

$$\text{რადგან } \left(\frac{1}{2} \bar{x}_2 - \frac{1}{2} \bar{x}_1 \right)^2 = \left(\frac{1}{2} \bar{x}_1 - \frac{1}{2} \bar{x}_2 \right)^2.$$

$$F = \frac{\frac{n}{2} (\bar{x}_1 - \bar{x}_2)^2}{\frac{1}{2} (\sigma_1^2 + \sigma_2^2)} = \frac{(\bar{x}_1 - \bar{x}_2)^2}{\frac{\sigma_1^2}{n} + \frac{\sigma_2^2}{n}} = \left[\frac{\bar{x}_1 - \bar{x}_2}{\sqrt{\frac{\sigma_1^2}{n} + \frac{\sigma_2^2}{n}}} \right]^2 = t^2.$$

ამრიგად, ორი ამონარჩევის შედარების დროს, სტიუდენტის კრიტერიუმი და დისპერსიული ანალიზი წარმოადგენს ერთი და იგივე კრიტერიუმს. თუ ამონარჩევების რაოდენობა $k > 2$, მაშინ ეს ასე არ არის.

10.9. მრავალგანზომილებიანი სისტემის ზოგიერთი
ჰიპოთეზების შემოწმება

ორი ვექტორის ტოლობის ჰიპოთეზა. ვთქვათ, მოცემულია საშუალოების ორი ვექტორი \bar{X}_i და \bar{X}_j , რომლებიც მიღებულია ორი ნორმალურად განაწილებული n -განზომილებიანი სისტემიდან. საჭიროა შევამოწმოთ \bar{X}_i და \bar{X}_j ვექტორების ტოლობის ნულოვანი ჰიპოთეზა, ე.ი. $H_0 : \bar{X}_i = \bar{X}_j$. ასეთი ნულოვანი ჰიპოთეზის შესამოწმებლად განვიხილოთ ჰოტელინგის კრიტერიუმი:

$$T^2 = \frac{m_i m_j}{m_i + m_j} (\bar{X}_i - \bar{X}_j)' S^{-1} (\bar{X}_i - \bar{X}_j),$$

სადაც, m_i, m_j შესაბამისად X_i და X_j ვექტორების განზომილებაა; S – გაერთიანებული კოვარიაციის მატრიცა, რომელიც გამოითვლება შემდეგნაირად:

$$S = \frac{1}{m_i + m_j - 2} [(m_i - 1)S_i + (m_j - 1)S_j].$$

თუ ნულოვანი ჰიპოთეზა სამართლიანია, მაშინ

$$F = \frac{m_i + m_j - n - 1}{(m_i + m_j - 2)n} T^2$$

სიდიდეს გააჩნია ფიშერის განაწილება $v_1 = n$ და $v_2 = m_i + m_j - n - 1$ თავისუფლების ხარისხებით. თუ $F < F_{\alpha; v_1, v_2}$, მაშინ ნულოვანი ჰიპოთეზა მიიღება, წინააღმდეგ შემთხვევაში, როცა $F \geq F_{\alpha; v_1, v_2}$, ნულოვანი ჰიპოთეზა უარყოფილია, ე.ი. განსხვავება ორ ვექტორს შორია სარწმუნოა.

კოვარიაციული მატრიცების ტოლობის ჰიპოთეზა. განვიხილოთ ორი კოვარიაციული მატრიცების ტოლობის ჰიპოთეზა $H_0 : S_1 = S_2$. ამისათვის უნდა გამოვთვალოთ შემდეგი სტატისტიკა:

$$W = b(-2 \ln V_1),$$

სადაც,

$$b = 1 - \left(\sum_{j=1}^2 \frac{1}{v_j} - \frac{1}{\sum_{j=1}^2 v_j} \right) \left(\frac{2n^2 + 3n - 1}{6(n+1)} \right),$$

$$-2 \ln V_1 = \left(\sum_{j=1}^2 v_j \right) \ln |S| - \sum_{j=1}^2 (v_j \ln |S_j|),$$

$$v_1 = m_1 - 1; \quad v_2 = m_2 - 1,$$

რომელსაც გააჩნია χ^2 განაწილება $v = \frac{n(n+1)}{2}$ თავისუფლების ხარისხით.

თუ $W < \chi_{\alpha;v}^2$, მაშინ ნულოვანი ჰიპოთეზა მიიღება და კოვარიაციული მატრიცები ერთმანეთის ტოლია, წინააღმდეგ შემთხვევაში, როცა $W \geq \chi_{\alpha;v}^2$, მაშინ კოვარიაციული მატრიცები განსხვავდება ერთმანეთისგან.

არტეფაქტური ვექტორის გამოვლენა. ვთქვათ, მოცემულია ნორმალურად განაწილებული X_1, X_2, \dots, X_n შემთხვევითი ვექტორთა სისტემა, რომლის სტრიქონები დაკვირვებათა ვექტორებია: $X_i = (x_{i1}, x_{i2}, \dots, x_{in})$, $i = 1, 2, \dots, m$.

არტეფაქტური ვექტორის გამოვლენის პროცედურა შემდეგია: თოთოეული X_i , $i = 1, 2, \dots, m$ დაკვირვების ვექტორისთვის გამოითვლება საშუალო არითმეტიკულის ვექტორი

$$\bar{X}_i = \frac{1}{n} \sum_{k=1}^n x_{ik}, \quad i = 1, 2, \dots, m$$

და კოვარიაციული მატრიცა S ყველა $(m - 1)$ დაკვირვების ვექტორით, გარდა X_i ვექტორისა. შემდეგ გამოითვლება მახალანობისის მანძილი X_i და \bar{X}_i ვექტორებს შორის S კოვარიაციული მატრიცის საშუალებით:

$$D_i^2 = (X - \bar{X}_i)' S^{-1} (X - \bar{X}_i).$$

ამის შემდეგ განისაზღვრება F_i მნიშვნელობა $k = m - 1$ სიდიდისთვის.

$$F_i = \frac{(m-n)m}{(m^2-1)n} D_i^2,$$

რომელსაც გააჩნია ფიშერის განაწილება $v_1 = n$ და $v_2 = m - n$ თავისუფლების ხარისხებით.

თუ $F_i > F_{\alpha; v_1, v_2}$, მაშინ X_i ვექტორი ითვლება არტეფაქტად და იგი უნდა გამოირიცხოს ამონარჩევიდან. პროცედურა გრძელდება დარჩენილ $(n - 1)$ დაკვირვებისთვის.

11. ჰიპოთეზის სტატისტიკური შემოწმების არაპარამეტრული მეთოდები

11.1. განაწილების კანონის შესახებ ჰიპოთეზის შემოწმება

ჩვენ განვიხილეთ ჰიპოთეზების შემოწმების პარამეტრული მეთოდები, რომლებიც ეხებოდა შემთხვევითი სიდიდის განაწილების პარამეტრებს, თანაც ითვლებოდა, რომ განაწილების კანონი ცნობილი იყო. მაგრამ ბევრ პრაქტიკულ ამოცანებში შესასწავლი შემთხვევითი სიდიდის განაწილების კანონი უცნობია, ე.ი. წარმოადგენს ჰიპოთეზას, რომელიც საჭიროებს სტატისტიკურ შემოწმებას.

ვთქვათ, საჭიროა შემოწმდეს ჰიპოთეზა იმის შესახებ, რომ X შემთხვევითი სიდიდე ემორჩილება $F(x)$ განაწილების კანონს. ამ ჰიპოთეზის შესამოწმებლად ავიღოთ ამონარჩევი x_1, x_2, \dots, x_n , რომელიც მიიღება X შემთხვევით სიდიდეზე n დამოუკიდებელი დაკვირვებით. ამონარჩევით შეიძლება აიგოს ემპირიული განაწილების ფუნქცია $F^*(x)$, მაგალითად, ჰისტოგრამის საშუალებით. მაშინ ნულოვან ჰიპოთეზას ექნება შემდეგი სახე: $H_0: F(x) = F^*(x)$. ემპირიული $F^*(x)$ და თეორიული $F(x)$ განაწილებების შედარება

ხდება პირსონის χ^2 თანხმობის კრიტერიუმით, კოლმოგოროვისა და სმირნოვის თანხმობის კრიტერიუმით და სხვ.

პირსონის თანხმობის კრიტერიუმი. იგი ძალიან ხშირად გამოიყენება პრაქტიკაში. ამისათვის, X შემთხვევითი სიდიდის ცვალებადობის მთელი დიაპაზონი დაეყოთ k რაოდენობის ინტერვალებად იგივე წესით, როგორც ჰისტოგრამის აგების დროს. შემდეგ თითოეული ინტერვალისთვის გამოვთვალოთ ემპირიული სიხშირეები და შევიტანოთ ეს მნიშვნელობები ცხრილში.

ინტერვალები	$[x^{(1)}; x^{(2)}[$	$[x^{(2)}; x^{(3)}[$...	$[x^{(k-1)}; x^{(k)}]$
ემპირიული სიხშირეები m_i	m_1	m_2	...	m_k
თეორიული სიხშირეები P_i	nP_1	nP_2	...	nP_k

თეორიული სიხშირეები გამოითვლება შემდეგნაირად: ცნობილია, რომ X შემთხვევითი სიდიდის $[x^{(i)}; x^{(i+1)}]$ ინტერვალში მოხვედრის P_i ალბათობა განისაზღვრება ფორმულით:

$$P_i = P(x^{(i)} \leq X \leq x^{(i+1)}) = \frac{1}{2} \left[\Phi \left(\frac{x^{(i+1)} - \bar{x}}{\sigma_x} \right) - \Phi \left(\frac{x^{(i)} - \bar{x}}{\sigma_x} \right) \right], i = 1, 2, \dots, k,$$

სადაც, $\Phi(\)$ – ლაპლასის ფუნქციაა. თუ მიღებულ P_i ალბათობებს გავამრავლებთ ამონარჩევის n განზომილებაზე, მაშინ მივიღებთ თითოეული ინტერვალისათვის თეორიული nP_i სიხშირეების მნიშვნელობებს.

თუ ემპირიული სიხშირეები მკვეთრად განსხვავდება თეორიულისგან, მაშინ ნულოვანი ჰიპოთეზა უარყოფილი იქნება, წინააღმდეგ შემთხვევაში, იგი მიიღება. ნულოვანი ჰიპოთეზის შესამოწმებლად გამოვთვალოთ სტატისტიკა

$$\chi^2 = \sum_{i=1}^k \frac{(m_i - nP_i)^2}{nP_i},$$

რომელსაც გააჩნია χ^2 განაწილება $\nu = k - r - 1$ თავისუფლების ხარისხით. აქ, r არის თეორიული $F(x)$ განაწილების პარამეტრების რაოდენობა, რომელიც ამონარჩევიდან გამოითვლება.

ნულოვანი ჰიპოთეზის შესამოწმებლად საჭიროა χ^2 განაწილების ცხრილიდან α მნიშვნელოვნების დონითა და ν თავისუფლების ხარისხით ვიპოვოთ $\chi_{\alpha;\nu}^2$ კრიტიკული მნიშვნელობა. თუ აღმოჩნდა, რომ $\chi^2 \geq \chi_{\alpha;\nu}^2$, მაშინ ნულოვანი ჰიპოთეზა უარყოფილია ალტერნატიულის სასარგებლოდ, ე.ი. ითვლება, რომ ემპირიული განაწილების ფუნქცია არ ემთხვევა თეორიულს. როცა $\chi^2 < \chi_{\alpha;\nu}^2$, მაშინ არა გვაქვს საფუძველი ნულოვანი ჰიპოთეზის უარსაყოფად.

მაგალითი. მოცემული $n = 200$ განზომილებიანი დაჯგუფებული სტატისტიკური X მწკრივისათვის

ინტერ- ვალები	[19;20[[20;21[[21;22[[22;23[[23;24[[24;25]
სიხში- რე m_i	10	26	56	64	30	14

შევამოწმოთ ნულოვანი ჰიპოთეზა ნორმალური განაწილების კანონის შესახებ.

ნულოვანი ჰიპოთეზის შესამოწმებლად განვსაზღვროთ a და s პარამეტრების წერტილოვანი შეფასებები

$$a = \bar{x} = 22,1; \quad s^2 = \sigma_x^2 = 1,52; \quad \sigma_x = 1,233.$$

გამოვთვალოთ X შემთხვევითი სიდიდის ინტერვალებში მოხვედრის P_i ალბათობები.

პირველი ინტერვალისთვის გვექნება:

$$\begin{aligned} P_1 &= P(19 < X < 20) = \frac{1}{2} \left[\Phi\left(\frac{20 - 22,1}{1,233}\right) - \Phi\left(\frac{19 - 22,1}{1,233}\right) \right] = \\ &= \frac{1}{2} [\Phi(-1,70) - \Phi(-2,51)] = \frac{1}{2} [\Phi(2,51) - \Phi(1,70)] = \\ &= \frac{1}{2} (0,9879 - 0,9109) = 0,0385. \end{aligned}$$

ანალოგიურად მივიღებთ: $P_2 = 0,142$, $P_3 = 0,281$, $P_4 = 0,299$, $P_5 = 0,171$, $P_6 = 0,052$. χ^2 სტატისტიკის გამოსათვლელად შევადგინოთ შემდეგი ცხრილი:

N_i	$x^{(i+1)} - x^{(i)}$	m_i	P_i	nP_i	$(m_i - nP_i)^2$	$\frac{(m_i - nP_i)^2}{nP_i}$
1	[19÷20[10	0,039	7,8	4,84	0,62
2	[20÷21[26	0,142	28,4	5,76	0,20
3	[21÷22[56	0,281	56,2	0,04	0,00
4	[22÷23[64	0,299	59,8	17,64	0,29
5	[23÷24[30	0,171	34,2	17,64	0,52
6	[24÷25]	14	0,052	10,4	12,96	1,25
Σ		200	0,913	197,3		$\chi^2 = 2,88$

$\alpha = 0,05$; $v = k - r - 1 = 6 - 2 - 1 = 3$; $\chi_{0,05;3}^2 = 7,815$.

რადგან $2,88 < 7,815$, ამიტომ ნულოვანი ჰიპოთეზა მიიღება, ე.ი. მოცემული ამონარჩევი ნორმალურადაა განაწილებული.

კოლმოგოროვ-სმირნოვის თანხმობის კრიტერიუმი. თუ ამონარჩევის მოცულობა მცირეა, მაშინ სპირმენის თანხმობის კრიტერიუმის გამოყენება შეუძლებელი ხდება. ამ შემთხვევაში, უნდა გამოვიყენოთ კოლმოგოროვ-სმირნოვის თანხმობის კრიტერიუმი. ამისათვის, ისევე როგორც სპირმენის კრიტერიუმის დროს, უნდა მოვახდინოთ მოცემული ამონარჩევის ინტერვალური დაჯგუფება. შემდეგ უნდა განისაზღვროს თითოეული ინტერვალისთვის სიხშირე m_i , დაგროვილი სიხშირე $\sum m_i$, ინტერვალში მოხვედრის P_i ალბათობა და დაგროვილი ალბათობა $\sum P_i$, რომელთა საშუალებით ხდება ემპირიული $F^*(x) = \frac{\sum m_i}{n}$ და თეორიული განაწილების $F(x) = \sum P_i$ ფუნქციების დადგენა. $H_0: F^*(x) = F(x)$ ჰიპოთეზის შესამოწმებლად უნდა განისაზღვროს კოლმოგოროვის λ სტატისტიკის მნიშვნელობა

$$\lambda = D\sqrt{n} = \max |F^*(x) - F(x)|\sqrt{n}.$$

α მნიშვნელოვნების დონით მოიძებნება λ_α კრიტიკული მნიშვნელობა, რომელიც მოცემულია შემდეგ ცხრილში:

α	0,20	0,10	0,05	0,02	0,01	0,001
λ_α	1,073	1,224	1,358	1,520	1,627	1,950

თუ აღმოჩნდება, რომ $\lambda \geq \lambda_\alpha$, მაშინ ნულოვანი ჰიპოთეზა უარყოფილია, ხოლო წინააღმდეგ შემთხვევაში, არ გვაქვს საფუძველი ნულოვანი ჰიპოთეზის უარსაყოფად.

მაგალითი. ზემოთ მოყვანილი მაგალითის მონაცემებისთვის გამოვიყენოთ კოლმოგოროვ-სმირნოვის კრიტერიუმი. ამისათვის შევადგინოთ შემდეგი ცხრილი:

ინტერვალები	სიხშირე m_i	დაგროვილი სიხშირე $\sum m_i$	ალბათობა P_i	დაგრ. ალბათობა $\sum P_i$	$F^*(x)$	$ F(x) - F^*(x) $
[19;20[10	10	0,039	0,039	0,05	0,011
[20;21[26	36	0,142	0,181	0,18	0,001
[21;22[56	92	0,281	0,462	0,46	0,002
[22;23[64	156	0,299	0,761	0,78	0,019
[23;24[30	186	0,171	0,932	0,93	0,002
[24;25]	14	200	0,052	0,984	1,00	0,016

$\lambda = 0,019\sqrt{200} = 0,269$. თუ $\alpha=0,05$, მაშინ $\lambda_{0,05} = 1,358$. რადგან $0,269 < 1,358$, ამიტომ ნულოვანი ჰიპოთეზა მიიღება, ე.ი. ამონარჩევი ნორმალურადაა განაწილებული.

ნორმალური განაწილების კანონის მიახლოებითი დადგენა შესაძლებელია როგორც ვიზუალურად, ასევე სტატისტიკური მახასიათებლებით. ვიზუალური შეფასებისთვის გამოიყენება ემპირიული განაწილების სიმკვრივის ფუნქციის გრაფიკული გამოსახულება (იხ. §6), ხოლო სტატისტიკური მახასიათებლებიდან – ასიმეტრიისა და ექსცესის კოეფიციენტები. ამისათვის,

ამონარჩევიდან, გარდა ასიმეტრიის A_x და ექსცესის E_x კოეფიციენტებისა, უნდა განისაზღვროს მათი სტატისტიკური შეცდომები ε_A და ε_E (იხ. §9.2). თუ $|A_x| < 3\varepsilon_A$ და $|E_x| < 3\varepsilon_E$, მაშინ მოცემული ამონარჩევი დაახლოებით ნორმალურად არის განაწილებული.

ნორმალური განაწილების კანონის მიახლოებითი შეფასება შესაძლებელია აგრეთვე ნორმალური ალბათობის გრაფიკით, რომელიც ისეა დაგრაფირებული, რომ თუ მასზე დავიტანთ ნორმალურად განაწილებული ამონარჩევის დაგროვილ სიხშირეებს, გამოსახულს პროცენტებში, მივიღებთ სწორ ხაზს.

11.2. ორი ამონარჩევის შედარების ჰიპოთეზის შემოწმება

U-კრიტერიუმი (უილკოქსონი-მანა-უიტნის კრიტერიუმი). ვთქვათ მოცემულია X და Y დამოუკიდებელი შემთხვევითი სიდიდეების ამონარჩევები x_1, x_2, \dots, x_n და y_1, y_2, \dots, y_m , რომელთა განაწილების კანონი უცნობია. საჭიროა შემოწმდეს ჰიპოთეზა, არის თუ არა განსხვავება ამ ორ ამონარჩევს შორის. ასეთი ნულოვანი ჰიპოთეზა შეიძლება შევამოწმოთ უილკოქსონის რანგული კრიტერიუმის საშუალებით. ზოგადად, U-კრიტერიუმით მოწმდება შემდეგი ნულოვანი ჰიპოთეზა: ორი დამოუკიდებელი ამონარჩევი მიეკუთვნება ერთი და იგივე გენერალურ ერთობლიობას და მათი განაწილების ფუნქციები ერთნაირია. ეს ჰიპოთეზა მოიცავს აგრეთვე განაწილების მდებარეობის მახასიათებლების, კერძოდ მედიანებისა და საშუალო მნიშვნელობების ტოლობას.

U-კრიტერიუმის გამოსათვლელად საჭიროა ამონარჩევები გავაერთიანოთ და დავალაგოთ ზრდადობით ერთ მწკრივში და შემდეგ გადავწომროთ. ამონარჩევის თითოეული მნიშვნელობა მიიღებს რიგით ნომერს, რომელსაც რანგი ეწოდება. თუ ამონარჩევში გვხვდება ერთნაირი სიდიდის რამდენიმე მაჩვენებელი, მაშინ თითოეულ მათგანს უნდა მივანიჭოთ საშუალო რანგი. შემდეგ გამოითვლება სიდიდეები

$$U_1 = R_x - \frac{n(n+1)}{2}, \quad U_2 = R_y - \frac{m(m+1)}{2},$$

სადაც, R_x არის პირველი ამონარჩევის რანგების ჯამი

$$R_x = \sum_{i=1}^n r(x_i),$$

ხოლო R_y – მეორე ამონარჩევის რანგების ჯამი

$$R_y = \sum_{i=1}^n r(y_i).$$

თუ ვიყენებთ ორმხრივ U -კრიტერიუმს, მაშინ α , n და m მნიშვნელობებით, სადაც, $n \geq m$, სპეციალური ცხრილიდან მოიძებნება $U_{\alpha;n,m}$ კრიტიკული მნიშვნელობა. თუ $n < m$, მაშინ n -ით აღნიშნავენ დიდი მოცულობის ამონარჩევს. ნულოვანი ჰიპოთეზა მიიღება, თუ $U = \min(U_1, U_2) \geq U_{\alpha;n,m}$, ე.ი. განსხვავება ორ ამონარჩევს შორის არ არსებობს. წინააღმდეგ შემთხვევაში, როცა $U < U_{\alpha;n,m}$, ნულოვანი ჰიპოთეზა უარყოფილია ალტერნატიულის სასარგებლოდ, ე.ი. განსხვავება ორ ამონარჩევს შორის სარწმუნოა.

თუ ამონარჩევის მოცულობები $m \rightarrow \infty$, $n \rightarrow \infty$, მაშინ U სტატისტიკა ასიმპტოტურად ნორმალურად არის განაწილებული $\frac{nm}{2}$ საშუალოთი და $\sigma^2 = \frac{1}{12}nm(n+m+1)$ დისპერსიით. მაშინ სიდიდეს

$$Z = \frac{U - \frac{1}{2}nm}{\sqrt{\frac{1}{12}nm(n+m+1)}}$$

გააჩნია ნორმალიზირებული ნორმალური განაწილება $N(0,1)$. როცა $n, m \geq 20$, მაშინ კრიტიკული წერტილის გამოსათვლელად შეიძლება გამოვიყენოთ მისი მიახლოებითი მნიშვნელობა, რომელიც განისაზღვრება შემდეგნაირად:

$$U_{\alpha;n,m} \approx \frac{1}{2}nm - \lambda_q \sqrt{\frac{1}{12}nm(n+m+1)},$$

სადაც, λ_q წარმოადგენს ნორმალიზირებული ნორმალური განაწილების კვანტილს და იგი განისაზღვრება α მნიშვნელოვნების დონით:

$$\lambda_q = \begin{cases} \lambda_{\frac{\alpha}{2}}, & \text{ორმხრივი კრიტერიუმისათვის,} \\ \lambda_{\alpha}, & \text{ცალმხრივი კრიტერიუმისათვის.} \end{cases}$$

λ_q მნიშვნელობები სხვადასხვა α -თვის მოცემულია შემდეგ ცხრილში:

α	0,1	0,05	0,025	0,01	0,005	0,001	0,0001
λ_{α}	1,282	1,645	1,960	2,326	2,576	3,090	3,719
$\lambda_{\frac{\alpha}{2}}$	1,645	1,960	1,241	2,576	2,807	3,291	3,891

მაგალითი. ცხრილში მოცემულია რეალური სისტემისა და მისი იმიტაციური მოდელის მუშაობის მაჩვენებლის შედეგები. საჭიროა დავადგინოთ, არის თუ არა განსხვავება რეალური სისტემის მუშაობასა და მის მოდელს შორის.

	რეალური სისტემა		იმიტაციური მოდელი	
	x	რანგი	y	რანგი
1	80	19	90	29,5
2	70	1	91	31
3	79	17,5	95	35,5
4	74	5	90	29,5
5	85	24	93	33
6	89	28	83	22
7	76	9	97	38
8	82	21	72	3
9	76	9	95	35,5
10	77	13,5	84	23
11	76	9	76	9
12	71	2	77	13,5
13	73	4	79	17,5
14	94	34	92	32
15	75	6	96	37
16	77	13,5	87	26
17	78	16	98	39
18	81	20	86	25
19			99	40

20			76	9
21			88	27
22			77	13,5
Σ	1413	251,3	1921	568,5

ე.ი. გვაქვს:

$$n = 18; m = 22; R_x = 251,3; R_y = 568,5;$$

$$U_1 = 251,3 - \frac{18 \cdot 19}{2} = 80,3; U_2 = 568,5 - \frac{22 \cdot 23}{2} = 315,5;$$

$$U = \min(80,3; 315,5) = 80,3;$$

$$U_{0,025} = \frac{1}{2} 18 \cdot 22 - 1,96 \sqrt{\frac{18 \cdot 22}{12} \cdot 41} = 191 - 1,96 \sqrt{1353} = 118,91.$$

რადგან $80,3 < 118,91$, ამიტომ განსხვავება რეალურ და იმიტაციური მოდელების მუშაობაში სარწმუნოა.

χ^2 კრიტერიუმი. თუ თავისებრივი მონაცემები მოცემულია ოთხუჯრედიანი ანუ (2x2) ცხრილის სახით, მაშინ ორ ამონარჩევს შორის განსხვავების ჰიპოთეზა შეიძლება შემონმდეს χ^2 კრიტერიუმის საშუალებით.

	ხდომილობა		სულ
	+	-	
პირველი ამონარ.	<i>a</i>	<i>b</i>	<i>a+b</i>
მეორე ამონარ.	<i>c</i>	<i>d</i>	<i>c+d</i>
სულ	<i>a+c</i>	<i>b+d</i>	<i>n</i>

განვიხილოთ სტატისტიკა

$$\chi^2 = \frac{n(|ad - bc| - 0,5n)^2}{(a+b)(c+d)(a+c)(b+d)},$$

რომელსაც გააჩნია χ^2 განაწილება $\nu = 1$ თავისუფლების ხარისხით. *a, b, c, d* - უჯრედებში მონაცემების მოხვედრის რაოდენობებია (სიხშირეები); $n = a + b + c + d$; თუ $\chi^2 \geq \chi_{\alpha, \nu}^2$, მაშინ ნულო-

ვანი ჰიპოთეზა უარყოფილია, ე.ი. განსხვავება ორ ამონარჩევს შორის სარწმუნოა. წინააღმდეგ შემთხვევაში, როცა $\chi^2 < \chi^2_{\alpha;v}$, მაშინ ნულოვანი ჰიპოთეზა მიიღება.

მაგალითი. ოთხუჯრედიან ცხრილში მოცემულია ახალი და ძველი პრეპარატების გამოყენების შედეგად მიღებული მონაცემები. გვანტერესებს, ახალი პრეპარატის გამოყენებით მკურნალობა ეფექტურია თუ არა.

მკურნალობა	გარდაიცვალა	გამოჯანმრთ.	სულ
ახალი	15	85	100
ძველი	4	77	81
სულ	19	162	181

$$\chi^2 = \frac{181(15 \cdot 77 - 4 \cdot 85 - 90,5)^2}{100 \cdot 81 \cdot 19 \cdot 162} = 3,81.$$

χ^2 განაწილების ცხრილიდან $\chi^2_{0,05;1} = 3,84$.

რადგან $3,81 < 3,84$, ამიტომ ნულოვანი ჰიპოთეზა მიიღება, ე.ი. ახალი პრეპარატით მკურნალობა რომ უფრო ეფექტურია, არ დასტურდება.

11.3. ამონარჩევების მრავლობითი შედარება

მრავლობითი შედარების პარამეტრული მეთოდები ადვილად ადაპტირდება არაპარამეტრულ მეთოდებში. კერძოდ, როდესაც ამონარჩევების განზომილება ერთნაირია, მაშინ შეგვიძლია გამოვიყენოთ ნიუმენ-კეილსის არაპარამეტრული (რანგული) კრიტერიუმი, ხოლო როდესაც გვაქვს სხვადასხვა განზომილებიანი ამონარჩევები – დანას კრიტერიუმი.

ვთქვათ, მოცემულია m რაოდენობის n -განზომილებიანი ამონარჩევი. მოვახდინოთ ამონარჩევების ერთდროული რანჟირება და ყოველ მათგანს მივანიჭოთ რანგი. მაშინ ნიუმენ-კეილსის რანგულ კრიტერიუმს აქვს შემდეგი სახე:

$$q = \frac{|R_i - R_j|}{\sqrt{\frac{n^2 l(nl + 1)}{12}}},$$

სადაც, R_i და R_j შესადარებელი i -ური და j -ური ამონარჩევების რანგების ჯამია. l – შედარების ინტერვალი, რომელიც განისაზღვრება ისევე, როგორც პარამეტრული მეთოდის დროს. α მნიშვნელოვნების დონით, $v = \infty$ და l სიდიდებით ცხრილიდან შეირჩევა $q_{\alpha;v;l}$ კრიტიკული მნიშვნელობა. თუ $q \geq q_{\alpha;v;l}$, მაშინ ნულოვანი ჰიპოთეზა უარყოფილია, ე.ი. განსხვავება ამონარჩევებს შორის სარწმუნოა. წინააღმდეგ შემთხვევაში, როცა $q < q_{\alpha;v;l}$, ნულოვანი ჰიპოთეზა მიიღება და ამონარჩევები ერთმანეთისგან არ განსხვავდება.

თუ ამონარჩევები სხვადასხვა განზომილებისაა, მაშინ მათი წყვილ-წყვილად შედარებისათვის უნდა გამოვიყენოთ **დანას კრიტერიუმი** (Q კრიტერიუმი):

$$Q = \frac{|\bar{R}_i - \bar{R}_j|}{\sqrt{\frac{N(N+1)}{12} \left(\frac{1}{n_i} + \frac{1}{n_j} \right)}},$$

სადაც, \bar{R}_i და \bar{R}_j i -ური და j -ური ამონარჩევების საშუალო რანგებია, n_i, n_j – ამონარჩევების განზომილებები, N — ყველა

ამონარჩევის ჯამური განზომილება, ე.ი. $N = \sum_{i=1}^m n_i$. Q კრიტიკული

მნიშვნელობები მოცემულია სპეციალურ ცხრილში (იხ. დანართი). α მნიშვნელოვნების დონითა და m სიდიდით ცხრილიდან შეირჩევა $Q_{\alpha;m}$ კრიტიკული მნიშვნელობა. თუ $Q \geq Q_{\alpha;m}$, მაშინ ამონარჩევები განსხვავდება ერთმანეთისგან, წინააღმდეგ შემთხვევაში, როცა $Q < Q_{\alpha;m}$, ისინი არ განსხვავდებიან ერთმანეთისგან. უნდა აღინიშნოს, რომ დანას კრიტერიუმი შეიძლება გამოვიყენოთ მაშინაც, როცა ამონარჩევებს გააჩნიათ ერთნაირი განზომილება.

მაგალითი. ცხრილში მოცემულია პულისის სიხშირის მნიშვნელობები 5 წლამდე ბავშვების სამი ჯგუფისთვის (x, y, z). დავადგინოთ, არის თუ არა განსხვავება ჯგუფებს შორის.

	x	y	z	R_x	R_y	R_z
1	112	90	110	2,5	2,5	9,5
2	116	78	117	16	1	17,5
3	108	109	112	6	7,5	12,5
4	120	92	115	19	4	14,5
5	109	100	118	7,5	5	19
6	90	115	117	2,5	14,5	17,5
7	110		125	9,5		20
8	111			11		
Σ				84,0	34,5	110,5

$N = 8 + 6 + 7 = 21$; საშუალო რანგებია $\bar{R}_x = 10,5$; $\bar{R}_y = 5,75$; $\bar{R}_z = 15,79$. კრიტიკული მნიშვნელობა $Q_{0,05;3} = 2,39$. დანას კრიტერიუმით მოვახდინოთ ჯგუფების წყვილ-წყვილად შედარება. მივიღებთ:

$$Q_{xy} = \frac{|10,5 - 5,75|}{\sqrt{\frac{21 \cdot 22}{12} \left(\frac{1}{8} + \frac{1}{6} \right)}} = 1,42; \quad Q_{xz} = \frac{|10,5 - 15,79|}{\sqrt{\frac{21 \cdot 22}{12} \left(\frac{1}{8} + \frac{1}{7} \right)}} = 1,65;$$

$$Q_{yz} = \frac{|5,75 - 15,79|}{\sqrt{\frac{21 \cdot 22}{12} \left(\frac{1}{6} + \frac{1}{7} \right)}} = 2,91.$$

როგორც ვხედავთ, მხოლოდ y და z ჯგუფები განსხვავდებიან ერთმანეთისგან, რადგან $2,91 > 2,39$. დანარჩენ ორ შემთხვევაში ჯგუფებს შორის განსხვავება არ შეიმჩნევა.

11.4. ორი დამოკიდებული ამონარჩევის შედარება

თუ ამონარჩევები x_1, x_2, \dots, x_n და y_1, y_2, \dots, y_n ერთმანეთის მიმართ დამოკიდებულია და მათი განაწილების კანონი უცნობია,

მაშინ უმჯობესია გამოვიყენოთ **უილკოქსონის T-კრიტერიუმი**. ამისათვის მოცემული ამონარჩევებიდან განვსაზღვროთ $d_i = x_i - y_i, i = 1, 2, \dots, n$ სხვაობები. მიღებული სხვაობებიდან უნდა გამოვრიცხოთ ნულოვანი მნიშვნელობები და დარჩენილების აბსოლუტური სიდიდით რანჟირების შემდეგ, მათ მივანიჭოთ რანგები. შემდეგ უნდა მოიძებნოს ცალკე დადებითი R_+ და ცალკე უარყოფითი R_- სხვაობების რანგების ჯამი, რომლებიც უნდა აკმაყოფილებდეს შემდეგ პირობას:

$$R_+ + R_- = \frac{n(n+1)}{2}.$$

ნულოვანი ჰიპოთეზა უარყოფილია, თუ $R = \min(R_+, R_-) \leq T_{\alpha, n}$, სადაც, $T_{\alpha, n}$ კრიტიკული მნიშვნელობა მოიძებნება სპეციალური ცხრილიდან (იხ. დანართი) ან ის შეგვიძლია განვსაზღვროთ (როცა $n > 25$) შემდეგი აპროქსიმაციის გამოყენებით:

$$T_{\alpha, n} \approx \frac{n(n+1)}{4} - \lambda_q \sqrt{\frac{1}{24} n(n+1)(2n+1)},$$

სადაც, λ_q ნორმალიზირებული ნორმალური განაწილების კვანტილია, რომელიც განისაზღვრება α მნიშვნელოვნების დონით (იხ. §11.2)

მაგალითი. ცხრილში მოცემულია ქერისა (x) და შვრიის (y) მოსავალი (ცენტნ/ჰ) შვიდი წლის განმავლობაში. გვაინტერესებს, არის თუ არა განსხვავება საშუალო მოსავლებს შორის.

	მოსავალი		d	R_d
	x	y		
1	7,7	8,26	-0,56	1
2	9,0	7,22	1,78	5
3	9,4	8,43	0,97	2
4	7,4	5,57	1,83	6
5	7,4	6,35	1,05	3
6	10,9	8,00	2,90	7
7	8,0	9,13	-1,13	4
Σ	59,8	52,96	6,84	
საშ.	8,54	7,56	0,98	

$$R_+ = (5 + 2 + 6 + 3 + 7) = 23 ; \quad R_- = 1 + 4 = 5.$$

$\alpha = 0,05$ და $n = 7$ მნიშვნელობებით T -კრიტერიუმის ცხრილიდან ვიღებთ $T_{0,05;7} = 3$. რადგან $T = \min(5; 23) = 5 > 3$, ამიტომ ნულოვანი ჰიპოთეზა მიიღება, ე.ი. განსხვავება საშუალო მოსავლებს შორის არ დასტურდება.

11.5. ფარდობითი სიხშირეების ტოლობის ჰიპოთეზის შემოწმება

თუ მოცემულია n_1 და n_2 განზომილებიანი ამონარჩევების ფარდობითი სიხშირეები $P_1^* = \frac{m_1}{n_1}$ და $P_2^* = \frac{m_2}{n_2}$ და გვინდა შევამოწმოთ $H_0: P_1^* = P_2^*$ ნულოვანი ჰიპოთეზა, მაშინ უნდა განვიხილოთ შემდეგი სტატისტიკა:

$$Z = \frac{|P_1^* - P_2^*|}{\sqrt{\frac{P_1^*(1-P_1^*)}{n_1} + \frac{P_2^*(1-P_2^*)}{n_2}}}.$$

თუ P_1^* და P_2^* წარმოადგენენ ერთი და იგივე P^* ფარდობითი სიხშირის შეფასებებს, მაშინ საშუალო კვადრატული გადახრა შეგვიძლია განვსაზღვროთ შემდეგი ფორმულით:

$$\sigma_p = \sqrt{P^*(1-P^*)}, \quad \text{სადაც, } P^* = \frac{m_1 + m_2}{n_1 + n_2}$$

და Z კრიტერიუმს ექნება შემდეგი სახე:

$$Z = \frac{|P_1^* - P_2^*|}{\sqrt{P^*(1-P^*) \left(\frac{1}{n_1} + \frac{1}{n_2} \right)}}.$$

Z კრიტერიუმს გააჩნია ნორმალიზებული ნორმალური განაწილება. ამიტომ მისი კრიტიკული მნიშვნელობის მოძებნა შეიძლება სტანდარტიზირებული ნორმალური განაწილების საშუალებით. რადგან სტიუდენტის განაწილება, როცა თავისუფლების ხარისხი იზრდება, სწრაფად გადადის ნორმალურ განაწილებაში, ამიტომ Z -ის კრიტიკული მნიშვნელობის მოსაძებნად, შეიძლება გამოვიყენოთ სტიუდენტის განაწილების ცხრილი α მნიშვნელოვნების დონითა და $\nu = \infty$ თავისუფლების ხარისხით. თუ $Z < t_{\alpha; \nu}$, მაშინ ნულოვანი ჰიპოთეზა მიიღება. წინააღმდეგ შემთხვევაში, იგი უარყოფილია ალტერნატიულის სასარგებლოდ.

რადგან Z -ს გააჩნია მიახლოებით ნორმალური განაწილება, ამიტომ P^* შეფასების სიდიდე შემცირებულია, რაც იწვევს ნულოვანი ჰიპოთეზის ძალზე ხშირად უარყოფას. ეს გამოწვეულია იმით, რომ Z იღებს მხოლოდ დისკრეტულ მნიშვნელობებს მაშინ, როცა ნორმალური განაწილება უწყვეტია. ამ ფაქტის კომპენსაცია შესაძლებელია, თუ გამოვიყენებთ **იეიტსის შესწორებას**. მაშინ Z გამოსახულებას ექნება შემდეგი სახე:

$$Z = \frac{|P_1^* - P_2^*| - \frac{1}{2} \left(\frac{1}{n_1} + \frac{1}{n_2} \right)}{\sqrt{P^*(1-P^*) \left(\frac{1}{n_1} + \frac{1}{n_2} \right)}}$$

იეიტსის შესწორება მცირეოდენად ამცირებს Z -ის მნიშვნელობას, რაც, თავის მხრივ, იწვევს ნორმალურ განაწილებასთან განსხვავების შემცირებას.

მაგალითი. ჰემოღიალიზზე მყოფი პაციენტების შუნტის თრომბოზის გამოკვლევა. მიღებულ იქნა ორი ჯგუფი: პირველ ჯგუფში პაციენტები იღებდნენ პლაცებოს, მეორეში – ასპირინს (ასპირინი თრომბოზის წარმოშობას ხელს უშლის). პირველ ჯგუფში 25 პაციენტიდან 18-ს განუვითარდა შუნტის თრომბოზი, მეორე ჯგუფში – 19-დან 6-ს. დავადგინოთ, რამდენად ეფექტურია ასპირინის გამოყენება.

გამოვთვალოთ ფარდობითი სიხშირეები $P_1^* = \frac{18}{25} = 0,72$;

$P_2^* = \frac{6}{19} = 0,38$. გაერთიანებული ფარდობითი სიხშირე ტოლია:

$P^* = \frac{6+18}{19+25} = 0,55$. რადგან $n_1 P_1^* = 18 > 5$, $n_1(1 - P_1^*) = 7 > 5$,

$n_2 P_2^* = 6 > 5$ და $n_2(1 - P_2^*) = 13 > 5$, ამიტომ Z კრიტერიუმის გამოყენება შესაძლებელია (იხ. §8.4).

$$Z = \frac{|0,72 - 0,32| - 0,05}{\sqrt{0,55(1 - 0,55)\left(\frac{1}{25} + \frac{1}{19}\right)}} = 2,33.$$

სტიუდენტის განაწილების ცხრილიდან $t_{0,05;\infty} = 1,96$.

რადგან $2,33 > 1,96$, ამიტომ ნულოვანი ჰიპოთეზა უარყოფილია, ე.ი. ასპირინის გამოყენება თრომბის წარმოშობის წინააღმდეგ ეფექტურია.

12. კორელაციური ანალიზი

12.1. კორელაციური ანალიზის არსი

პრაქტიკულ კვლევებში ძალიან ხშირად საჭირო ხდება ორ შემთხვევით სიდიდეს შორის დამოკიდებულების გამოკვლევა და შემდეგ მისი შეფასება. ორი სიდიდე ერთმანეთის მიმართ შეიძლება იყოს ფუნქციონალურ ან სხვა სახის დამოკიდებულებაში, რომელსაც სტატისტიკურს უწოდებენ, ან კიდევ ერთმანეთის მიმართ იყვნენ დამოუკიდებელი.

ტექნიკასა და საბუნებისმეტყველო მეცნიერებაში საქმე ეხება X და Y ცვლადებს შორის ფუნქციონალურ დამოკიდებულებას, როცა X -ის ყოველ შესაძლო მნიშვნელობას ცალსახად შეესაბამება Y -ის მნიშვნელობა. რეალურ სამყაროში ბუნების

მრავალი მოვლენა მიმდინარეობს მრავალრიცხოვანი ფაქტორების გარემოცვაში, ამიტომ კავშირი ცალსახა აღარ გამოდის. აქ შეიძლება ლაპარაკი მხოლოდ სტატისტიკურ კავშირზე.

სტატისტიკური კავშირი მდგომარეობს იმაში, რომ ერთი შემთხვევითი სიდიდე რეაგირებს მეორე შემთხვევითი სიდიდის ცვლილებაზე. შემთხვევით სიდიდეთა შორის ამგვარ დამოკიდებულებას ეწოდება კორელაცია (ლათინური სიტყვა *correlatio*-დან, რაც ნიშნავს თანაფარდობას, კავშირს).

კორელაციური ანალიზის ძირითადი ამოცანაა შემთხვევით სიდიდეთა შორის კავშირის გამოვლენა. კორელაციური ანალიზის მოთხოვნებია: ცვლადები უნდა იყვნენ შემთხვევითი სიდიდეები და შემთხვევით სიდიდეებს უნდა ჰქონდეთ ერთობლივი ნორმალური განაწილება.

ორ ცვლადს შორის წრფივი კავშირის დასადგენად საჭიროა, გამოვთვალოთ კორელაციის კოეფიციენტი, ხოლო არა-წრფივი კავშირისათვის – კორელაციური ფარდობა η (იხ. §13.5). კორელაციის კოეფიციენტი, რომელიც r_{xy} ან *Corr*(X , Y) სიმბოლოებით აღინიშნება, განყენებული რიცხვია და იცვლება $-1 \leq r_{xy} \leq 1$ ფარგლებში. თუ $r_{xy} = 0$, მაშინ კავშირი X და Y შემთხვევით სიდიდეებს შორის არ არსებობს. რაც უფრო ძლიერია კავშირი შემთხვევით სიდიდეებს შორის, მით უფრო დიდია კორელაციის კოეფიციენტი. თუ $r_{xy} = \pm 1$, მაშინ ცვლადებს შორის არსებობს ფუნქციონალური კავშირი.

დადებითი ანუ პირდაპირი წრფივი კავშირის დროს, როცა ერთი ცვლადის ზრდისას იზრდება მეორე ცვლადიც, კორელაციის კოეფიციენტს აქვს დადებითი ნიშანი და იცვლება ნულსა და ერთს შორის. უარყოფითი, ანუ უკუკავშირის დროს, როცა ერთი ცვლადის ზრდისას მეორე მცირდება, კორელაციის კოეფიციენტს გააჩნია უარყოფითი ნიშანი და იგი იცვლება 0-სა და -1 შორის.

კორელაციის კოეფიციენტი არ არის დამოკიდებული ათვლის წერტილსა და გაზომვის ერთეულზე, ე.ი. X და Y სიდიდეები შეიძლება რამდენჯერმე გავზარდოთ ან შევამციროთ, აგრეთვე მივუმატოთ ან გამოვაკლოთ რაიმე რიცხვი, ამით კორელაციის კოეფიციენტის სიდიდე არ შეიცვლება.

როცა $r_{xy} = 0$, მაშინ X და Y შემთხვევითი სიდიდეები არაკორელირებულია. არაკორელირების ცნება არ უნდა ავურიოთ დამოუკიდებლობის ცნებაში. დამოუკიდებელი სიდიდეები ყოველთვის არაკორელირებულია, მაგრამ შეებრუნებული მტკიცება

არასწორია. არაკორელირებული სიდიდეები შეიძლება დამოკიდებული იყვნენ, თანაც ფუნქციონალურად, თუმცა ეს კავშირი არანრფივია.

12.2. კორელაციის კოეფიციენტის განსაზღვრის პარამეტრული მეთოდი

პირსონის კორელაციის კოეფიციენტი. ვთქვათ, X და Y შემთხვევითი სიდიდეებია, რომლებსაც გააჩნიათ ერთობლივი ნორმალური განაწილება და მოცემულია მათი ამონარჩევები x_1, x_2, \dots, x_n და y_1, y_2, \dots, y_n . მაშინ მათ შორის ნრფივი კავშირი შეიძლება აღინეროს კორელაციის კოეფიციენტით, რომელიც გამოითვლება შემდეგი ექვივალენტური ფორმულებით:

$$r_{xy} = \frac{\sum_{i=1}^n (x_i - \bar{x})(y_i - \bar{y})}{n\sigma_x\sigma_y};$$

$$r_{xy} = \frac{\sum_{i=1}^n (x_i - \bar{x})(y_i - \bar{y})}{\sqrt{\sum_{i=1}^n (x_i - \bar{x})^2 \sum_{i=1}^n (y_i - \bar{y})^2}}; \quad (12.1)$$

$$r_{xy} = \frac{\sum_{i=1}^n x_i y_i - \frac{1}{n} \left(\sum_{i=1}^n x_i \right) \left(\sum_{i=1}^n y_i \right)}{\sqrt{\left[\sum_{i=1}^n x_i^2 - \frac{1}{n} \left(\sum_{i=1}^n x_i \right)^2 \right] \left[\sum_{i=1}^n y_i^2 - \frac{1}{n} \left(\sum_{i=1}^n y_i \right)^2 \right]}}$$

სადაც, \bar{x}, \bar{y} – საშუალო არითმეტიკულებია, σ_x, σ_y – საშუალო კვადრატული გადახრები, n – ამონარჩევების განზომილება. r_{xy} კოეფიციენტს ეწოდება კორელაციის ემპირიული კოეფიციენტი ან პირსონის კორელაციის კოეფიციენტი.

პრაქტიკულად, კორელაციის კოეფიციენტი ჩვეულებრივ უცნობია. ამონარჩევების მიხედვით ჩვენ შეგვიძლია ვიპოვოთ მისი

ნერტილოვანი შეფასება, რომლის სტატისტიკური შეცდომა გამოითვლება შემდეგი ფორმულით:

$$\varepsilon_r = \sqrt{\frac{1-r_{xy}^2}{n-2}}$$

იმისათვის, რომ გამოვარკვიოთ, იმყოფებიან თუ არა შემთხვევითი სიდიდეები კორელაციურ დამოკიდებულებაში, უნდა შემოწმდეს $H_0: r_{xy} = 0$ ნულოვანი ჰიპოთეზა. ამისათვის საჭიროა

გამოვთვალოთ შემდეგი სტატისტიკა: $t = \frac{r_{xy}}{\varepsilon_r}$ ან $t = |r_{xy}| \sqrt{\frac{n-2}{1-r_{xy}^2}}$,

რომელსაც გააჩნია სტიუდენტის განაწილება $v=n-2$ თავისუფლების ხარისხით. სტიუდენტის განაწილების ცხრილიდან α მნიშვნელოვნების დონითა და v თავისუფლების ხარისხით მოიძებნება $t_{\alpha;v}$ კრიტიკული მნიშვნელობა. თუ $t \geq t_{\alpha;v}$, მაშინ ნულოვანი ჰიპოთეზა უარყოფილია, ე.ი. კავშირი X და Y ცვლადებს შორის სარწმუნოა. თუ $t < t_{\alpha;v}$, მაშინ არა გვაქვს საფუძველი ნულოვანი ჰიპოთეზის უარსაყოფად, ე.ი. კავშირი X და Y ცვლადებს შორის არ არსებობს.

აღმოჩნდა, რომ როცა ამონარჩევის განზომილება მცირეა ($n < 30$), კორელაციის კოეფიციენტის სიდიდე, გამოთვლილი (12.1) ფორმულით, იძლევა გენერალური ერთობლიობის კორელაციის კოეფიციენტის შემცირებულ მნიშვნელობას. ამ შემთხვევაში,

სასურველია r_{xy} სიდიდის კორექტირება $\left[1 + \frac{1-r_{xy}^2}{2(n-3)}\right]$ სიდიდით,

რომელზედაც უნდა გამრავლდეს გამოთვლილი კორელაციის კოეფიციენტი, ე.ი.

$$r_{xy}^* = r_{xy} \left[1 + \frac{1-r_{xy}^2}{2(n-3)}\right].$$

მაგალითი. ცხრილში მოცემულია 15 სტიუდენტის სიმაღლე (x) სმ-ში და წონა (y) კგ-ში. პირსონის კორელაციის კოეფიციენტით დავადგინოთ, არსებობს თუ არა კავშირი ამ ორ სიდიდეს შორის.

№	x_i	y_i	x_i^2	y_i^2	$x_i y_i$
1	1,65	72,9	2,72	5314,41	120,29
2	1,71	48,4	2,92	2342,5	82,76
3	1,82	66,3	3,31	4395,69	120,67
4	1,65	64,1	2,72	4108,81	105,77
5	1,83	62,7	3,35	3931,29	114,74
6	1,80	76,0	3,24	5776,0	136,80
7	1,83	72,8	3,35	5299,84	133,22
8	1,66	50,6	2,76	2560,36	83,99
9	1,73	52,3	2,99	2735,29	90,48
10	1,84	68,6	3,39	4705,96	126,22
11	1,68	52,6	2,82	2766,76	88,37
12	1,64	72,8	2,69	5299,84	119,39
13	1,70	61,6	2,89	3794,56	104,72
14	1,74	66,8	3,03	4462,24	116,23
15	1,72	56,5	2,96	3192,25	97,18
Σ	26,0	945,0	45,14	60685,25	1640,83

კორელაციის კოეფიციენტი ტოლია:

$$r_{xy} = \frac{1640,83 - \frac{26 \cdot 945}{15}}{\sqrt{\left(45,14 - \frac{26^2}{15}\right) \left(60685,25 - \frac{945^2}{15}\right)}} = 0,32.$$

შესწორებული კორელაციის კოეფიციენტი ტოლია:

$$r_{xy}^* = 0,32 \left[1 + \frac{1 - 0,32^2}{2 \cdot 12} \right] = 0,33.$$

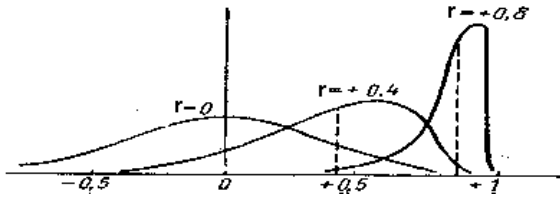
შევამოწმოთ $H_0 : 0,33 = 0$ ნულოვანი ჰიპოთეზა.

$$\varepsilon_r = \sqrt{\frac{1 - 0,33^2}{13}} = 0,26; \quad t = \frac{0,33}{0,26} = 1,27.$$

სტიუდენტის განაწილების ცხრილიდან $v = n - 2 = 13$ და $\alpha = 0,05$ მნიშვნელობებისათვის ვღებულობთ $t_{0,05;13} = 1,77$. რადგან $1,27 < 1,77$, ამიტომ ნულოვანი ჰიპოთეზა მიიღება, ე.ი. კავშირი

სიმაღლესა და წონას შორის არ შეიმჩნევა. კორელაციის კოეფიციენტის მნიშვნელობა ჩაინერება შემდეგნაირად: $r_{xy} = 0,33 \pm 0,26$.

როგორც აღვნიშნეთ, კორელაციური ანალიზი გამოიყენება იმ შემთხვევაში, როდესაც ამონარჩევები ნორმალურად არიან განაწილებული. მათემატიკური სტატისტიკიდან ცნობილია, რომ ორ ცვლადს შორის ძლიერი სტატისტიკური კავშირის დროს ($r_{xy} > 0,5$) კორელაციის კოეფიციენტის განაწილება მცირე ამონარჩევების დროს მნიშვნელოვნად განსხვავდება ნორმალურისგან, როგორც ეს შემდეგ ნახაზზეა ნაჩვენები:



ნახაზზე წარმოდგენილია $n = 12$ განზომილებიანი ამონარჩევების ემპირიული კორელაციის კოეფიციენტების განაწილების მრუდები გენერალური პარამეტრის $r = 0; 0,4$ და $0,8$ მნიშვნელობებისთვის. როგორც ნახაზიდან ჩანს, კორელაციის კოეფიციენტის განაწილებას, როდესაც მისი სიდიდე ერთთან ახლოა, გააჩნია ძლიერი ასიმეტრია. ამიტომ შერჩევითი კორელაციის კოეფიციენტი, რომლის სიდიდე $0,5$ -ზე მეტია და გამოთვლილია მცირე ამონარჩევით, არ იქნება გენერალური ერთობლიობის კორელაციის კოეფიციენტის ზუსტი შეფასება. გაითვალისწინა რა ეს, რ. ფიშერმა r_{xy} სიდიდის მაგივრად შემოგვთავაზა Z სიდიდე, რომელიც ასე გამოითვლება:

$$Z = \frac{1}{2} \ln \frac{1+r_{xy}}{1-r_{xy}} \quad \text{ან} \quad Z = 1,15129 \lg \frac{1+r_{xy}}{1-r_{xy}}$$

Z სიდიდის განაწილება თითქმის არ იცვლება, რადგან იგი ნაკლებადაა დამოკიდებული ამონარჩევის მოცულობაზე. აქვე უნდა აღვნიშნოთ, რომ თუ კორელაციის კოეფიციენტი იცვლება -1 -დან $+1$ -მდე, Z სიდიდე იცვლება $-\infty$ -დან $+\infty$ -მდე და მისი განაწილება სწრაფად უახლოვდება ნორმალურს.

$H_0: Z = 0$ ნულოვანი ჰიპოთეზის შესამოწმებლად, უნდა გამოვთვალოთ $t = Z\sqrt{n-3}$ სიდიდე, რომელსაც გააჩნია სტიუდენტის

ტის განაწილება $v = n - 2$ თავისუფლების ხარისხით. როცა $t < t_{\alpha;v}$, მაშინ ნულოვანი ჰიპოთეზა მიიღება, წინააღმდეგ შემთხვევაში იგი უარყოფილი იქნება.

კორელაციის კოეფიციენტის ზუსტი შეფასებისთვის ანუ სტატისტიკური კავშირი რომ იყოს სარწმუნო, შეგვიძლია გამოვთვალოთ ამონარჩევის მინიმალური განზომილება ფორმულით:

$$n = \frac{t^2}{Z^2} + 3,$$

სადაც, t სიდიდე აიღება სტიუდენტის განაწილების ცხრილიდან $\alpha = 0,01$ და $v = \infty$ სიდიდეებისთვის.

კორელაციის კერძო კოეფიციენტი. მრავალგანზომილებიანი X_1, X_2, \dots, X_m სისტემის დროს, ორ პარამეტრს შორის დამოკიდებულება შეიძლება წარმოიშვას იმ მიზეზითაც, რომ ორივე ეს პარამეტრი ცოტად თუ ბევრად იმყოფებიან სხვა პარამეტრთან ან პარამეტრებთან ისეთ ძლიერ დამოკიდებულებაში, რომ ეს დამოკიდებულება იწვევდეს მოცემულ პარამეტრებს შორის კავშირს, როცა რეალურად ეს კავშირი შეიძლება არც კი არსებობდეს ან არ იწვევდეს კავშირს, როცა რეალურად ეს კავშირი არსებობს.

იმისათვის, რომ გამოვრიცხოთ სხვა პარამეტრების გავლენა, შემოაქვთ კორელაციის კერძო კოეფიციენტის ცნება, რომლის განსაზღვრისათვის წინასწარ უნდა განისაზღვროს კორელაციური მატრიცა:

$$R = \begin{bmatrix} 1 & r_{12} & \dots & r_{1m} \\ r_{21} & 1 & \dots & r_{2m} \\ \dots & \dots & \dots & \dots \\ r_{m1} & r_{m2} & \dots & 1 \end{bmatrix},$$

რომლის ელემენტებს პირსონის კორელაციის კოეფიციენტები წარმოადგენენ. ზოგადი სახით კორელაციის კერძო კოეფიციენტი გამოითვლება ფორმულით:

$$r_{ij(1,2,\dots,p)} = -\frac{R_{ij}}{\sqrt{R_{ii}R_{jj}}},$$

სადაც, R_{ij} – კორელაციური მატრიცის r_{ij} ელემენტის ალგებრული დამატებაა; R_{ii} , R_{jj} – შესაბამისად r_{ii} , r_{jj} ელემენტების ალგებრული დამატებებია.

კერძოდ, თუ მოცემულია სამი X , Y და Z შემთხვევითი სიდიდის ამონარჩევები, მაშინ კორელაციის კერძო კოეფიციენტები გამოითვლება შემდეგი ფორმულებით:

$$r_{xy(z)} = \frac{r_{xy} - r_{xz}r_{yz}}{\sqrt{(1-r_{xz}^2)(1-r_{yz}^2)}};$$

$$r_{xz(y)} = \frac{r_{xz} - r_{xy}r_{zy}}{\sqrt{(1-r_{xy}^2)(1-r_{zy}^2)}};$$

$$r_{yz(x)} = \frac{r_{yz} - r_{yx}r_{zx}}{\sqrt{(1-r_{yx}^2)(1-r_{zx}^2)}}.$$

კორელაციის კერძო კოეფიციენტი იცვლება -1 -დან $+1$ -მდე და მისი სარწმუნოება ფასდება ისევე, როგორც პირსონის კორელაციის კოეფიციენტის დროს.

მაგალითი. მოცემული კორელაციური მატრიცისათვის ($n=10$)

$$R = \begin{bmatrix} 1 & 0,865 & 0,853 \\ & 1 & 0,950 \\ & & 1 \end{bmatrix}.$$

გამოვთვალოთ კორელაციის კერძო კოეფიციენტები. ე.ი. გვაქვს: $r_{xy} = 0,865$, $r_{xz} = 0,853$ და $r_{yz} = 0,950$, მაშინ მივიღებთ:

$$r_{xy(z)} = \frac{0,865 - 0,853 \cdot 0,95}{\sqrt{(1-0,853^2)(1-0,95^2)}} = \frac{0,055}{\sqrt{0,027}} = 0,335;$$

$$r_{xz(y)} = \frac{0,853 - 0,865 \cdot 0,95}{\sqrt{(1-0,865^2)(1-0,95^2)}} = \frac{0,032}{\sqrt{0,024}} = 0,20;$$

$$r_{yz(x)} = \frac{0,95 - 0,865 \cdot 0,853}{\sqrt{(1-0,865^2)(1-0,853^2)}} = \frac{0,212}{\sqrt{0,06886}} = 0,809.$$

ყველაზე მაღალი კორელაციის კოეფიციენტი აღმოჩნდა $r_{yz(x)}$. შემდგომში ნულოვანი ჰიპოთეზა $H_0 : 0,809 = 0$. ამისათვის გამოვთვალოთ სტატისტიკა

$$t = 0,809 \sqrt{\frac{10-2}{1-0,809^2}} = 0,809 \sqrt{23,15} = 3,89; \alpha = 0,05; \nu = 10 - 2 = 8.$$

სტიუდენტის განაწილების ცხრილიდან $t_{0,05;8} = 1,86$. რადგან $3,89 > 1,86$, ამიტომ ნულოვანი ჰიპოთეზა უარყოფილია, ე.ი. კავშირი y და z ცვლადებს შორის სარწმუნოა, ხოლო x და y და x და z ცვლადებს შორის – ნაკლებად სარწმუნო. მართლაც,

$$t = 0,335 \sqrt{\frac{8}{1-0,335^2}} = 1,01 < 1,86.$$

მრავლობითი კორელაციის კოეფიციენტი. პრაქტიკაში ხშირად საინტერესოა შეფასდეს ერთი პარამეტრის კავშირი სხვა დანარჩენ პარამეტრებთან. ეს შეიძლება გაკეთდეს კორელაციის მრავლობითი კოეფიციენტის საშუალებით, რომელიც გამოითვლება შემდეგი ფორმულით:

$$r_{j:(1,2,\dots,p)} = \sqrt{1 - \frac{|R|}{R_{jj}}},$$

სადაც, $|R|$ – კორელაციური მატრიცის დეტერმინანტია, R_{jj} – კორელაციური მატრიცის r_{jj} ელემენტის ალგებრული დამატება.

კორელაციის მრავლობითი კოეფიციენტი დადებითი სიდიდეა და იცვლება ნულსა და ერთს შორის. მისი სარწმუნოების შეფასებისთვის განვიხილოთ $H_0: r_{j:(1,2,\dots,p)} = 0$ ნულოვანი ჰიპოთეზა. ამისათვის საჭიროა გამოითვალოს შემდეგი სტატისტიკა:

$$F = \frac{r_{j:(1,2,\dots,p)}^2 (n-p-1)}{(1-r_{j:(1,2,\dots,p)}^2) p},$$

სადაც n ამონარჩევის დანზომილებაა. F სიდიდეს გააჩნია ფიშერის განაწილება $\nu_1 = p$ და $\nu_2 = n - p - 1$ თავისუფლების ხარისხებით. თუ $F \geq F_{\alpha; \nu_1, \nu_2}$, მაშინ, ნულოვანი ჰიპოთეზა უარყოფილია ალტერნატიულის სასარგებლოდ, ხოლო თუ $F < F_{\alpha; \nu_1, \nu_2}$, მაშინ ნულოვანი ჰიპოთეზა მიიღება.

მრავლობითი კორელაციის კერძო მარტივ შემთხვევას წარმოადგენს დამოკიდებულება სამ X, Y, Z ცვლადს შორის. მაშინ კავშირი, მაგ. X -სა და დანარჩენ ორ Y და Z ცვლადებს შორის განისაზღვრება ფორმულით:

$$r_{x(yz)} = \sqrt{\frac{r_{xy}^2 + r_{xz}^2 - 2r_{xy}r_{xz}r_{yz}}{1 - r_{yz}^2}},$$

სადაც, r_{xy} , r_{xz} და r_{yz} პირსონის კორელაციის კოეფიციენტებია.

მაგალითი. ზემოთ განხილული მაგალითისთვის გვექნება:

$$r_{x(yz)} = \sqrt{\frac{0,865^2 + 0,853^2 - 2 \cdot 0,865 \cdot 0,863 \cdot 0,95}{1 - 0,95^2}} = \sqrt{\frac{0,0739}{0,0975}} = 0,871;$$

$$F = \frac{0,871^2(10 - 2 - 1)}{1(1 - 0,871^2)2} = \frac{5,31}{0,48} = 11,06; F_{0,05;2,7} = 4,739.$$

რადგან $11,06 > 4,729$, ამიტომ კავშირი X -სა და დანარჩენ Y და Z ცვლადებს შორის სარწმუნოა.

12.3. კორელაციის კოეფიციენტის განსაზღვრის არაპარამეტრული მეთოდები

პრაქტიკაში ხშირად ამონარჩევის განაწილების კანონი უცნობია ან განსხვავდება ნორმალურისგან. ამ შემთხვევაში პარამეტრული მეთოდების გამოყენება შეუძლებელია. გარდა ამისა, თუ მონაცემები თვისებრივი ხასიათისაა, მაშინ მათ შორის კავშირის დადგენა პარამეტრული მეთოდებით შეუძლებელია. განვიხილოთ ორ ცვლადს შორის კავშირის დადგენის ის არაპარამეტრული მეთოდები, რომლებიც უფრო ხშირად გამოიყენება პრაქტიკულ კვლევებში.

სპირმენის რანგული კორელაციის კოეფიციენტი. ზოგჯერ გვხვდება ისეთი შემთხვევები, როცა პარამეტრები რაოდენობრივ შეფასებებს არ ექვემდებარება. მაგალითად, ვთქვათ, საჭიროა შევაფასოთ თანაფარდობა მოსწავლეთა ჯგუფის მუსიკალურ და მათემატიკურ ნიჭს შორის. ამ შემთხვევაში „ნიჭის დონე“ არის ცვლადი სიდიდე იმ აზრით, რომ იგი იცვლება ერთი

ინდივიდუმიდან მეორემდე. მისი გაზომვა შესაძლებელია, თუ ყოველ მოსწავლეს დავუწეროთ ნიშანს. მაგრამ ასეთ მეთოდს გააჩნია არაობიექტურობა, რადგან სხვადასხვა გამომცდელს ერთი და იგივე მოსწავლე შეუძლია სხვადასხვანაირად შეაფასოს. სუბიექტურობის ელემენტი შეიძლება გამოირიცხოს, თუ მოსწავლეები რანჟირებულნი იქნებიან ნიჭის დონის მიხედვით და ყოველ მათგანს მივანიჭებთ რანგს.

რანგებს შორის კორელაცია უფრო ზუსტად ასახავს თანაფარდობას მუსიკალურ და მათემატიკურ ნიჭს შორის. ერთ-ერთ ყველაზე უფრო გავრცელებულ მაჩვენებელს წარმოადგენს სპირმენის რანგული კორელაციის კოეფიციენტი, რომელიც ამონარჩევების დამოუკიდებლად რანჟირების შემდეგ განისაზღვრება შემდეგნაირად:

$$r_s = 1 - \frac{6 \sum_{i=1}^n d_i^2}{n^3 - n},$$

სადაც, d არის X და Y შეუღლებულ მნიშვნელობათა რანგებს შორის სხვაობა, ე.ი. $d_i = R_{x_i} - R_{y_i}$. n - ამონარჩევის მოცულობაა. თუ რანგებს შორის სხვაობა $d_i = 0$, $i = 1, 2, \dots, n$, მაშინ $r_s = 1$.

როდესაც რანგულ მსკრივი გვხვდება ერთნაირი სიდიდის რანგები, მაშინ სპირმენის რანგული კორელაციის კოეფიციენტი, მიზანშეწონილია, გამოითვალოს შემდეგი ფორმულით:

$$r_s = 1 - \frac{6 \sum_{i=1}^n d_i^2}{(n^3 - n) - (T_x + T_y)},$$

სადაც,

$$T_x = \frac{1}{2} \sum (t_x^3 - t_x); \quad T_y = \frac{1}{2} \sum (t_y^3 - t_y);$$

t_x და t_y წარმოადგენენ ერთნაირი რანგების წევრთა რაოდენობებს.

სპირმენის რანგული კორელაციის სარწმუნოება მოწმდება ისევე, როგორც პირსონის კორელაციის კოეფიციენტისა. რანგული კორელაციის კოეფიციენტი იცვლება $-1 \leq r_s \leq 1$ და ის ძირითადად გამოიყენება იმ შემთხვევებში, როცა საჭიროა სწრაფად შეფასდეს პარამეტრებს შორის კავშირი, თუ პარამეტრების რან-

ჟირება შესაძლებელია და ამ პარამეტრებს არ გააჩნიათ ნორმა-
ლური განაწილება.

მაგალითი. მოცემულია 8 სტუდენტის შეფასებები მათე-
მატიკასა (x) და უცხო ენაში (y). გვინტერესებს, არსებობს თუ არა
კავშირი მათემატიკასა და უცხო ენას მოსწრებას შორის. შე-
ვადგინოთ შემდეგი ცხრილი:

№	x	y	R_x	R_y	d	d^2
1	5	4	1	2	-1	1
2	4	2	3	7	-4	16
3	4	5	3	1	2	4
4	4	3	3	4	-1	1
5	3	2	5,5	7	-1,5	2,25
6	3	3	5,5	4	1,5	2,25
7	2	2	7,5	7	0,5	0,25
8	2	3	7,5	4	3,5	12,25

$$\sum_{i=1}^8 d_i^2 = 39; \quad T_x = \frac{1}{2}[(3^3 - 3) + (2^3 - 2) + (2^3 - 2)] = 18;$$

$$T_y = \frac{1}{2}[(3^3 - 3) + (3^3 - 3)] = 24;$$

$$r_s = 1 - \frac{6 \cdot 39}{(8^3 - 8) - (18 + 24)} = 0,49; \quad t = 0,49 \sqrt{\frac{6}{1 - (0,49)^2}} = 1,38.$$

$t_{0,05;6} = 2,015$. რადგან $t < t_{0,05;6}$, ამიტომ მათემატიკასა და უცხო
ენას მოსწრებას შორის კავშირი არ არსებობს.

ასოციაციის კოეფიციენტი. ორი A და B თვისებრივ
მაჩვენებელს შორის დამოკიდებულება შეიძლება განისაზღვროს
ასოციაციის კოეფიციენტით ანუ ტეტრაქონული კავშირის
მაჩვენებლით. თუ მოცემულია (2×2) რიგის ცხრილი, რომელიც
შევსებულია A და B მაჩვენებლების არსებობის ან არარსებობის
(\bar{A}, \bar{B}) საფუძველზე, მაშინ ასოციაციის კოეფიციენტი გან-
ისაზღვრება ფორმულით:

$$r_A = \frac{|ad - bc| - 0,5n}{\sqrt{(a+b)(c+d)(a+c)(b+d)}},$$

	B	\bar{B}	
A	a	b	$a+b$
\bar{A}	c	d	$c+d$
	$a+c$	$b+d$	n

სადაც, a, b, c, d – (2×2) ცხრილის უჯრედებში მოხვედრის რაოდენობებია (სიბშირეებია), $n = a + b + c + d$.

ასოციაციის r_A კოეფიციენტი იცვლება $-1 \leq r_A \leq 1$ შუალედში და χ^2 კრიტერიუმთან არსებობს შემდეგი დამოკიდებულება:

$$r_A = \sqrt{\frac{\chi^2}{n}}.$$

ასოციაციის კოეფიციენტის სარწმუნოება განისაზღვრება $H_0: r_A = 0$ ნულოვანი ჰიპოთეზის შემოწმებით. ნულოვანი ჰიპოთეზა უარყოფილი იქნება, როცა $nr_A^2 \geq \chi_{\alpha, \nu}^2$. წინააღმდეგ შემთხვევაში, იგი მიიღება. $\chi_{\alpha, \nu}^2$ კრიტიკული წერტილი α მნიშვნელოვნების დონითა და $\nu = 1$ თავისუფლების ხარისხით აიღება χ^2 განაწილების ცხრილიდან.

r_A კოეფიციენტის სარწმუნოება შეიძლება შემოწმდეს აგრეთვე სტიუდენტის კრიტერიუმით, ისევე როგორც პირსონის კორელაციის კოეფიციენტის დროს.

მაგალითი. საჭიროა დავადგინოთ, არსებობს თუ არა კავშირი ქოლერის საწინააღმდეგო აცრასა და ქოლერით დაავადებას შორის. მონაცემები მოყვანილია შემდეგ ცხრილში:

	არ დაავადნენ	დაავადნენ	სულ
აცრით	276	3	279
აცრის გარეშე	473	86	559
ჯამი	749	89	838

$$r_A = \frac{276 \cdot 86 - 3 \cdot 473 - 0,5 \cdot 838}{\sqrt{279 \cdot 559 \cdot 749 \cdot 89}} = \frac{21898}{101963,31} = 0,215.$$

$$\chi_{0,05;1}^2 = 3,84$$

რადგან $838 \cdot (0,215)^2 = 38,74 > 3,84$, ამიტომ კავშირი ქოლერის საწინააღმდეგო აცრასა და ქოლერით დაავადებას შორის საწინააღმდეგოა.

ურთიერთშეუღლებულობის კოეფიციენტი. თუ თვისებრივ პარამეტრთა რაოდენობა ორზე მეტია, მაშინ გამოიყენება ურთიერთშეუღლებულობის კოეფიციენტი ანუ, როგორც მას უწოდებენ, პოლიქორიული კავშირის მაჩვენებელი, რომელიც პირსონმა შემოგვთავაზა და ჩუპროვმა გაუკეთა მოდერნიზაცია:

$$k = \sqrt{\frac{\varphi^2}{\varphi^2 + 1}} = \sqrt{\frac{\varphi^2}{(n_x - 1)(n_y - 1)}},$$

სადაც, $\varphi^2 = \left(\sum_{i=1}^m \frac{m_{xy}^2}{m_x m_y} \right) - 1$ ეწოდება პირსონის კონტინგენციის

კოეფიციენტი; m_{xy} – ცხრილის უჯრედების სიხშირეა; m_x, m_y – ცხრილის სტრიქონებისა და სვეტების სიხშირეების ჯამი; n_x, n_y – ცხრილის სტრიქონებისა და სვეტებში ჯგუფების რაოდენობა.

ურთიერთშეუღლებულობის კოეფიციენტი იცვლება ნულიდან ერთამდე. ნულოვანი ჰიპოთეზა $H_0: k = 0$ უარყოფილი იქნება, თუ

$$\chi^2 = n\varphi^2 \geq \chi_{\alpha;v}^2,$$

სადაც, n ამონარჩევის მოცულობაა, $v = (n_x - 1)(n_y - 1)$.

მაგალითი. შეისწავლეს დამოკიდებულება თმის ფერსა და თვალის ფერს შორის. მონაცემები მოყვანილია შემდეგ ცხრილში:

თვალის ფერი	თმის ფერი			სულ
	ქერა	ნაბლისფერი	წითური	
ცისფერი	170	80	5	255
ნაცრისფერი	70	152	8	230
თაფლისფერი	68	340	7	415
სულ	308	572	20	900

$$\phi^2 = \frac{170^2}{255 \cdot 308} + \frac{80^2}{255 \cdot 572} + \frac{5^2}{255 \cdot 20} + \frac{70^2}{230 \cdot 308} + \frac{152^2}{230 \cdot 572} + \frac{8^2}{230 \cdot 20} + \frac{68^2}{415 \cdot 308} + \frac{340^2}{415 \cdot 572} + \frac{7^2}{415 \cdot 20} - 1 = 1,205 - 1 = 0,205;$$

$$k = \sqrt{\frac{0,205}{\sqrt{(3-1)(3-1)}}} = 0,226; \chi_{0,01;4}^2 = 13,28; \nu = (3-1)(3-1) = 4;$$

$$\chi^2 = 900 \cdot 0,205 = 184,5.$$

რადგან $184,5 > 13,28$, ამიტომ ნულოვანი ჰიპოთეზა უარყოფილია, ე.ი. დამოკიდებულება თმის ფერსა და თვალის ფერს შორის სარწმუნოა.

12.4. კონკორდაციის კოეფიციენტი

პრაქტიკაში ხშირად იქმნება ისეთი სიტუაცია, როდესაც საჭირო ხდება ექსპერტთა ჯგუფის გამოყენება ამა თუ იმ საკითხის გადასაჭრელად (კონსილიუმის მოწვევა დიაგნოზის დასასმელად, სხვადასხვა ნორმატივების დადგენა, პრეპარატების შეფასება და სხვ.). აქედან გამომდინარე, სასურველია დავადგინოთ, რამდენად ემთხვევა ექსპერტთა აზრი ერთი და იმავე საკითხის განხილვას. თუ გვყავს მხოლოდ ორი ექსპერტი, მაშინ თანხმობის ზომად შეგვიძლია მივიღოთ სპირმენის რანგული კორელაციის კოეფიციენტის სიდიდე. მაგრამ როდესაც ექსპერტთა რაოდენობა დიდია, მაშინ სპირმენის რანგული კორელაციის კოეფიციენტის გამოყენება მიზანშეწონილი არ არის.

დავუშვათ, გვაქვს ობიექტების n რაოდენობა და გვყავს m ექსპერტი, რომლებიც აფასებენ ამ ობიექტებს და ახდენენ მათ რანჟირებას (უკეთესიდან უარესისკენ). რანჟირების შედეგად ვიღებთ ასეთ ცხრილს:

ობიექტი ექსპერტი	1	2	3	...	n
1	x_{11}	x_{12}	x_{13}	...	x_{1n}
2	x_{21}	x_{22}	x_{23}	...	x_{2n}
⋮	⋮	⋮	⋮	⋮	⋮
j	x_{j1}	x_{j2}	x_{j3}	...	x_{jn}
⋮	⋮	⋮	⋮	⋮	⋮
m	x_{m1}	x_{m2}	x_{m3}	...	x_{mn}

ექსპერტთა თანხმობის ზომის დასადგენად, შეგვიძლია გამოვიყენოთ კენდელის მიერ შემოთავაზებული თანხმობის ანუ კონკორდაციის კოეფიციენტი W , რომელიც ასე განისაზღვრება:

$$W = \frac{12 \sum_{i=1}^n S_i^2}{m^2 (n^3 - n)},$$

სადაც, S არის სხვაობა ობიექტების რანგების ჯამსა და რანგების საერთო საშუალო არითმეტიკულს შორის, ე.ი.

$$S_i = R_i - \bar{R}, \quad R_i = \sum_{k=1}^m x_{ki}, \quad i = 1, 2, \dots, n.$$

თუ ექსპერტების რანჟირებულ მნკრივში გვხვდება ერთი და იგივე რანგის მნიშვნელობა (ეს ის შემთხვევაა, როცა ექსპერტი ვერ ანიჭებს უპირატესობას), მაშინ კონკორდაციის კოეფიციენტი გამოითვლება შემდეგი ფორმულით:

$$W = \frac{\sum_{i=1}^n S_i^2}{\frac{1}{12} m^2 (n^3 - n) - m \sum_{j=1}^m T_j},$$

სადაც, $T_j = \frac{1}{12} \sum_j (t_j^3 - t_j)$, t_j - j -ურ მწკრივში ერთნაირი

რანგების რაოდენობაა.

ზოგადად, $0 \leq W \leq 1$. თუ ექსპერტთა შეფასებები ერთმანეთს მთლიანად ემთხვევა, მაშინ $W = 1$, ხოლო თუ მათი აზრები მკვეთრად განსხვავდებიან ერთმანეთისგან, მაშინ $W = 0$.

კონკორდაციის კოეფიციენტის სარწმუნოების დასადგენად უნდა შევამოწმოთ $H_0: W = 0$ ნულოვანი ჰიპოთეზა. ასეთი ნულოვანი ჰიპოთეზის შესამოწმებლად უნდა გამოვთვალოთ შემდეგი სტატისტიკა: $\chi^2 = Wm(n-1)$, რომელსაც გააჩნია χ^2 განაწილება $\nu = n - 1$ თავისუფლების ხარისხით.

თუ $\chi^2 < \chi_{\alpha, \nu}^2$, მაშინ ნულოვანი ჰიპოთეზა მიიღება, ე.ი. ექსპერტთა აზრები განსხვავდებიან ერთმანეთისაგან, ხოლო როცა $\chi^2 \geq \chi_{\alpha, \nu}^2$, მაშინ ექსპერტთა აზრები ერთმანეთს ემთხვევა.

მაგალითი. 6 ექსპერტი აფასებს 4 ფარმაცევტული ფირმის მიერ გამოშვებულ პრეპარატს. შედეგები მოყვანილია შემდეგ ცხრილში:

ფირმები ექსპერტები	1	2	3	4	სულ
1	1	3	2	4	
2	2	1	4	3	
3	1	3	2	4	
4	3	2	1	4	
5	1	2	4	3	
6	2	3	1	4	
R	10	14	14	22	60
S	-5	-1	1	7	
S^2	25	1	1	49	76

$$\bar{R} = \frac{60}{4} = 15; \quad W = \frac{12 \cdot 76}{36(4^3 - 4)} = 0,42; \quad \chi^2 = 0,42 \cdot 6 \cdot 3 = 7,56;$$

$$\chi_{0,05;5}^2 = 11,07$$

რადგან $7,56 < 11,07$, ამიტომ ნულოვანი ჰიპოთეზა მიიღება, ე.ი. ექსპერტთა აზრები განსხვავდება ერთმანეთისგან.

13. რეგრესიული ანალიზის საფუძვლები

13.1. რეგრესიული ანალიზის არსი

რეგრესიულ ანალიზში განიხილება კავშირი ერთ დამოკიდებულ y ცვლადსა და ერთ ან რამდენიმე დამოუკიდებელ x_1, x_2, \dots, x_n ცვლადებს შორის, რომლებიც ურთიერთდამოუკიდებელი უნდა იყვნენ. ეს კავშირი, კორელაციური ანალიზისაგან განსხვავებით, წარმოდგენილია რეგრესიის განტოლების $\hat{y} = f(x_1, x_2, \dots, x_n)$ სახით, რომელიც აკავშირებს დამოკიდებულ ცვლადს დამოუკიდებელ ცვლადებთან. დამოუკიდებელ ცვლადებს ზოგჯერ **პრედიქტორებს** ან **რეგრესორებს** უწოდებენ.

რეგრესიის განტოლება წარმოადგენს ყველაზე უფრო გავრცელებულ სტატისტიკურ მოდელს, რადგან, გარდა დამოკიდებულების აღწერისა, იგი შეიძლება გამოვიყენოთ ინფორმაციული პარამეტრების შესარჩევად და პროგნოზირების ამოცანის გადასაწყვეტად.

რეგრესიის განტოლების შედგენა გულისხმობს ორი ძირითადი ამოცანის გადაწყვეტას. პირველი მდგომარეობს ისეთი დამოუკიდებელი ცვლადების შერჩევაში, რომლებიც მნიშვნელოვნად მოქმედებენ დამოკიდებულ ცვლადზე და მეორე – რეგრესიის განტოლების სახის შერჩევაში.

განიხილოთ ყველაზე მარტივი შემთხვევა, როდესაც ჩვენ გვინდა აღვწეროთ ორ X და Y ცვლადებს შორის დამოკიდებულების ფუნქცია. დავუშვათ, რომ თეორიულად ამ ცვლადებს შორის არსებობს მარტივი წრფივი დამოკიდებულება

$$y = \alpha + \beta x, \quad (13.1)$$

სადაც, α და β უცნობი მუდმივი პარამეტრებია (კოეფიციენტებია), x – დამოუკიდებელი, ხოლო y – დამოკიდებული ცვლადებია. პრაქტიკულად, y და x შორის დამოკიდებულება ცალსახად არაა მკაცრი, რადგან y -ზე მოქმედებს სხვადასხვა

გაუთვალისწინებელი ფაქტორები, შეშფოთებები, ხმაური და სხვა. ამიტომ (13.1) განტოლება შეიძლება ასე წარმოვადგინოთ:

$$y = \alpha + \beta x + \varepsilon, \quad (13.2)$$

სადაც, ε ცვლადი სიდიდეა, რომელიც ახასიათებს თეორიული ნირიდან გადახრას.

იმისათვის, რომ რეგრესიის (13.2) განტოლება შერჩეული იყოს სწორად, ε ცდომლება უნდა აკმაყოფილებდეს შემდეგ პირობებს:

1. ε სიდიდე უნდა იყოს შემთხვევითი;
2. ε ცდომლების მათემატიკური ლოდინი უნდა იყოს ნულის ტოლი;
3. ε ცდომლების დისპერსია უნდა იყოს მუდმივი სიდიდე;
4. ε ცდომლების შესაძლო მნიშვნელობები უნდა იყვნენ ერთმანეთის მიმართ დამოუკიდებელი (არაკორელირებულები).

ამრიგად, რეგრესიის განტოლების შედგენისას, მიიღება ჰიპოთეზა იმის შესახებ, რომ ყოველი i -ური დაკვირვებისთვის სრულდება შემდეგი დამოკიდებულება:

$$y_i = \alpha + \beta x_i + \varepsilon_i, \quad i=1,2,\dots,n.$$

შემთხვევითი ε ცდომლების მათემატიკური ლოდინი, დისპერსია და კოვარიაცია ტოლია:

$$M(\varepsilon_i) = 0; \quad M(\varepsilon_i \varepsilon_j) = \begin{cases} 0, & \text{როცა } i \neq j \\ s^2, & \text{როცა } i = j \end{cases} \quad i, j = 1, 2, \dots, n.$$

აქედან გამომდინარე, $M(\varepsilon_i, \varepsilon_j)$ არის შემთხვევითი ცდომილების დისპერსია, ხოლო $M(\varepsilon_i, \varepsilon_j)$ – კოვარიაცია.

ამრიგად, თუ მოცემულია დამოუკიდებელი ცვლადის x_i და მისი შესაბამისი დამოკიდებული y_i ცვლადების მნიშვნელობები, მაშინ ამოცანა მდგომარეობს α და β პარამეტრების მოძებნაში. ამ პარამეტრების რეალური მნიშვნელობების განსაზღვრა შეუძლებელია, რადგან ჩვენ საქმე გვაქვს სასრული რაოდენობის ამონარჩევთან. ამიტომ საჭიროა მოიძებნოს ამ პარამეტრების შეფასებები. თუ შეფასებებს აღვნიშნავთ a და b სიმბოლოებით, მაშინ გვექნება შემდეგი სახის წრფივი რეგრესიის განტოლება:

$$\hat{y} = a + bx,$$

ხოლო ზოგადად, თუ საქმე გვაქვს მრავალგანზომილებიან რეგრესიულ ანალიზთან, მაშინ

$$\hat{y} = a_1x_1 + a_2x_2 + \dots + a_nx_n.$$

ამრიგად, რეგრესიული ანალიზის მეთოდი ზოგადად მდგომარეობს დამოკიდებულ y ცვლადსა და დამოუკიდებელ X_1, X_2, \dots, X_n ცვლადებს შორის კავშირის ფორმის დადგენაში. ეს კავშირი აღინერება

$$y = f(X_1, X_2, \dots, X_n; a_1, a_2, \dots, a_n) + \varepsilon$$

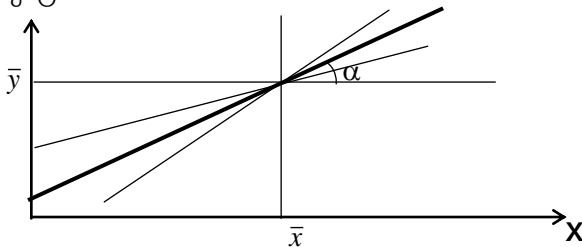
მათემატიკური მოდელის სახით, სადაც, a_i უცნობი (საძიებელი) პარამეტრებია, ε – შემთხვევითი ცდომლება.

თუ რეგრესიის f ფუნქცია წრფივია $a_i, i=1, 2, \dots, n$ პარამეტრების მიმართ (მაგრამ არ არის აუცილებელი X_1, X_2, \dots, X_n დამოუკიდებელი ცვლადების მიმართ), მაშინ ასეთ რეგრესიულ ანალიზს წრფივი ეწოდება. თუ რეგრესიის ფუნქცია არაწრფივია a_i პარამეტრების მიმართ, მაშინ საქმე გვაქვს არაწრფივ რეგრესიულ ანალიზთან.

წრფივი რეგრესიის განტოლების გეომეტრიული ინტერპრეტაციისთვის განვიხილოთ ორ ცვლადს შორის დამოკიდებულება

$$\hat{y} = a + bx,$$

სადაც, a თავისუფალი წევრია და გვიჩვენებს რეგრესიის განტოლების y ღერძთან გადაკვეთის წერტილს. b პარამეტრი გვიჩვენებს რეგრესიის წირის დახრილობას x ღერძის მიმართ და იგი ტოლია $b = tg\alpha$. ანალიზურ გეომეტრიაში ამ პარამეტრს უწოდებენ კუთხის კოეფიციენტს, ხოლო ბიომეტრიაში – რეგრესიის კოეფიციენტს.



გრაფიკზე წარმოდგენილი რეგრესიის წირი AB , რომელიც გაივლის $M(\bar{x}, \bar{y})$ წერტილში, შეესაბამება X და Y ცვლადების ფუნქციონალურ დამოკიდებულებას, როდესაც $r_{xy} = 1$. რაც უფრო ძლიერია კავშირი X და Y ცვლადებს შორის, მით უფრო

უახლოვდება რეგრესიის წირი AB წირს და პირიქით, რაც უფრო სუსტია ეს კავშირი, მით უფრო შორდება AB წირს. თუ კავშირი ცვლადებს შორის არ არსებობს, მაშინ რეგრესიის განტოლების წირი X ღერძთან ადგენს 90° კუთხეს.

რადგან რეგრესია განაპირობებს ცვლადების ორმხრივ კორელაციურ კავშირს, ამიტომ რეგრესიის განტოლება შეიძლება ასე ჩავენროთ:

$$\hat{y}_x = a_{yx} + b_{yx}x \text{ და } \hat{x}_y = a_{xy} + b_{xy}y.$$

აქვე უნდა შევნიშნოთ, რომ არსებობს კავშირი რეგრესიის კოეფიციენტსა და კორელაციის კოეფიციენტს შორის

$$r_{xy} = \sqrt{b_{yx} \cdot b_{xy}},$$

სადაც, რეგრესიის კოეფიციენტი შეიძლება გამოვთვალოთ შემდეგნაირად:

$$b_{yx} = \frac{\sum_{i=1}^n (y_i - \bar{y})(x_i - \bar{x})}{\sum_{i=1}^n (x_i - \bar{x})^2} \text{ ან } b_{xy} = \frac{\sum_{i=1}^n (y_i - \bar{y})(x_i - \bar{x})}{\sum_{i=1}^n (y_i - \bar{y})^2}.$$

თუ კორელაციის კოეფიციენტი ცნობილია, მაშინ რეგრესიის კოეფიციენტი შეგვიძლია ასე განვსაზღვროთ:

$$b_{yx} = r_{xy} \sqrt{\frac{\sum_{i=1}^n (y_i - \bar{y})^2}{\sum_{i=1}^n (x_i - \bar{x})^2}}, \quad b_{xy} = r_{xy} \sqrt{\frac{\sum_{i=1}^n (x_i - \bar{x})^2}{\sum_{i=1}^n (y_i - \bar{y})^2}},$$

ან

$$b_{yx} = r_{xy} \frac{\sigma_y}{\sigma_x}; \quad b_{xy} = r_{xy} \frac{\sigma_x}{\sigma_y},$$

სადაც, σ_x და σ_y საშუალო კვადრატული გადახრებია.

პრაქტიკულად, y_i ცვლადის რეალური მნიშვნელობები არ დაემთხვევა რეგრესიის განტოლებით გამოთვლილ \hat{y}_i მნიშვნელობებს, რადგან თვით რეგრესიის განტოლება აღწერს ამ დამოკიდებულების საშუალო მნიშვნელობებს.

13.2. უმცირეს კვადრატთა მეთოდი

თუ Y დამოკიდებული და X_1, X_2, \dots, X_n ურთიერთდამოუკიდებელი ცვლადები ნორმალურად არიან განაწილებულნი, მაშინ უცნობი a_1, a_2, \dots, a_n პარამეტრების შეფასებები შეიძლება მოიძებნოს უმცირეს კვადრატთა მეთოდით, რომელიც, სხვა მეთოდებთან შედარებით, იძლევა ალბათური აზრით საუკეთესო შედეგებს. უმცირეს კვადრატთა მეთოდის არსი მდგომარეობს შემდეგში: რეგრესიის f ფუნქცია ისე უნდა შევარჩიოთ, რომ მოცემული y_i მნიშვნელობების რეგრესიის განტოლებით გამოთვლილი \hat{y}_i მნიშვნელობებთან გადახრის $(y - \hat{y}_i)$ კვადრატების ჯამი იყოს მინიმალური, ე.ი.

$$Q = \sum_{i=1}^n [y_i - f(x_i, a_1, a_2, \dots, a_n)]^2 = \sum_{i=1}^n (y_i - \hat{y}_i)^2 = \sum_{i=1}^n e_i^2 = \min, \quad (13.3)$$

სადაც, $e_i = y_i - \hat{y}_i$ უწოდებენ ნარჩენ მნიშვნელობებს, ანუ ნაშთებს.

მოვძებნოთ a_1, a_2, \dots, a_n უცნობი კოეფიციენტები, რომლებიც (13.3) გამოსახულების მარცხენა მხარეს გადააქცევენ მინიმუმად. ამისათვის ეს გამოსახულება გავანარმოოთ a_1, a_2, \dots, a_n -ით და გავუტოლოთ ნულს, მაშინ გვექნება:

$$\left. \begin{aligned} \frac{\partial Q}{\partial a_1} &= \sum_{i=1}^n [y_i - f(x_i, a_1, a_2, \dots, a_n)] \left(\frac{\partial f}{\partial a_1} \right)_i = 0 \\ \frac{\partial Q}{\partial a_2} &= \sum_{i=1}^n [y_i - f(x_i, a_1, a_2, \dots, a_n)] \left(\frac{\partial f}{\partial a_2} \right)_i = 0 \\ &\dots\dots \\ \frac{\partial Q}{\partial a_n} &= \sum_{i=1}^n [y_i - f(x_i, a_1, a_2, \dots, a_n)] \left(\frac{\partial f}{\partial a_n} \right)_i = 0 \end{aligned} \right\}, \quad (13.4)$$

სადაც, $\left(\frac{\partial f}{\partial a_j} \right)_i = f_{a_j}(x_i, a_1, a_2, \dots, a_n) - f$ ფუნქციის კერძო ნარ-

მოებულებაა a_j პარამეტრით x_i წერტილში. (13.4) სისტემა, რომელსაც **ნორმალურ განტოლებათა სისტემას** უწოდებენ, შეიცავს იმდენივე განტოლებას, რამდენი უცნობიცაა. აქვე უნდა აღინიშნოს, რომ (13.4) სისტემის ამოხსნა a_j პარამეტრების გან-

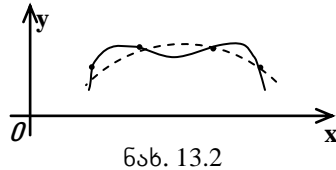
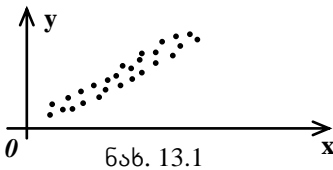
საზღვრისათვის ზოგადი სახით არ შეიძლება, ამისათვის აუცილებელია რეგრესიის ფუნქციის კონკრეტული სახე.

უმცირეს კვადრატთა მეთოდს გააჩნია შემდეგი დადებითი თვისებები:

1. ამ მეთოდით შეფასებული პარამეტრები გადაუადგილებადია;
2. პარამეტრთა შეფასებები საფუძვლიანია;
3. პარამეტრთა შეფასებები ეფექტურია იმ გაგებით, რომ მათ გააჩნიათ შესაძლო მინიმალური დისპერსია, სხვა მეთოდით გამოთვლილ დისპერსიასთან შედარებით.

13.3. წრფივი რეგრესია

დავუშვათ, რომ ცდის შედეგად მივიღეთ ექსპერიმენტალურ წერტილთა ერთობლიობა და ავაგეთ y -ის x -თან დამოკიდებულების გრაფიკი (ნახ. 13.1).



ცნობილია, რომ ნებისმიერ (x_i, y_i) კოორდინატებიან n წერტილზე შესაძლებელია გავავლოთ მრუდი, რომელიც ანალიზურად გამოისახება $(n - 1)$ ხარისხის პოლინომით ისე, რომ ზუსტად გაიაროს თითოეულ წერტილზე. საკითხის ასეთი გადაწყვეტა, ჩვეულებრივ, არ ითვლება დამაკმაყოფილებლად, რადგან რეგრესიის განტოლება გამოდის მაღალი რიგის და, რაც მთავარია, ასეთი რეგრესიის განტოლება პროცესს აღწერს დამახინჯებით (ნახ. 13.2). ამიტომ, მიზანშეწონილია, რომ ექსპერიმენტალური წერტილები ისე დავამუშავოთ, რომ რაც შეიძლება ზუსტად აღიწეროს x -სა და y -ს შორის დამოკიდებულების ტენდენცია.

ნახ. 13.1-ზე წარმოდგენილი წერტილები მიგვანიშნებს, რომ y -სა და x -ს შორის დამოკიდებულება შეიძლება გამოვსახოთ წრფივი განტოლების სახით. მაშინ რეგრესიის განტოლებას ექნება შემდეგი სახე:

$$\hat{y} = a + bx.$$

უმცირეს კვადრატთა მეთოდის საშუალებით შევადგათ a და b პარამეტრები. როგორც ვიცით,

$$Q = \sum_{i=1}^n (y_i - \hat{y}_i)^2 = \sum_{i=1}^n [y_i - (a + bx_i)]^2 = \min.$$

გავანარმოთ ეს გამოსახულება a -თი და b -თი, მაშინ გვექნება:

$$\left. \begin{aligned} \frac{\partial Q}{\partial a} &= -2 \sum_{i=1}^n (y_i - a - bx_i) = 0 \\ \frac{\partial Q}{\partial b} &= -2 \sum_{i=1}^n (y_i - a - bx_i) x_i = 0 \end{aligned} \right\}.$$

თუ მოვახდენთ ელემენტარულ გარდაქმნებს, მაშინ მივიღებთ შემდეგი სახის ნორმალურ განტოლებათა სისტემას:

$$\left. \begin{aligned} \sum_{i=1}^n y_i &= an + b \sum_{i=1}^n x_i \\ \sum_{i=1}^n x_i y_i &= a \sum_{i=1}^n x_i + b \sum_{i=1}^n x_i^2 \end{aligned} \right\},$$

რომლის ამოხსნის შემდეგ მივიღებთ:

$$b = \frac{\sum_{i=1}^n x_i y_i - \frac{1}{n} \sum_{i=1}^n x_i \cdot \sum_{i=1}^n y_i}{\sum_{i=1}^n x_i^2 - \frac{1}{n} \left(\sum_{i=1}^n x_i \right)^2},$$

$$a = \frac{1}{n} \sum_{i=1}^n y_i - b \frac{1}{n} \sum_{i=1}^n x_i = \bar{y} - b\bar{x}.$$

მას შემდეგ, რაც რეგრესიის განტოლება შერჩეულია, საჭიროა შევადგათ როგორც რეგრესიის განტოლება, ასევე რეგრესიის განტოლების კოეფიციენტებიც. ყველაზე მნიშვნელოვან მომენტს რეგრესიულ ანალიზში წარმოადგენს რეგრესიის განტოლების ვარგისიანობის შეფასება, ანუ დადგენა, შეესაბამება თუ არა შერჩეული განტოლება ექსპერიმენტალურ მონაცემებს. განტოლების ვარგისიანობა, ანუ ადეკვატურობა მონმდება ფიშერის კრიტერიუმით

$$F = \frac{\sigma_y^2}{\sigma_{\text{ნარ}}^2},$$

სადაც, $\sigma_{\text{ნარ}}^2 = \frac{1}{n-2} \sum_{i=1}^n (y_i - \hat{y}_i)^2$ წარმოადგენს ნარჩენ დისპერ-

სიას; $\sigma_y^2 = \frac{1}{n-2} \sum_{i=1}^n (\hat{y}_i - \bar{y})^2$.

α მნიშვნელოვნების დონითა და $v_1 = 1$, $v_2 = n - 2$ თავისუფლების ხარისხებით ფიშერის განაწილების ცხრილიდან მოიძებნება $F_{\alpha; v_1, v_2}$ კრიტიკული წერტილი. თუ $F \geq F_{\alpha; v_1, v_2}$, მაშინ რეგრესიის განტოლება ითვლება ადეკვატურად, ანუ ის კარგად აღწერს ორ ცვლადს შორის დამოკიდებულებას. თუ $F < F_{\alpha; v_1, v_2}$, მაშინ რეგრესიის განტოლება არაადეკვატურია და საჭიროა მისი შეცვლა.

რეგრესიის a და b კოეფიციენტების ნულთან ტოლობის ჰიპოთეზის შესამოწმებლად განვიხილოთ სტატისტიკა:

$$t_a = \frac{a}{\sigma_a}, \quad t_b = \frac{b}{\sigma_b},$$

სადაც,

$$\sigma_a = \frac{\sigma_{\text{ნარ}}}{\sqrt{n-2}}, \quad \sigma_b = \frac{\sigma_{\text{ნარ}}}{\sigma_x \sqrt{n-2}}, \quad \sigma_x^2 = \frac{1}{n-1} \left[\sum_{i=1}^n x_i^2 - \frac{1}{n} \left(\sum_{i=1}^n x_i \right)^2 \right].$$

t_a სტატისტიკას გააჩნია სტიუდენტის განაწილება $v = n - 2$ თავისუფლების ხარისხით. α მნიშვნელოვნების დონითა და v სიდიდით სტიუდენტის განაწილების ცხრილიდან მოიძებნება $t_{\alpha; v}$ კრიტიკული წერტილი. თუ $|t_a| \geq t_{\alpha; v}$, მაშინ ნულოვანი ჰიპოთეზა უარყოფილია და რეგრესიის a კოეფიციენტი სარწმუნოა. თუ $|t_a| < t_{\alpha; v}$, მაშინ ნულოვანი ჰიპოთეზა მიიღება, ე.ი. რეგრესიის კოეფიციენტის სიდიდე უნდა ჩაითვალოს უმნიშვნელოდ, რომელიც შემთხვევით განსხვავდება ნულისაგან მცირე ამონარჩევის ან სხვა მიზეზის გამო. ანალოგიურად მოწმდება b კოეფიციენტიც.

მაგალითი. x და y ცვლადების გაზომვების შედეგები მოცემულია ქვემოთ მოყვანილ ცხრილში. შევარჩიოთ $\hat{y} = a + bx$ სახის რეგრესიის განტოლება.

№	x_i	y_i	x_i^2	y_i^2	$x_i y_i$	\hat{y}_i
1	32	20	1024	400	640	24,43
2	30	24	900	576	720	23,34
3	36	28	1296	784	1008	26,60
4	40	30	1600	900	1200	28,78
5	44	31	1681	961	1271	29,32
6	47	33	2209	1089	1551	32,58
7	56	34	3136	1156	1904	37,47
8	54	37	2916	1369	1994	36,38
9	60	38	3600	1444	2280	39,65
10	55	40	3025	1600	2200	36,93
11	61	41	3721	1681	2501	40,19
12	67	43	4469	1849	2881	43,45
13	69	45	4761	2025	3105	44,54
14	76	48	5576	2304	3648	48,34
Σ	724	492	40134	18138	26907	

ცხრილის გაგრძელება

№	$(y_i - \hat{y}_i)$	$(y_i - \hat{y}_i)^2$	$(\hat{y}_i - \bar{y})$	$(\hat{y}_i - \bar{y})^2$
1	-4,43	18,84	-10,71	114,70
2	-0,66	0,44	-11,80	739,24
3	1,4	1,96	-8,54	72,93
4	1,22	1,49	-6,36	40,45
5	1,68	2,82	-5,82	33,87
6	0,42	0,18	-2,56	6,55
7	-3,47	12,04	2,33	5,43
8	0,62	0,38	1,24	1,54
9	-1,65	2,72	4,51	20,34
10	3,07	9,43	1,79	3,20
11	0,81	0,66	5,05	25,50
12	-0,45	0,20	8,31	69,06
13	0,46	0,21	9,40	88,36
14	-0,34	0,12	13,20	154,24
Σ		51,41		1375,27

$$\bar{x} = 51,71; \quad \bar{y} = 35,14;$$

$$b = \frac{26907 - \frac{1}{14} \cdot 492 \cdot 724}{40134 - \frac{1}{14} \cdot 724 \cdot 724} = 0,5435; \quad a = 35,14 - 0,5535 \cdot 51,71 = 7,0356;$$

$\hat{y} = 7,0356 + 0,5435x$. შევამოწმოთ რეგრესიის განტოლება:

$$\sigma_{\text{ვარ}}^2 = \frac{51,41}{12} = 4,28; \quad \sigma_y^2 = \frac{1375,27}{12} = 114,64; \quad F = \frac{114,64}{4,28} = 26,78;$$

$$F_{0,05;1,12} = 4,75.$$

რადგან $F > F_{0,05;1,12}$, რეგრესიის განტოლება ადეკვატურია. შევამოწმოთ რეგრესიის კოეფიციენტები:

$$\sigma_b = \frac{2,07}{14,39 \cdot 3,46} = 0,042; \quad \sigma_a = \frac{2,07}{3,46} = 0,598;$$

$$t_b = \frac{0,5435}{0,042} = 12,94; \quad t_a = \frac{7,0356}{0,598} = 11,77; \quad t_{0,05;12} = 1,78.$$

რადგან t_a და t_b მეტია 1,78-ზე, ამიტომ რეგრესიის a და b კოეფიციენტები სარწმუნოა.

13.4. წრფივი რეგრესიის განტოლების ნდობის ინტერვალი

განვიხილოთ \hat{y} მნიშვნელობისთვის ნდობის ინტერვალის განსაზღვრის საკითხი, ანუ ისეთი ინტერვალისა, სადაც მოცემული ალბათობით დამოკიდებული y პარამეტრის რეალური მნიშვნელობა იმყოფება.

ნდობის ინტერვალის განსაზღვრისთვის ვიპოვოთ \hat{y} მნიშვნელობის დისპერსია, რომელიც შედგება a და b პარამეტრების დისპერსიების ჯამისგან. ნორმალურ განტოლებათა სისტემიდან გვაქვს $a = \bar{y} - b\bar{x}$. ჩავსვათ ეს მნიშვნელობა რეგრესიის განტოლებაში, მაშინ

$$\hat{y} = a + bx = \bar{y} - b\bar{x} + bx = \bar{y} + b(x - \bar{x})$$

და დისპერსია ტოლია:

$$\sigma_{\hat{y}}^2 = \frac{\sigma_{\text{ნარ}}^2}{n} + \frac{\sigma_{\text{ნარ}}^2}{\sum_{i=1}^n (x_i - \bar{x})^2} (x_p - \bar{x})^2 = \sigma_{\text{ნარ}}^2 \left[\frac{1}{n} + \frac{(x_p - \bar{x})^2}{\sum_{i=1}^n (x_i - \bar{x})^2} \right], \quad (13.5)$$

სადაც, $x_p - \bar{x}$ ცვლადის ის მნიშვნელობაა, რომლის შესაბამისი y პარამეტრისთვის გვაინტერესებს ნდობის ინტერვალის სიდიდე. (13.5) გამოსახულებიდან გამომდინარეობს, რომ $\sigma_{\hat{y}}^2$ მინიმალურ მნიშვნელობას იღებს, როცა $(x_p - \bar{x}) = 0$. ამ შემთხვევაში

$\sigma_{\hat{y}}^2 = \frac{\sigma_{\text{ნარ}}^2}{n}$. ამრიგად, თუ ვიცით \hat{y} მაჩვენებლის დისპერსია, მაშინ ადვილია განისაზღვროს ნდობის ინტერვალი

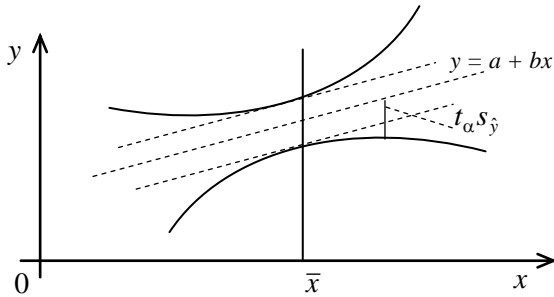
$$\hat{y} \pm t_{\alpha, \nu} \sigma_{\hat{y}},$$

სადაც, $t_{\alpha, \nu}$ სიდიდე განისაზღვრება სტიუდენტის განაწილების ცხრილიდან α მნიშვნელოვნების დონითა და $\nu = n - 2$ თავისუფლების ხარისხით.

აქვე უნდა აღვნიშნოთ, რომ ნდობის ეს ინტერვალი განსაზღვრავს რეგრესიის წირის (ე.ი. \hat{y} საშუალო მნიშვნელობის) მდებარეობას და არ არის დამოკიდებული ცვლადის თითოეული შესაძლო მნიშვნელობების მდებარეობაზე, რომლებიც რაღაც მანძილით გადახრილნი არიან რეგრესიის წირიდან. ამიტომ, თუ ჩვენ გვინდა y დამოკიდებული ცვლადის ცალკეული მნიშვნელობისათვის განისაზღვროს ნდობის ინტერვალი, მაშინ (13.5) ფორმულას უნდა დაემატოს კიდევ ერთი ნარჩენი დისპერსიის მნიშვნელობა, ე.ი. $y_i = a + bx_i + e_i$ განტოლებას შეესაბამება:

$$\sigma_p^2 = \frac{\sigma_{\text{ნარ}}^2}{n} + \frac{\sigma_{\text{ნარ}}^2}{\sum_{i=1}^n (x_i - \bar{x})^2} (x_p - \bar{x})^2 + \sigma_{\text{ნარ}}^2 = \sigma_{\text{ნარ}}^2 \left[1 + \frac{1}{n} + \frac{(x_p - \bar{x})^2}{\sum_{i=1}^n (x_i - \bar{x})^2} \right].$$

წრფივი რეგრესიის ნდობის ინტერვალი გრაფიკულად ასე შეიძლება წარმოვადგინოთ:



13.5. არანრფივი რეგრესია

მოვლენათა შორის ყოველთვის არ არსებობს ნრფივი დამოკიდებულება და ამიტომ ამ შემთხვევაში ნრფივი რეგრესიული მოდელის გამოყენება შეუძლებელია. არანრფივი დამოკიდებულების შემთხვევაში, რეგრესიის განტოლების შესადგენად უნდა გამოვიყენოთ არანრფივი ფუნქციები.

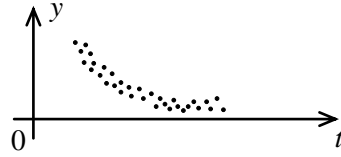
ანსხვავებენ არანრფივი რეგრესიის ორ სახეს. პირველს მიეკუთვნება ისეთი რეგრესიის განტოლებები, რომლებიც არანრფივია ცვლადების მიმართ, მაგრამ ნრფივია შესაფასებელი პარამეტრების (კოეფიციენტების) მიმართ. ზოგჯერ, ასეთ არანრფივ რეგრესიის განტოლებებს კვაზინრფივ რეგრესიას უწოდებენ. ასეთი ტიპის რეგრესიის განტოლების კოეფიციენტები მოიძებნება უმცირეს კვადრატთა მეთოდით.

მეორე ტიპის არანრფივ რეგრესიას მიეკუთვნებიან ისეთი განტოლებები, რომლებთაც ახასიათებთ არანრფივობა შესაფასებელი პარამეტრების ანუ კოეფიციენტების მიმართ. ამ შემთხვევაში უმცირეს კვადრატთა მეთოდის გამოყენება არ შეიძლება, ე.ი. უნდა გამოვიყენოთ არანრფივი უმცირეს კვადრატთა მეთოდი. ჩვენ განვიხილავთ მხოლოდ პირველი ტიპის რეგრესიის განტოლებებს, რადგან ისინი უფრო ხშირად გვხვდება პრაქტიკულ კვლევებში.

დავუშვათ, რომ მივიღეთ ექსპერიმენტალურ ნერტილთა ერთობლიობა (x_i, y_i) , $i=1, 2, \dots, n$ და ავაგეთ დამოკიდებულება $y=f(x)$.



ნახ. 13.3



ნახ. 13.4

13.3 ნახაზზე წარმოდგენილი წერტილები მიგვანიშნებს, რომ y და x შორის დამოკიდებულება შეიძლება გამოვსახოთ მეორე რიგის პარაბოლის (პარაბოლის) სახით

$$\hat{y} = a + bx + cx^2,$$

ხოლო 13.4 ნახაზზე წარმოდგენილი – $\hat{y} = ae^{-bx}$ მაჩვენებლიანი ან

$\hat{y} = a + \frac{b}{x}$ პირველი რიგის ჰიპერბოლის სახით. გარდა ამისა,

დამოკიდებულება შეიძლება იყოს ლოგარითმული, ხარისხოვანი, ტრიგონომეტრიული ფუნქციები და სხვა. მაგალითისთვის განვიხილოთ პარაბოლური, მაჩვენებლიანი და ჰიპერბოლური დამოკიდებულებები.

მეორე რიგის პარაბოლის a , b და c კოეფიციენტები განისაზღვრებიან შემდეგი ნორმალურ განტოლებათა სისტემის ამოხსნით:

$$\left. \begin{aligned} \sum_i y_i &= an + b \sum_i x_i + c \sum_i x_i^2 \\ \sum_i y_i x_i &= a \sum_i x_i + b \sum_i x_i^2 + c \sum_i x_i^3 \\ \sum_i y_i x_i^2 &= a \sum_i x_i^2 + b \sum_i x_i^3 + c \sum_i x_i^4 \end{aligned} \right\}. \quad (13.6)$$

(13.6) სისტემა ადვილად ამოიხსნება, თუ მოვახდენთ დამოუკიდებელი ცვლადის ცენტრირებას, ანუ თუ განვიხილავთ x_i ცვლადების სხვაობას მათ საშუალო არითმეტიკულთან, ე.ი.

$x_i^* = x_i - \bar{x}$. მაშინ გვექნება:

$$a = \frac{1}{D} \left(\sum_i y_i \sum_i (x_i^*)^4 - \sum_i (x_i^*)^2 \sum_i y_i (x_i^*)^2 \right), b = \frac{\sum_i y_i x_i^*}{\sum_i (x_i^*)^2},$$

$$c = \frac{1}{D} \left(n \sum_i y_i (x_i^*)^2 - \sum_i (x_i^*)^2 \sum_i y_i \right),$$

სადაც,
$$D = n \sum_i (x_i^*)^4 - \left(\sum_i (x_i^*)^2 \right)^2.$$

თუ გავალოგარიტმებთ მაჩვენებლიანი დამოკიდებულების განტოლებას $y = ae^{bx}$, იგი გადაიქცევა სწორი ხაზის განტოლებად და შესაძლებელი ხდება უმცირეს კვადრატთა მეთოდის გამოყენება, ჩვენ შემთხვევაში გვექნება:

$$\ln y = \ln a + bx$$

და ნორმალურ განტოლებათა სისტემას აქვს შემდეგი სახე:

$$\left. \begin{aligned} \sum_i \ln y_i &= n \ln a + b \sum_i x_i \\ \sum_i \ln y_i x_i &= \ln a \sum_i x_i + b \sum_i x_i^2 \end{aligned} \right\}.$$

ამ სისტემის ამოხსნის შემდეგ მივიღებთ:

$$\ln a = \frac{1}{D} \left[\sum_i \ln y_i \sum_i x_i^2 - \sum_i x_i \ln y_i \sum_i x_i \right], \quad a = e^{\ln a},$$

სადაც,

$$D = n \sum_i x_i^2 - \left(\sum_i x_i \right)^2.$$

ჰიპერბოლური დამოკიდებულების შემთხვევაში $y = a + \frac{b}{x}$, გვექნება შემდეგი ნორმალურ განტოლებათა სისტემა:

$$\left. \begin{aligned} \sum_i y_i &= an + b \sum_i \frac{1}{x_i} \\ \sum_i \frac{y_i}{x_i} &= a \sum_i \frac{1}{x_i} + b \sum_i \frac{1}{x_i^2} \end{aligned} \right\},$$

რომლის ამოხსნის შემდეგ მივიღებთ:

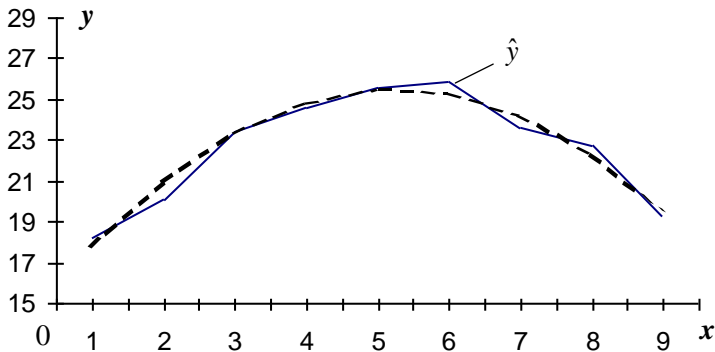
$$a = \frac{1}{D} \left[\sum_i y_i \sum_i \frac{1}{x_i^2} - \sum_i \frac{y_i}{x_i} \sum_i \frac{1}{x_i} \right],$$

$$b = \frac{1}{D} \left[n \sum_i \frac{y_i}{x_i} - \sum_i y_i \sum_i \frac{1}{x_i^2} \right],$$

სადაც,

$$D = n \sum_i \frac{1}{x_i^2} - \left(\sum_i \frac{1}{x_i} \right)^2.$$

მაგალითი. დაკვირვებებმა გვიჩვენა, რომ ძროხის წველადობა ლაქტაციის პერიოდის ცვლილებისას იცვლება ისე, როგორც ეს ნახაზზეა წარმოდგენილი:



შევადგინოთ შემდეგი ცხრილი:

№	ლაქტაცია (თვე) x_i	წველადობა (ც) y_i	$x_i y_i$	x_i^2	$x_i^2 y_i$	x_i^3
1	1	18,2	18,2	1	18,2	1
2	2	20,1	40,2	4	80,4	8
3	3	23,4	70,2	9	210,6	27
4	4	24,6	98,4	16	393,6	64
5	5	25,6	128,0	25	640,0	125
6	6	25,9	155,4	36	932,4	216
7	7	23,6	166,5	49	1156,4	343
8	8	22,7	181,6	64	1452,8	512
9	9	19,2	172,8	81	1555,2	729
Σ	45	203,3	1030,0	285	6439,6	2025

ცხრილის გაგრძელება

№	x_i^4	\hat{y}_i	$(y_i - \hat{y}_i)$	$(y_i - \hat{y}_i)^2$	$(\hat{y}_i - \bar{y})$	$(\hat{y}_i - \bar{y})^2$
1	1	17,6	0,60	0,36	-4,99	24,9
2	16	20,9	-0,8	0,64	-1,69	2,86
3	81	23,3	0,1	0,01	0,71	0,504
4	256	24,8	-0,2	0,04	2,21	4,884
5	625	25,5	0,1	0,01	2,91	8,468
6	1296	25,3	0,6	0,36	2,71	7,344
7	2401	24,2	-0,6	0,36	1,61	2,592
8	4096	22,2	0,5	0,25	-0,39	0,152
9	6561	19,4	0,2	0,04	-3,19	10,176
Σ	15333	203,3		2,07		61,88

შევარჩიოთ შემდეგი სახის რეგრესიის განტოლება:

$$y = a + bx + cx^2.$$

ცხრილის დახმარებით შევადგინოთ ნორმალურ განტოლებათა სისტემა

$$\left. \begin{aligned} 9a + 45b + 285c &= 203,3 \\ 45a + 285b + 2025c &= 1030,0 \\ 285a + 2025b + 15333c &= 6439,6 \end{aligned} \right\},$$

რომლის ამოხსნის შემდეგ მივიღებთ:

$$a = 13,466; \quad b = 4,587 \quad \text{და} \quad c = -0,436.$$

ამრიგად,

$$\hat{y} = 13,466 + 4,587x - 0,436x^2.$$

ჩავსვათ ამ განტოლებაში x -ის მნიშვნელობები

$$x = 1 \quad \hat{y}_1 = 13,466 + 4,587 \cdot 1 - 0,436 \cdot 1 = 17,6;$$

$$x = 2 \quad \hat{y}_2 = 13,466 + 4,587 \cdot 2 - 0,436 \cdot 4 = 20,9$$

და ა.შ. ეს მნიშვნელობები მოცემულია ზემოთ მოყვანილ ცხრილში და გრაფიკზე ნაჩვენებია წყვეტილი მრუდის სახით.

შევამოწმოთ რეგრესიის განტოლების ადეკვატურობა:

$$\bar{y} = 22,59; \sigma_{\text{ნაშ}}^2 = \frac{1}{7} \cdot 2,07 = 0,296; \sigma_y^2 = \frac{1}{7} \cdot 61,88 = 8,84;$$

$$F = \frac{8,84}{0,296} = 29,87; \quad F_{0,05;1;7} = 5,591.$$

რადგან $29,87 > 5,591$, რეგრესიის განტოლება ადეკვატურია.

რეგრესიის განტოლების ადეკვატურობის ზომად შეიძლება გამოვიყენოთ **დეტერმინაციის** კოეფიციენტი, რომელიც გამოითვლება შემდეგი ფორმულით:

$$R^2 = \frac{\sum_{i=1}^n (\hat{y}_i - \bar{y})^2}{\sum_{i=1}^n (y_i - \bar{y})^2} \quad \text{ან} \quad R^2 = \frac{\sum_{i=1}^n (y_i - \hat{y}_i)^2}{\sum_{i=1}^n (y_i - \bar{y})^2}.$$

რაც უფრო დიდია R^2 მნიშვნელობა, მით უფრო ძლიერია ადეკვატურობის ზომა. კერძოდ, თუ $R^2 > 0,95$, მაშინ ითვლება, რომ რეგრესიის მოდელი კარგად აღწერს მოვლენას. თუ $0,8 < R^2 < 0,95$, მაშინ რეგრესიის მოდელი დამაკმაყოფილებად აღწერს მოვლენას და თუ $R^2 < 0,6$, მაშინ ითვლება, რომ მოდელი არაადეკვატურია და საჭიროა მისი შეცვლა.

გვახსოვდეს, რომ R^2 -ს გააჩნია ნაკლი, რომელიც იმაში მდგომარეობს, რომ დეტერმინაციის კოეფიციენტის მაღალი მნიშვნელობა შეიძლება მივიღოთ მცირე ამონარჩევის წყალობით. ამიტომ საჭიროა R^2 სიდიდის კორექტირება

$$\tilde{R}^2 = 1 - \frac{n-1}{n-(m+1)}(1-R^2),$$

სადაც, n - ამონარჩევის განზომილებაა, m - შესაფასებელი პარამეტრების (კოეფიციენტების) რაოდენობა.

დადებითი ფესვი დეტერმინაციის კოეფიციენტიდან გვაძლევს **კორელაციურ ფარდობას η** და თუ X და Y ცვლადებს შორის დამოკიდებულება წრფივია, მაშინ იგი დაემთხვევა კორელაციის კოეფიციენტს, ე.ი. $\eta = r_{xy}$. წინააღმდეგ შემთხვევაში, დამოკიდებულება არაწრფივია.

13.6. ნაშთთა ანალიზი

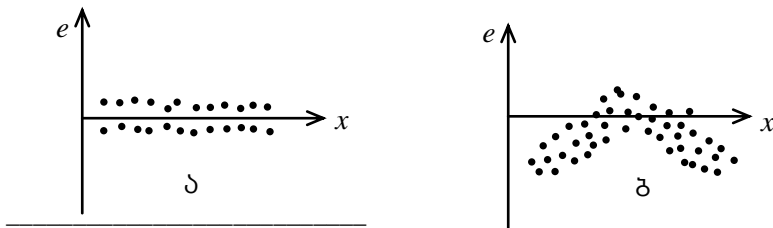
რეგრესიის განტოლების ადეკვატურობისა და კოეფიციენტების შემოწმების გარდა, ხშირად მიმართავენ ნაშთთა ანალიზს. როგორც ვიცით, მოცემულ y_i მნიშვნელობების რეგრესიის განტოლებით გამოთვლილ \hat{y}_i მნიშვნელობებთან გადახრას

$$e_i = y_i - \hat{y}_i, \quad i = 1, 2, \dots, n$$

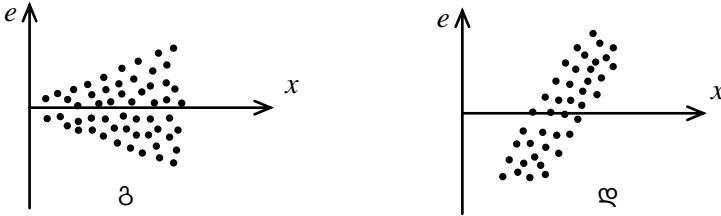
ენოდება ნარჩენი მნიშვნელობები ანუ ნაშთები.

ნაშთთა ანალიზით შესაძლებელია შემოწმდეს ის ძირითადი დაშვებები გადახრების (შეცდომების) მიმართ, რომლებსაც ეფუძნება წრფივი რეგრესია. ჩვენ დავუშვით, რომ რეგრესიის მრუდი წრფივია, გადახრები e_i არიან დამოუკიდებელი, გააჩნიათ ნულოვანი საშუალო, ერთნაირი (მუდმივი) დისპერსია და განაწილებულნი არიან ნორმალურად. თუ შერჩეული რეგრესიის მოდელი ადეკვატურია, მაშინ იგი მეტ-ნაკლებად უნდა აკმაყოფილებდეს დაშვების პირობებს. სწორედ ეს იდეა უდევს საფუძვლად ნაშთთა ანალიზს.

ნაშთთა ანალიზი ხორციელდება გრაფიკულად. აიგება $e = f(x)$ გრაფიკული გამოსახულება და თუ მიღებულ წერტილთა ერთობლიობა დაახლოებით ერთნაირადაა განლაგებული x ღერძის მიმართ (თანაბრად ღერძის ზემოთ და ქვემოთ), მაშინ რეგრესიის განტოლება ადეკვატურია და ყველა დაშვება დაახლოებით შესრულებულია (ნახ. ა).



ე. ყუბანიევილი – ბიომეტრია



თუ დარღვეულია რეგრესიის მრუდის წრფივობის დაშვება, მაშინ ნაშთთა გრაფიკს დაახლოებით ექნება ისეთი სახე, როგორც ეს ნაჩვენებია ბ ნახაზზე. თუ დარღვეულია დისპერსიის მუდმივობა, მაშინ დაახლოებით გვექნება ბ ნახაზზე წარმოდგენილი სურათი. თუ გადახრები დამოკიდებულია x_i -ზე, მაშინ გვექნება დაახლოებით დ ნახაზზე წარმოდგენილი სურათი.

14. მრავლობითი რეგრესიული ანალიზი

14.1. მრავლობითი რეგრესიის მოდელი

წრფივი მოდელი. ორცვლადიანი რეგრესიის დროს y მნიშვნელობა დამოკიდებულია მხოლოდ ერთ x ცვლადზე. საზოგადოდ, y შეიძლება იყოს მრავალი ცვლადის ფუნქცია. განვიხილოთ ეს შემთხვევა. ვთქვათ, დამოკიდებულ Y -სა და დამოუკიდებელ X_1, X_2, \dots, X_n ცვლადებს შორის არსებობს ასეთი წრფივი დამოკიდებულება:

$$Y = a_1X_1 + a_2X_2 + \dots + a_nX_n,$$

სადაც, Y, X_1, X_2, \dots, X_n ვექტორებია. წარმოვადგინოთ ეს დამოკიდებულება მატრიცული სახით $Y = AX$, სადაც, $A = [a_i], i = 1, 2, \dots, n$ საძიებელი კოეფიციენტების ვექტორია, $Y = [y_i], i = 1, 2, \dots, m$ დამოკიდებული ცვლადის ვექტორია, ხოლო $X = [x_{ij}], i = 1, 2, \dots, m, j = 1, 2, \dots, n$ დამოუკიდებელი ცვლადების მატრიცაა.

a_i კოეფიციენტები მოძებნოთ უმცირეს კვადრატთა

მეთოდით, რომლის თანახმად $Q = \sum_{i=1}^m [y_i - \hat{y}_i]^2 = \min.$

ჩვენერთ ეს გამოსახულება მატრიცული სახით

$$Q = (Y - XA)'(Y - XA) = Y'Y - A'X'Y - Y'XA + A'X'XA.$$

რადგან $A'X'Y = Y'XA$, ამიტომ გვექნება:

$$Q = Y'Y - 2A'X'Y + A'X'XA.$$

გავანარმოთ ეს გამოსახულება

$$\frac{\partial Q}{\partial a} = -2X'Y + 2(X'X)A = 0,$$

მაშინ ნორმალურ განტოლებათა სისტემას ექნება შემდეგი სახე:

$$X'Y = X'XA,$$

საიდანაც, $A = (X'X)^{-1}X'Y$.

რადგან X მატრიცა შეიცავს n წრფივად დამოუკიდებელ ვექტორ-სვეტებს, ამიტომ, როგორც ცნობილია, X მატრიცის რანგი $(m-1)$ -ის ტოლია. ეს, თავის მხრივ, მიგვანიშნებს იმაზე, რომ $|X'X| \neq 0$, ე.ი. მატრიცას გააჩნია შებრუნებული მატრიცა $(X'X)^{-1}$.

ხშირად რეგრესიის განტოლება შეიცავს თავისუფალ წევრს

$$\hat{Y} = a_0 + a_1X_1 + a_2X_2 + \dots + a_nX_n.$$

მაშინ a_0 პარამეტრის შეფასებისთვის საჭიროა X მატრიცას დაემატოს ერთეულოვანი ვექტორ-სვეტი

$$X = \begin{bmatrix} 1 & x_{11} & x_{12} & \dots & x_{1n} \\ 1 & x_{21} & x_{22} & \dots & x_{2n} \\ \dots & \dots & \dots & \dots & \dots \\ 1 & x_{m1} & x_{m2} & \dots & x_{mn} \end{bmatrix},$$

მაშინ გვექნება:

$$XX' = \begin{bmatrix} 1 & 1 & \dots & 1 \\ x_{11} & x_{21} & \dots & x_{m1} \\ \dots & \dots & \dots & \dots \\ x_{1n} & x_{2n} & \dots & x_{mn} \end{bmatrix} \begin{bmatrix} 1 & x_{11} & x_{12} & \dots & x_{1n} \\ 1 & x_{21} & x_{22} & \dots & x_{2n} \\ \dots & \dots & \dots & \dots & \dots \\ 1 & x_{m1} & x_{m2} & \dots & x_{mn} \end{bmatrix} =$$

$$= \begin{bmatrix} n & \sum_i x_i & \dots & \sum_i x_{im} \\ \sum_i x_i & \sum_i x_i^2 & \dots & \sum_i x_i x_{im} \\ \dots & \dots & \dots & \dots \\ \sum_i x_{im} & \sum_i x_{i1} x_{im} & \dots & \sum_i x_{im}^2 \end{bmatrix}.$$

ნარჩენი დისპერსია გამოითვლება შემდეგი ფორმულით:

$$\sigma_{\text{ნარ}}^2 = \frac{1}{m-n-1} \sum_{i=1}^m (y_i - \hat{y}_i)^2,$$

სადაც, n - ცვლადების რაოდენობაა, m - დაკვირვებათა რაოდენობა.

რეგრესიის განტოლების ადეკვატურობის შესამოწმებლად საჭიროა გამოითვალოს შემდეგი სტატისტიკა:

$$F = \frac{Q}{\sigma_{\text{ნარ}}^2}, \quad Q = \frac{1}{n} \sum_{i=1}^m (\hat{y}_i - \bar{y})^2.$$

α მნიშვნელოვნების დონითა და $v_1 = n$, $v_2 = m-n-1$ თავისუფლების ხარისხებით ფიშერის განაწილების ცხრილიდან შეირჩევა $F_{\alpha; v_1; v_2}$ კრიტიკული მნიშვნელობა. როცა $F \geq F_{\alpha; v_1; v_2}$, მაშინ რეგრესიის განტოლება ადეკვატურია, წინააღმდეგ შემთხვევაში, როცა $F < F_{\alpha; v_1; v_2}$, განტოლება არაადეკვატურია.

რეგრესიის განტოლების კოეფიციენტების შემოწმება სარწმუნოებაზე ხდება შემდეგი $H_0: a_i = 0$ ნულოვანი ჰიპოთეზის საშუალებით. ამისათვის საჭიროა გამოითვალოს სტატისტიკა:

$$t_i = \frac{a_i}{\sigma_i}, \quad \text{სადაც, } \sigma_i = \sqrt{\sigma_{\text{ნარ}}^2 \cdot s_{ii}}, \quad i=1,2,\dots,n.$$

s_{ii} წარმოადგენს $(X'X)^{-1}$ მატრიცის მთავარ დიაგონალზე მყოფი ელემენტის მნიშვნელობას. α მნიშვნელოვნების დონითა და $v = m - n - 1$ თავისუფლების ხარისხით სტიუდენტის განაწილების ცხრილიდან მოიძებნება $t_{\alpha; v}$ კრიტიკული მნიშვნელობა. თუ $|t_i| < t_{\alpha; v}$ მაშინ ნულოვანი ჰიპოთეზა მიიღება, ე.ი. i -ური კოეფიციენტის მნიშვნელობა ახლოსაა ნულთან და შესაძლებელია მისი გამო-რიცხვა რეგრესიის განტოლებიდან. თუ $|t_i| \geq t_{\alpha; v}$, მაშინ a_i კოეფიციენტი სარწმუნოა.

მაგალითი. ცხრილში მოცემულია y დამოკიდებული და x_1, x_2 დამოუკიდებელი ცვლადების მნიშვნელობები.

i	y_i	x_{1i}	x_{2i}	\hat{y}	$(y - \hat{y})$	$(y - \hat{y})^2$	$(\hat{y} - \bar{y})$	$(\hat{y} - \bar{y})^2$
1	10	2	1	10,256	-0,256	0,066	-3,444	11,861
2	12	2	2	10,868	1,132	1,281	-2,832	8,020
3	17	8	10	16,532	0,468	0,219	2,832	8,020
4	13	2	4	12,091	0,909	0,826	-1,609	2,589
5	15	6	8	15,052	-0,052	0,003	1,352	1,828
6	10	3	4	12,22	-2,22	4,928	-1,480	2,190
7	14	5	7	14,312	-0,312	0,098	0,612	0,375
8	12	3	3	11,608	0,392	0,154	-2,092	4,377
9	16	9	10	16,661	-0,661	0,437	2,961	8,768
10	18	10	11	17,401	0,599	0,359	3,701	13,697
Σ	137	50	60			8,371		61,725

$$\bar{y} = 13,7; \bar{x}_1 = 5,0; \bar{x}_2 = 6,0.$$

შევარჩიოთ რეგრესიის განტოლება $\hat{y} = a_0 + a_1x_1 + a_2x_2$.

ამრიგად გვაქვს:

$$Y = \begin{bmatrix} 10 \\ 12 \\ 17 \\ \dots \\ 18 \end{bmatrix} \quad X = \begin{bmatrix} 1 & 2 & 1 \\ 1 & 2 & 2 \\ 1 & 8 & 10 \\ \dots & \dots & \dots \\ 1 & 10 & 11 \end{bmatrix};$$

$$X'X = \begin{bmatrix} 1 & 1 & 1 & \dots & 1 \\ 2 & 2 & 8 & \dots & 10 \\ 1 & 2 & 10 & \dots & 11 \end{bmatrix} \begin{bmatrix} 1 & 2 & 1 \\ 1 & 2 & 2 \\ 1 & 8 & 10 \\ \dots & \dots & \dots \\ 1 & 10 & 11 \end{bmatrix} = \begin{bmatrix} 10 & 50 & 60 \\ 50 & 336 & 398 \\ 60 & 398 & 480 \end{bmatrix};$$

$$(X'X)^{-1} = \begin{bmatrix} 0,40168 & -0,01676 & -0,03631 \\ -0,01676 & 0,16760 & -0,13687 \\ -0,03631 & -0,13687 & 0,12011 \end{bmatrix};$$

$$XY = \begin{bmatrix} 1 & 1 & 1 & \dots & 1 \\ 2 & 2 & 8 & \dots & 10 \\ 1 & 2 & 10 & \dots & 11 \end{bmatrix} \begin{bmatrix} 10 \\ 12 \\ 17 \\ \dots \\ 18 \end{bmatrix} = \begin{bmatrix} 137 \\ 756 \\ 908 \end{bmatrix};$$

$$A = (XX)^{-1}XY = \begin{bmatrix} 0,40168 & -0,01676 & -0,03631 \\ -0,01676 & 0,16760 & -0,13687 \\ -0,03631 & -0,13687 & 0,12011 \end{bmatrix} \begin{bmatrix} 137 \\ 756 \\ 908 \end{bmatrix} = \begin{bmatrix} 9,3872 \\ 0,1285 \\ 0,6174 \end{bmatrix}.$$

ამრიგად, რეგრესიის განტოლებას აქვს შემდეგი სახე:

$$\hat{y} = 9,3872 + 0,1285x_1 + 0,6174x_2.$$

თუ x_1 და x_2 სიდიდეები იზომება ერთი და იგივე ფიზიკური ერთეულით, მაშინ ადვილი შესამჩნევია, რომ x_2 -ის ზეგავლენა y -ზე x_1 -თან შედარებით, დაახლოებით ხუთჯერ უფრო ძლიერია. თუ x_1 და x_2 სხვადასხვა ფიზიკურ ერთეულებშია წარმოდგენილი, ასეთი შედარება უაზრობაა.

გამოვთვალოთ ნარჩენი დისპერსია

$$\sigma_{\text{ნარ}}^2 = \frac{8,371}{10 - 2 - 1} = 1,1959; \quad \sigma_y^2 = \frac{61,725}{2} = 30,863.$$

შევამოწმოთ რეგრესიის განტოლების ადეკვატურობა

$$F = \frac{30,863}{1,1959} = 25,807; \quad F_{0,05;2;7} = 4,74.$$

რადგან $25,807 > 4,74$, რეგრესიის განტოლება ადეკვატურია. კოეფიციენტების შესამოწმებლად განვიხილოთ სტატისტიკა:

$$t_i = \frac{a_i}{\sigma_i}; \quad i = 0,1,2; \quad \sigma_0 = \sqrt{1,1959 \cdot 0,40168} = 0,698;$$

$$\sigma_1 = \sqrt{1,1959 \cdot 0,1676} = 0,4478; \quad \sigma_2 = \sqrt{1,1959 \cdot 0,12011} = 0,379;$$

$$t_0 = \frac{9,3872}{0,6931} = 13,54; \quad t_1 = \frac{0,1285}{0,4477} = 0,287; \quad t_2 = \frac{0,6174}{0,379} = 1,629;$$

$t_{0,05;7} = 2,37$. როგორც ვხედავთ, გარდა a_0 -სა, დანარჩენი a_1 და a_2 კოეფიციენტები სარწმუნო არ არიან.

არანრფივი მოდელი. პრაქტიკულ კვლევებში, ზოგჯერ, შეუძლებელი ხდება მრავლობითი ნრფივი რეგრესიის განტოლების გამოყენება მისი არაადეკვატურობის გამო. ამ შემთხვევაში, უნდა გამოვიყენოთ ისეთი არანრფივი მოდელი, რომელიც არანრფივია ცვლადების მიმართ, მაგრამ ნრფივია რეგრესიის კოეფიციენტების მიმართ. ამ შემთხვევაში, კოეფიციენტების შეფასება ხდება საკმაოდ ადვილად. კერძოდ, დამოკიდებული ცვლადები შეიცვლება ახალი პირველი ხარისხის ცვლადებით და მის მიმართ გამოიყენება ნრფივი უმცირეს კვადრატთა მეთოდი. მაგალითისთვის განვიხილოთ შემდეგი არანრფივი რეგრესიის განტოლება:

$$y = a_0 + a_1x_1 + a_2x_2 + a_3x_1^2 + a_4x_2^2.$$

შემოვიტანოთ ახალი ცვლადები $z_1 = x_1, z_2 = x_2, z_3 = x_1^2, z_4 = x_2^2$, მაშინ გვექნება:

$$y = a_0 + a_1z_1 + a_2z_2 + a_3z_3 + a_4z_4.$$

უმცირეს კვადრატთა მეთოდის გამოყენებით მივიღებთ შემდეგ ნორმალურ განტოლებათა სისტემას:

$$Z'Y = Z'ZA, \text{ საიდანაც } A = (Z'Z)^{-1}Z'Y.$$

მიღებული რეგრესიის განტოლებისა და კოეფიციენტების შეფასებები ხდება ისევე, როგორც მრავლობითი ნრფივი რეგრესიის დროს.

14.2. მრავლობითი ნრფივი რეგრესიის განტოლების დოზის ინტერპოლაცია

განვსაზღვროთ მრავალგანზომილებიანი ნრფივი რეგრესიის განტოლების **ნდობის ინტერვალი**. ამისათვის გამოვთვალოთ $\hat{y} = a_0 + a_1x_1 + a_2x_2 + \dots + a_nx_n$ გამოსახულების დისპერსია

$$D(\hat{y}) = D(a_0 + a_1x_1 + a_2x_2 + \dots + a_nx_n).$$

თუ გავიხსენებთ დამოკიდებული ცვლადების ჯამის დისპერსიის თვისებას, მივიღებთ:

$$\sigma_{\hat{y}}^2 = D(\hat{y}) = D(a_0) + x_1^2 D(a_1) + x_2^2 D(a_2) + \dots + x_n^2 D(a_n) + 2x_1 \text{cov}(a_0 a_1) + 2x_1 x_2 \text{cov}(a_1 a_2) + \dots + 2x_{n-1} x_n \text{cov}(a_{n-1} a_n) \quad (14.1)$$

გავიხსენოთ, რომ $\text{cov}(a_i, a_j) = M(a_i, a_j)$. ჩავწეროთ (17.1) გამოსახულება მატრიცული სახით

$$D(\hat{Y}) = X'_p \text{cov}(a) X_p, \quad (14.2)$$

სადაც, $X_p = (1, X_{p_1}, X_{p_2}, \dots, X_{p_n})$ მოცემული დამოუკიდებელი ცვლადის ვექტორია. $\text{cov}(a) - a$ შეფასების კოვარიაციული მატრიცაა. ვიპოვოთ a პარამეტრის დისპერსია

$$\text{cov}(a) = M[(a - \alpha)(a - \alpha)'],$$

$$a = (X'X)^{-1} X'Y = (X'X)^{-1} X'(X\alpha + \varepsilon) = \alpha + (X'X)^{-1} X'\varepsilon.$$

აქედან,

$$a - \alpha = (X'X)^{-1} X'\varepsilon,$$

$$\text{cov}(a) = M[(X'X)^{-1} X'\varepsilon\varepsilon'X(X'X)^{-1}] = M(\varepsilon\varepsilon')(X'X)^{-1} \quad (14.3)$$

$M(\varepsilon\varepsilon')$ წარმოადგენს მატრიცას, რომლის ყველა ელემენტი, გარდა მთავარ დიაგონალზე მყოფისა, ნულის ტოლია, რადგან ჩვენი დაშვებით, შემთხვევითი ცდომილებები ერთმანეთთან არ არიან კორელირებული, ე.ი. $M(\varepsilon\varepsilon') = 0$. რაც შეეხება მთავარ დიაგონალზე მყოფ ელემენტებს, ისინი წარმოადგენენ დისპერსიებს და რადგან ჩვენივე დაშვებით, შემთხვევით ცდომილებებს გააჩნიათ ერთი და იგივე დისპერსია, ამიტომ

$$M(\varepsilon\varepsilon') = \sigma_{\text{ნარ}}^2 I,$$

სადაც, I - ერთეულოვანი მატრიცაა. მიღებული შედეგი ჩავსვათ (14.3)-ში:

$$\text{cov}(a) = \sigma_{\text{ნარ}}^2 (X'X)^{-1},$$

ხოლო ეს უკანასკნელი კი - (17.2)-ში:

$$D(\hat{y}) = \sigma_{\text{ნარ}}^2 X'_p (X'X)^{-1} X_p.$$

ამრიგად, \hat{y} მნიშვნელობისთვის გვექნება შემდეგი ნდობის ინტერვალი:

$$\hat{y} \pm t_{\alpha;v} \sigma_{\text{ნარ}} \sqrt{X'_p (X'X)^{-1} X_p}.$$

14.3. ცვლადების შერჩევა

რეგრესიული ანალიზი საშუალებას გვაძლევს მოცემულ X_1, X_2, \dots, X_n დამოუკიდებელ ცვლადებიდან შევარჩიოთ ის ცვლადები, რომლებიც კავშირშია დამოკიდებულ Y ცვლადთან და გამოვრიცხოთ რეგრესიის განტოლებიდან არაინფორმატიული ანუ ის ცვლადები, რომლებიც ნაკლებად ან სულაც არ არიან კავშირში Y ცვლადთან. განვიხილოთ რამდენიმე მეთოდი.

ყველა შესაძლო რეგრესიათა მეთოდი. თუ დამოუკიდებელ ცვლადთა რაოდენობა n არც ისე დიდია, მაშინ რეკომენდებულია, მოისინჯოს ყველანაირი კომბინაცია რეგრესიის წრფივი მოდელის ასაგებად და რაიმე კრიტერიუმით არჩეულ იქნეს მათ შორის საუკეთესო. მაგალითად, თუ $n = 4$, მაშინ უნდა მოისინჯოს $2^4 = 16$ მოდელი. კერძოდ, 4 – ერთცვლადიანი, 6 – ორცვლადიანი, 4 – სამცვლადიანი, 1 – ოთხცვლადიანი, და ბოლოს ერთი, რომელიც არ შეიცავს არც ერთ ცვლადს, გარდა თავისუფალი წევრისა. საუკეთესოდ ითვლება ის ადეკვატური მოდელი, რომელსაც აქვს უმცირესი ნარჩენი დისპერსია ან, რაც იგივეა, უდიდესი დეტერმინაციის კოეფიციენტი.

გამორიცხვის მეთოდი. თუ დამოუკიდებელ ცვლადთა რაოდენობა საკმაოდ დიდია, მაშინ ყველა შესაძლო რეგრესიათა მეთოდის გამოყენება შეუძლებელია. არსებობს რამდენიმე ალტერნატიული მეთოდი, მათ შორის ცვლადების გამორიცხვის ანუ ცვლადების შერჩევის მიმდევრობითი მეთოდი, რომლის არსი შემდეგში მდგომარეობს: დასაწყისში განიხილება მოდელი, რომელიც შეიცავს ყველა განსახილველ ცვლადებს. რეგრესიის განტოლების კოეფიციენტების შემონმება სარწმუნოებაზე ხდება $H_0 : a_i = 0$ ნულოვანი ჰიპოთეზის საშუალებით. ამისათვის, თითოეული კოეფიციენტისათვის უნდა განისაზღვროს შემდეგი

სტატისტიკა: $t_i = \frac{a_i}{\sigma_i}, i = 1, 2, \dots, n$, სადაც $\sigma_i = \sqrt{\sigma_{\text{ნარ}}^2 \cdot s_{ii}}$, s_{ii}

ნარმოადგენს $(XX)^{-1}$ მატრიცის მთავარ დიაგონალზე მყოფი ელემენტის მნიშვნელობას. თუ გამოთვლილ სტატისტიკებს შორის მინიმალური t_i ფარდობის აბსოლუტური სიდიდე $|t_i| < t_{\alpha;v}$, სადაც $t_{\alpha;v}$ სიდიდე აღებულია სტიუდენტის განაწილების ცხრილიდან α

მნიშვნელოვნების დონითა და $v = m - n - 1$ თავისუფლების ხარისხით, მაშინ ნულოვანი ჰიპოთეზა მიიღება, ე.ი. i -ური კოეფიციენტის მნიშვნელობა ახლოსაა ნულთან და შესაძლებელია მისი გამორიცხვა რეგრესიის განტოლებიდან. შემდეგ ბიჯზე განიხილება მოდელი დარჩენილი $(n - 1)$ ცვლადით და მონდება რეგრესიის განტოლების ადეკვატურობა. თუ აღმოჩნდება, რომ რეგრესიის განტოლება არაადეკვატურია, მაშინ გამორიცხული i -ური კოეფიციენტი, ანუ X_i პარამეტრი უნდა დავაბრუნოთ რეგრესიის განტოლებაში.

პროცედურა გაგრძელდება მანამ, სანამ ყველა გამოთვლილი t_i ფარდობის აბსოლუტური სიდიდე არ აღმოჩნდება $\geq t_{\alpha;v}$. საბოლოოდ, რეგრესიის განტოლებაში რჩება ის ცვლადები, რომლებიც არ გამოირიცხა პროცედურის ჩატარებისას.

ჩართვის მეთოდი. გამორიცხვის მეთოდის ალტერნატივაა ჩართვის მეთოდი, რომელსაც ბიჯურ რეგრესიას უწოდებენ. პირველ ბიჯზე განიხილება ერთცვლადიანი მოდელები. თითოეულ

მოდელში ხდება a_0 და a_i კოეფიციენტების შეფასება და $t_i = \frac{a_i}{\sigma_i}$

სტატისტიკის განსაზღვრა ისე, როგორც ჩართვის მეთოდის დროს. ის ცვლადი, რომელსაც შეესაბამება $|t_i|$ სიდიდის მაქსიმალური მნიშვნელობა, ჩაირთვება მოდელში მხოლოდ იმ პირობით, რომ ეს მაქსიმალური მნიშვნელობა აღემატება $t_{\alpha;v}$ კრიტიკულ მნიშვნელობას. დაუშვათ, ასეთი ცვლადია X_1 . ამის შემდეგ განიხილება ორცვლადიანი მოდელი X_1 ცვლადის დანარჩენ ცვლადებთან: $(X_1, X_2), (X_1, X_3), \dots, (X_1, X_n)$. ხდება თითოეული მოდელის შეფასება და დარჩენილ X_2, X_3, \dots, X_n ცვლადებიდან მოდელში ჩაირთვება ის ცვლადი, რომლისთვისაც $|t_i|$ სიდიდე მაქსიმალურია და აღემატება $t_{\alpha;v}$. დაუშვათ ასეთი ცვლადია X_2 . შემდეგ განიხილება $(X_1, X_2, X_3), (X_1, X_2, X_4), \dots, (X_1, X_2, X_n)$ სამცვლადიანი მოდელები, სადაც ტარდება იგივე პროცედურა და ა.შ. პროცედურა გრძელდება მანამ, სანამ რომელიმე ბიჯზე არ აღმოჩნდება, რომ ყველა t_i ფარდობის აბსოლუტური სიდიდე $\leq t_{\alpha;v}$ მნიშვნელობაზე.

არსებობს ჩართვის სხვა მეთოდებიც. მაგალითად, კორელაციის კერძო კოეფიციენტების გამოყენების მეთოდი. ამ შემ-

თხვევაში განისაზღვრება Y ცვლადის ყველა დამოუკიდებელ X_1, X_2, \dots, X_n ცვლადებს შორის კორელაციის კერძო კოეფიციენტები. ის დამოუკიდებელი ცვლადი, რომელსაც გააჩნია Y -თან უდიდესი კერძო კორელაციის კოეფიციენტი, ჩაირთვება რეგრესიის მოდელში და რეგრესიის განტოლება მონმდება ადეკვატურობაზე. თუ აღმოჩნდება, რომ განტოლება არაადეკვატურია, მაშინ მოდელში ჩაირთვება შემდეგი ცვლადი, რომელსაც დამოუკიდებელ Y ცვლადთან გააჩნია უდიდესი კორელაციის კერძო კოეფიციენტი, კვლავ მონმდება განტოლების ადეკვატურობა და ა.შ. მანამ, სანამ არ მივიღებთ რეგრესიის ადეკვატურ განტოლებას.

კომბინირებული მეთოდი. იგი წარმოადგენს ჩართვისა და გამორიცხვის მეთოდების კომბინაციას. პროცედურა იწყება ისევე, როგორც ჩართვის მეთოდი, ე.ი. ცვლადების მიმდევრობით ჩართვით მოდელში, ოღონდ ყოველი ახალი ცვლადის დამატების შემდეგ მონმდება ადრე ჩართული ცვლადები, კერძოდ, ხომ არ მოიძებნა მათ შორის გამოსარიცხი. მაგალითად, ვთქვათ, მოდელში ჩართულია X_3, X_5, X_6 და X_7 ცვლადები, რომელთაგან X_7 წარმოადგენს ამ ბიჯზე დამატებულ ცვლადს, მაშინ მონმდება X_3, X_5 და X_6 ცვლადები ისევე, როგორც გამორიცხვის მეთოდშია აღწერილი. თუ ამ ცვლადებიდან რომელიმე აღმოჩნდა ისეთი, რომ სრულდება $|t_i| < t_{\alpha;v}$ პირობა, მაშინ ეს ცვლადი მოდელიდან გამორიცხება და ა.შ.

მიუხედავად ცვლადების შერჩევის მრავალფეროვნებისა, უნდა გვახსოვდეს, რომ არ არსებობს იმის გარანტია, რომ ყოველ ცალკეულ შემთხვევაში მიიღება ჩვენთვის სასურველი ადეკვატური მოდელი. აქედან გამომდინარე, სასურველია თავიდან გამორიცხოს ის დამოუკიდებელი ცვლადები, რომელთა ჩართვა რეგრესიის მოდელში, ამა თუ იმ მოსაზრებით, აზრს მოკლებულია.

14.4. ნარჩენების ავტოკორელაციისა და მულტიკოლინეარობის პრობლემა

რეგრესიულ ანალიზში, უმცირეს კვადრატთა მეთოდის გამოყენებისას, ჩვენ დაუშვით, რომ ε_i ცდომილებები შემთხვევითი დამოუკიდებელი (არაკორელირებული) სიდიდეებია ნულო-

ვანი საშუალოთი. პრაქტიკაში ამ მოთხოვნის შესრულება ძნელია, განსაკუთრებით დროითი მწკრივებისათვის.

აღმოჩნდა, რომ თუ e_i ნარჩენები ერთმანეთში კორელირებენ, მაშინ ამბობენ, რომ ადგილი აქვს ცდომილების ავტოკორელაციას. მიუხედავად იმისა, რომ უმცირეს კვადრატთა მეთოდი ამ შემთხვევის დროსაც გვაძლევს გადაუადგილებად და საფუძვლიან შეფასებებს, რეგრესიის პარამეტრების განსაზღვრისას ნდობის ინტერვალის განსაზღვრა კარგავს აზრს მისი არასაიმედობის გამო. ამიტომ, თუ აღმოჩნდება ნარჩენების ავტოკორელაციის ეფექტი, საჭიროა გადაისინჯოს რეგრესიის განტოლების მოდელი.

არსებობს ავტოკორელაციის აღმოჩენის მთელი რიგი მეთოდები. ჩვენ განვიხილავთ პირველი რიგის ავტოკორელაციის არსებობის ჰიპოთეზის შემოწმების შედარებით მარტივ და საკმაოდ საიმედო მეთოდს, რომელიც შემოგვთავაზებს დარბინმა და უიტსონმა. ამ ჰიპოთეზის შესამოწმებლად გამოვთვალოთ შემდეგი სტატისტიკა:

$$d = \frac{\sum_{i=1}^n (e_i - e_{i+1})^2}{\sum_{i=1}^n e_i^2}.$$

n -ის დიდი რაოდენობის დროს $\sum_{i=1}^n e_i \approx \sum_{i=1}^n e_{i-1}$, მაშინ

$$d \approx \frac{2 \sum_{i=2}^n e_i^2 - 2 \sum_{i=2}^n e_i e_{i-1}}{\sum_{i=1}^n e_i^2} = 2 \left(1 - \frac{\sum_{i=2}^n e_i e_{i-1}}{\sum_{i=1}^n e_i^2} \right) = 2(1 - R),$$

სადაც, R წარმოადგენს პირველი რიგის ავტოკორელაციის კოეფიციენტს და თუ იგი ნულის ტოლია, მაშინ ავტოკორელაცია არ არსებობს. თუ ავტოკორელაცია მთლიანად არსებობს, მაშინ

$R = \pm 1$. აქედან გამომდინარე, თუ ავტოკორელაცია არ არსებობს, მაშინ d სიდიდის მნიშვნელობა მიახლოებით 2-ის ტოლია და მთლიანი ავტოკორელაციის არსებობის დროს იგი 0-ის ან 4-ის ტოლია.

ავტოკორელაციის სარწმუნოების დასადგენად სპეციალური ცხრილიდან დამოკიდებული n პარამეტრების რაოდენობისა და m დაკვირვებათა რაოდენობის საშუალებით მოიძებნება d კრიტერიუმის ქვედა d_e და ზედა d_u ზღვრების მნიშვნელობები. თუ გამოთვლილი d სტატისტიკა მოთავსებულია d_u და $(4 - d_u)$ საზღვრებში, მაშინ ჰიპოთეზა ავტოკორელაციის არარსებობის შესახებ მიიღება. თუ d მოთავსებულია d_e და d_u შორის ან $(4 - d_u)$ და $(4 - d_e)$ შორის, მაშინ ჩვენ არა გვაქვს საფუძველი ჰიპოთეზის უარსაყოფად და არც მისაღებად, ე.ი. საქმე გვაქვს განუსაზღვრელობასთან. თუ $d < d_e$, საქმე გვაქვს დადებით ავტოკორელაციასთან, ხოლო როცა $d > (4 - d_e)$ – უარყოფით ავტოკორელაციასთან.

რეგრესიის განტოლების ფორმირებისას, ხშირად ვაწყდებით მულტიკოლინეარობის პრობლემას. როგორც აღვნიშნეთ, დამოუკიდებელი ცვლადები X_1, X_2, \dots, X_n რეგრესიის განტოლებაში უნდა იყვნენ ურთიერთდამოუკიდებელი, მაგრამ ამ პირობის შესრულება პრაქტიკულად საკმაოდ რთულია, განსაკუთრებით, ბიოსამედიცინო კვლევებში. ამ მოვლენას მულტიკოლინეარობა ეწოდება. თუ ცვლადებს შორის დამოკიდებულება ფუნქციონალურია, მაშინ საქმე გვაქვს მკაცრ მულტიკოლინეარობასთან, ხოლო თუ დამოკიდებულება არც ისე მკაცრია და გამოვლინდება მიახლოებით, მაშინ მულტიკოლინეარობა არ არის მკაცრი.

უმცირეს კვადრატთა მეთოდის ერთ-ერთი მოთხოვნა ისაა, რომ დამოუკიდებელ ცვლადებს შორის არ უნდა არსებობდეს წრფივი კავშირი. მულტიკოლინეარობის არსებობა იწვევს ამ მოთხოვნის დარღვევას. ფორმალურად რეგრესიის განტოლება ამ შემთხვევაში მიიღება, მაგრამ იგი არ არის საიმედო იმ ვაგებით, რომ სანყისი მონაცემების უმნიშვნელო ცვლილებამ შეიძლება გამოიწვიოს პარამეტრთა შეფასების მკვეთრი ცვლილებები.

მულტიკოლინეარობის აღმოჩენის მეთოდებიდან შეგვიძლია განვიხილოთ კორელაციური მატრიცის მეთოდი. კორელაციის კოეფიციენტები, რომლებიც ახლოს არიან ± 1 სიდიდესთან, მიგვანიშნებენ მულტიკოლინეარობის არსებობაზე. უფრო საიმედო მეთოდია $X'X$ მატრიცის დეტერმინანტის განსაზღვრა. თუ ეს სიდიდე ნულთან ახლოსაა, მაშინ საქმე გვაქვს მულტიკოლინეარობასთან. მულტიკოლინეარობის აღმოსაჩენად შეგვიძლია გამოვიყენოთ შემდეგი სტატისტიკა:

$$\chi^2 = - \left[m - 1 - \frac{1}{6}(2n + 5) \right] \lg(\det[\tilde{X}\tilde{X}]),$$

რომელსაც გააჩნია χ^2 განაწილება $v = \frac{n(n-1)}{2}$ თავისუფლების ხარისხით. აქ, m —დაკვირვებათა რაოდენობაა, n —დამოუკიდებელი ცვლადების რაოდენობა. $[\tilde{X}\tilde{X}]$ მატრიცის ელემენტები განისაზღვრება სანყისი $[XX]$ მატრიციდან შემდეგნაირად:

$$\tilde{x}_{ik} = \frac{x_{ik} - \bar{x}_k}{\sigma_k \sqrt{m}},$$

სადაც, \bar{x}_k , σ_k — შესაბამისად i -ური ცვლადის საშუალო არითმეტიკული და საშუალო კვადრატული გადახრაა.

თუ $\chi^2 \geq \chi_{\alpha;v}^2$, მაშინ მულტიკოლინეარობა არ არსებობს, ნინაალმდეგ შემთხვევაში, მისი არსებობა სარწმუნოა.

მულტიკოლინეარობის გამორიცხვა შესაძლებელია რეგრესიის განტოლების სტრუქტურის გადასინჯვით, კერძოდ, ორი დამოკიდებული ცვლადიდან უნდა გამოირიცხოს ერთი. მეორე გზაა ცვლადების გარდაქმნა ისე, რომ ისინი გახდნენ ურთიერთდამოუკიდებელი, მაგალითად, მთავარი კომპონენტების მეთოდის გამოყენებით.

15. დისპერსიული ანალიზის საფუძვლები

15.1. მეთოდის არსი

ორი ამონარჩევის შედარების მარტივი მეთოდისგან განსხვავებით, პრაქტიკაში ხშირად გვხვდება ისეთი შემთხვევები, სადაც საჭიროა ერთმანეთს ერთდროულად შევადაროთ რამდენიმე ამონარჩევი, რომლებიც გაერთიანებულია ერთიან სტატისტიკურ კომპლექსში. ამ შემთხვევაში საშუალოების ნყვილ-ნყვილი შედარება, თუ ამონარჩევების (ჯგუფების) რაოდენობა დიდია, მი-

ზანშეუნონელია არა მარტო მეთოდოლოგიური შეცდომების (იხ. §10.5), არამედ დიდი გამოთვლითი სამუშაოს გამო. მაგალითად, თუ გვაქვს 7 ამონარჩევი, მაშინ ჩასატარებელი იქნება $C_7^2 = 21$ საშუალოების შედარება. ცხადია, რომ ჯგუფების რაოდენობის ზრდისას, შესადარებელ წყვილთა რაოდენობა ძალზე სწრაფად იზრდება. ასე მაგალითად, 14 ჯგუფისათვის გვექნება 91 შედარება.

გაითვალისწინა რა ეს პრობლემა, რ. ფიშერმა შემოგვთავაზა საშუალოების კომპლექსური შედარების მეთოდი, რომელსაც დისპერსიული ანალიზი (ANOVA – *Analysis of Variance*) ეწოდება. დისპერსიული ანალიზის მეთოდი ეფუძნება საერთო დისპერსიის დაშლას დამოუკიდებელ მდგენელებად, რომლებიც გამოწვეულია როგორც რეგულირებადი, ისე არარეგულირებადი ფაქტორებით. შემოვიტანოთ რამდენიმე განმარტება.

პარამეტრებს, რომლებიც რაიმე მიზეზით იცვლებიან, ეწოდებათ შედეგობრივი. მიზეზებს, რომლებიც იწვევენ შედეგობრივი პარამეტრების ცვლილებას, ეწოდებათ ფაქტორები. მაგალითად, სხეულის მასა, მოსავლიანობა, მოსწავლეთა მოსწრება და ა.შ. წარმოადგენს შედეგობრივ პარამეტრებს, რომლებზედაც სხვადასხვა ფაქტორები ახდენენ ზემოქმედებას: წამლის ან ტოქსიკური ნივთიერებათა დოზები, სასუქის რაოდენობა, კვების რეჟიმი, ფიზიკური და გონებრივი სავარჯიშოები და სხვა. არსებობს ერთი და იმავე პარამეტრზე მოქმედი მრავალი ფაქტორი, რომელთაგან ცდაში (ექსპერიმენტში) რეგულირებადია მხოლოდ ზოგიერთი მათგანი და მათ რეგულირებადი ფაქტორები ეწოდებათ. ისეთ ფაქტორებს, რომლებიც არ ექვემდებარებიან რეგულირებას, ეწოდებათ არარეგულირებადი, თუმცა, მათაც გააჩნიათ გარკვეული ზემოქმედება შედეგობრივ პარამეტრებზე. არარეგულირებად ფაქტორებს მიეკუთვნებიან აგრეთვე სხვა დაუფიქსირებადი ანუ გაუთვალისწინებელი ფაქტორები. ჩვეულებრივ, ყოველი რეგულირებადი ფაქტორი წარმოდგენილია დამოუკიდებელ გრადაციებად (ჯგუფებად), რომელთა რაოდენობა დამოკიდებულია ცდის (ექსპერიმენტის) პირობებზე.

თუ რეგულირებადი ფაქტორი იწვევს მნიშვნელოვან ზეგავლენას შედეგობრივ პარამეტრზე, მაშინ ეს ზეგავლენა აისახება ჯგუფურ საშუალოებზე, რომლებიც სტატისტიკურად ერთმანეთისგან განსხვავებული იქნებიან. თითოეული ჯგუფის შიგნითაც აღინიშნება ცვალებადობა, რომელიც გამოწვეულია არარეგულირებადი ფაქტორით ან ფაქტორებით. ცვალებადობის ეს დამოკიდებულება შეიძლება ასე წარმოვადგინოთ:

$$\sigma^2 = \sigma_f^2 + \sigma_e^2,$$

სადაც, σ^2 – მთელი კომპლექსის საერთო დისპერსიის შეფასებაა; σ_f^2 – ჯგუფთაშორისო (დონეთაშორისო) დისპერსიის შეფასებაა, რომელიც გამონვეულია რეგულირებადი ფაქტორის ზეგავლენით; σ_e^2 – შიგაჯგუფური (ნარჩენი) დისპერსიის შეფასებაა, რომელიც გამონვეულია გაუთვალისწინებელი ანუ არარეგულირებადი ფაქტორებით.

რეგულირებადი ფაქტორის ზეგავლენის დასადგენად საჭიროა ფიშერის კრიტერიუმით შემონმდეს $H: \sigma_f^2 = \sigma_e^2$ ნულოვანი ჰიპოთეზა. თუ ნულოვანი ჰიპოთეზა უარყოფილია α მნიშვნელოვნების დონით, მაშინ რეგულირებადი ფაქტორის ზეგავლენა შედეგობრივ პარამეტრზე სარწმუნოა, წინააღმდეგ შემთხვევაში რეგულირებადი ფაქტორის ზეგავლენა შედეგობრივ პარამეტრზე არ შეიმჩნევა ან უმნიშვნელოა.

ამრიგად, დისპერსიული ანალიზი, გარდა საშუალოების ერთდროული შედარებისა (იხ. §10.4), შეისწავლის დაკვირვებებზე (ექსპერიმენტზე) მოქმედი სხვადასხვა ფაქტორების ზეგავლენას, მნიშვნელოვანი ფაქტორების ამორჩევასა და მათი მოქმედებების შეფასებას.

ანალიზში ჩართული ფაქტორების რაოდენობის მიხედვით დისპერსიული ანალიზი იყოფა ერთფაქტორიან, ორფაქტორიან და მრავალფაქტორიან ანალიზებად. დისპერსიული ანალიზის ჩასატარებლად საჭიროა, რომ დაკვირვებები იყოს შემთხვევითი სიდიდეები, რომელთაც გააჩნიათ ნორმალური განაწილება და ერთნაირი დისპერსია. მხოლოდ ამ შემთხვევაშია შესაძლებელი დისპერსიისა და მათემატიკური ლოდინის სარწმუნოების შეფასება და ნდობის ინტერვალების დადგენა და საერთოდ, დისპერსიული ანალიზის ჩატარება.

15.2. ერთფაქტორიანი დისპერსიული ანალიზი

ვთქვათ, გვაქვს მხოლოდ ერთი A ფაქტორი, რომელიც იყოფა m დონედ (გრადაციად). მოცემულია აგრეთვე ყოველი დონისათვის დაკვირვებათა n_i როდენობა. დავუშვათ, რომ $n_1 = n_2 =$

$\dots = n_m = n$, მაშინ ყველა დაკვირვება შეიძლება წარმოვადგინოთ შემდეგი ცხრილის სახით:

დონე	დაკვირვებები						საშ.
	1	2	...	j	...	n	
1	x_{11}	x_{12}	...	x_{1j}	...	x_{1n}	$\bar{x}_{1\bullet}$
2	x_{21}	x_{22}	...	x_{2j}	...	x_{2n}	$\bar{x}_{2\bullet}$
...
m	x_{m1}	x_{m2}	...	x_{mj}	...	x_{mn}	$\bar{x}_{m\bullet}$

ერთფაქტორიანი დისპერსიული ანალიზის შედეგების წარმოდგენის მოდელს აქვს შემდეგი სახე:

$$x_{ij} = \mu + \gamma_i + \varepsilon_{ij},$$

სადაც, x_{ij} – i დონის j დაკვირვების შედეგია;

μ – საერთო საშუალო;

γ_i – ეფექტი, რომელიც გამონვეულია i -ური დონის A ფაქტორის მიერ;

ε_{ij} – ცდომილება, გამონვეული დონის შიგნით შედეგების ვარიაციით, რომელიც გამონვეულია სხვა დაუფიქსირებელი (გაუთვალისწინებელი) ფაქტორების ზეგავლენით.

საჭიროა განისაზღვროს A ფაქტორის ზეგავლენა X ცვლადზე. ამრიგად, x_{ij} ცვლადის საერთო დისპერსია შეიძლება დავყოთ ორ ნაწილად: ერთი, რომელიც ხასიათდება A ფაქტორის ზეგავლენით და მეორე – დაუფიქსირებელი ფაქტორების ზეგავლენით. ამისათვის საჭიროა განისაზღვროს საერთო საშუალოს \bar{x} შეფასება და თოთოეული დონისთვის შესაბამისი საშუალოების $\bar{x}_{i\bullet}$ შეფასებები:

$$\bar{x} = \frac{1}{nm} \sum_{i=1}^m \sum_{j=1}^n x_{ij} = \frac{1}{m} \sum_{i=1}^m \bar{x}_{i\bullet},$$

$$\bar{x}_{i\bullet} = \frac{1}{n} \sum_{j=1}^n x_{ij}, \quad i=1,2,\dots,m.$$

განვსაზღვროთ x_{ij} და \bar{x} შორის სხვაობის კვადრატების ჯამი

$$\begin{aligned} Q &= \sum_{i=1}^m \sum_{j=1}^n (x_{ij} - \bar{x})^2 = \sum_{i=1}^m \sum_{j=1}^n (x_{ij} - \bar{x}_{i\bullet} + \bar{x}_{i\bullet} - \bar{x})^2 = \sum_{i=1}^m \sum_{j=1}^n (x_{ij} - \bar{x}_{i\bullet})^2 + \\ &+ \sum_{i=1}^m \sum_{j=1}^n (\bar{x}_{i\bullet} - \bar{x})^2 + 2 \sum_{i=1}^m \sum_{j=1}^n (x_{ij} - \bar{x}_{i\bullet})(\bar{x}_{i\bullet} - \bar{x}). \end{aligned} \quad (15.1)$$

რადგან $\sum_{j=1}^n (x_{ij} - \bar{x}_{i\bullet}) = 0$, ამიტომ

$$\sum_{i=1}^m \sum_{j=1}^n (x_{ij} - \bar{x}_{i\bullet})(\bar{x}_{i\bullet} - \bar{x}) = \sum_{j=1}^n (x_{ij} - \bar{x}_{i\bullet}) \sum_{i=1}^m (\bar{x}_{i\bullet} - \bar{x}) = 0.$$

გარდა ამისა, (15.1) გამოსახულების მეორე წევრი შეიძლება ასე წარმოვადგინოთ:

$$\sum_{i=1}^m \sum_{j=1}^n (\bar{x}_{i\bullet} - \bar{x})^2 = n \sum_{i=1}^m (\bar{x}_{i\bullet} - \bar{x})^2.$$

მაშინ (15.1) შეიძლება ასე წარმოვადგინოთ:

$$\underbrace{\sum_{i=1}^m \sum_{j=1}^n (x_{ij} - \bar{x})^2}_Q = n \underbrace{\sum_{i=1}^m (\bar{x}_{i\bullet} - \bar{x})^2}_{Q_A} + \underbrace{\sum_{i=1}^m \sum_{j=1}^n (x_{ij} - \bar{x}_{i\bullet})^2}_{Q_e},$$

ე.ი.

$$Q = Q_A + Q_e,$$

სადაც, Q_A არის დონეებს შორის გადახრების კვადრატების ჯამი, რომელსაც ზოგჯერ **დევიატას** (ლათ. სიტყვა *Devio* – გადახრა) უწოდებენ და ახასიათებს განსხვავებას დონეებს შორის. მას ხშირად უწოდებენ გაფანტვას ფაქტორების მიმართ ანუ გაფანტვას, გამოწვეულს A ფაქტორის ზეგავლენით. Q_e არის თი-

თოეულ დაკვირვებებსა და i -ურ დონეს შორის სხვაობების კვადრატების ჯამი. ამ ჯამს უწოდებენ დონეებს შიგნით გადახრების კვადრატების ჯამს და იგი ახასიათებს განსხვავებას i -ურ დონის დაკვირვებებს შორის. Q_e -ს აგრეთვე უწოდებენ ნარჩენ დისპერსიას ანუ გაფანტვას, გამოწვეულს გაუთვალისწინებელი ფაქტორების ზეგავლენით. დაბოლოს Q -ს უწოდებენ საერთო, ანუ თითოეული მონაცემების საერთო საშუალოსთან გადახრის კვადრატების ჯამს.

ამრიგად, თუ ცნობილია Q , Q_A და Q_e დევიატები, მაშინ შეიძლება შევადგინოთ შესაბამისი დასპერსიები: საერთო, დონეთაშორისო (ჯგუფთაშორისო) და დონეთაშიგნითი:

$$\sigma^2 = \frac{Q}{mn-1}, \quad \sigma_A^2 = \frac{Q_A}{m-1}, \quad \sigma_e^2 = \frac{Q_e}{m(n-1)}.$$

თუ A ფაქტორის ზეგავლენა ყველა დონისათვის ერთნაირია, მაშინ σ_A^2 და σ_e^2 არის საერთო დისპერსიის შეფასებები. იმისათვის, რომ შევამოწმოთ A ფაქტორის ზეგავლენის სარწმუნოება, საჭიროა შემოწმდეს შემდეგი ნულოვანი ჰიპოთეზა $H_0: \sigma_A^2 = \sigma_e^2$. როგორც ვიცით, ასეთი ჰიპოთეზის შემოწმებისთვის საჭიროა განისაზღვროს ფიშერის კრიტერიუმი

$$F = \frac{\sigma_A^2}{\sigma_e^2}$$

$v_1 = m - 1$ და $v_2 = m(n - 1)$ თავისუფლების ხარისხებით. თუ აღმოჩნდება, რომ $F \geq F_{\alpha; v_1, v_2}$, მაშინ ნულოვანი ჰიპოთეზა უარყოფილია და ვაკეთებთ დასკვნას A ფაქტორის მნიშვნელოვანი ზეგავლენის შესახებ. წინააღმდეგ შემთხვევაში, როცა $F < F_{\alpha; v_1, v_2}$, A ფაქტორის ზეგავლენა დაკვირვებებზე მეტად უმნიშვნელოა ან საერთოდ ადგილი არა აქვს.

ერთფაქტორიანი დისპერსიული ანალიზი უმჯობესია წარმოვადგინოთ შემდეგი ცხრილის სახით:

დის-პერსია	კვადრატების ჯამი	თავისუფ-ლების ხარისხი	დისპერ-სიების შეფასება	F ფარ-დობა	F კრიტ
ფაქტორია-ლური (დონეებს შორის)	$Q_A = n \sum_{i=1}^m (x_{i\cdot} - \bar{x})^2$	$v_1 = m-1$	$\sigma_A^2 = \frac{Q_A}{v_1}$	$F = \frac{\sigma_A^2}{\sigma_e^2}$	$F_{\alpha; v_1, v_2}$
ნარჩენი (დონეებს შიგნით)	$Q_e = \sum_{i,j} (x_{ij} - \bar{x}_{i\cdot})^2$	$v_2 = m(n-1)$	$\sigma_e^2 = \frac{Q_e}{v_2}$		
საერთო	$Q = Q_A + Q_e$	$v = mn-1$	$\sigma^2 = \frac{Q}{v}$		

მას შემდეგ, როცა დადგინდება შედეგობრივ პარამეტრზე ფაქტორის ზეგავლენა, შეგვიძლია ამ ზეგავლენის სიძლიერის განსაზღვრა, მაგალითად სნედეკორის მეთოდის გამოყენებით. მაშინ გვექნება:

$$h^2 = \frac{\tilde{\sigma}_A^2}{\tilde{\sigma}_A^2 + \sigma_e^2},$$

სადაც, $\tilde{\sigma}_A^2 = \frac{1}{n}(\sigma_A^2 - \sigma_e^2)$, ან

$$h^2 = \frac{\frac{\sigma_A^2 - \sigma_e^2}{n}}{\frac{\sigma_A^2 - \sigma_e^2}{n} + \sigma_e^2} = \frac{\sigma_A^2 - \sigma_e^2}{\sigma_A^2 + (n-1)\sigma_e^2}.$$

ინტერპრეტაციისათვის h^2 სიდიდე უმჯობესია პროცენტებში გამოვსახოთ.

მაგალითი. შესწავლილ იქნა ხორბლის ექვსი ჯიშის მოსავლიანობა. ცდები ჩატარებულ იქნა ოთხჯერ. გვანტირებებს, მოქმედებს თუ არა მოსავლიანობაზე ხორბლის სხვადასხვა ჯიში? მონაცემები მოცემულია შემდეგ ცხრილში:

ხორბლის ჯიში	მოსავლიანობა (ცენტ./ჰექტ)				საშუალო მოსავალი
	1	2	3	4	
1	26,1	29,3	30,0	27,3	28,2
2	25,0	24,3	28,5	29,0	26,7
3	27,2	26,4	31,0	26,4	27,8
4	23,6	27,2	25,2	24,8	25,2
5	0,0	33,0	36,0	29,8	32,2
6	23,0	26,0	26,0	24,8	25,0

გამოთვლების შედეგი წარმოდგენილია შემდეგ დისპერსიულ ცხრილში:

დისპერსია	კვადრა- ტების ჯამი	თავისუ- ფლების ხარისხი	დისპერ- სიების შეფასება	F ფარ- დობა	F კრიტ.
ფაქტორია- ლური (დონეებს შორის)	$Q_A = 140$	$v_1 = 5$	$\sigma_A^2 = 28$	$F = 6,4$	$F_{0,05;5;18} = 3,7$
ნარჩენი (დონეებს შიგნით)	$Q_e = 79,6$	$v_2 = 18$	$\sigma_e^2 = 4,4$		
საერთო	$Q = 219,6$	$v = 25$	$\sigma^2 = 87,8$		

რადგან $F > F_{0,05;5;18}$, ამიტომ ნულოვანი ჰიპოთეზა უარყოფილია, ე.ი. ხორბლის ჯიშებს შორის მოსავლიანობის განსხვავება სარწმუნოა. განვსაზღვროთ მოსავლიანობაზე ხორბლის ჯიშის გავლენის სიძლიერე:

$$h^2 = \frac{28 - 4,4}{28 + (4 - 1)4,4} = 0,573, \quad \text{ანუ } 57,3\%.$$

ამრიგად, ხორბლის ჯიშის ზეგავლენა მოსავლიანობაზე შეადგენს 57,3%-ს.

15.3. ორფაქტორიანი დისპერსიული ანალიზი

ვთქვათ, ექსპერიმენტალურ შედეგებზე მოქმედებს ორი A და B ფაქტორი, რომელთაც გააჩნიათ შესაბამისად r და v დონეები. დაკვირვებათა მატრიცა შეიძლება ასე წარმოვადგინოთ:

$A \backslash B$	B_1	B_2	...	B_j	...	B_v	საშ.
A_1	x_{11}	x_{12}	...	x_{1j}	...	x_{1v}	$\bar{x}_{1\bullet}$
A_2	x_{21}	x_{22}	...	x_{2j}	...	x_{2v}	$\bar{x}_{2\bullet}$
...
A_i	x_{i1}	x_{i2}	...	x_{ij}	...	x_{iv}	$\bar{x}_{i\bullet}$
...
A_r	x_{r1}	x_{r2}	...	x_{rj}	...	x_{rv}	$\bar{x}_{r\bullet}$
$\bar{x}_{\bullet j}$	$\bar{x}_{\bullet 1}$	$\bar{x}_{\bullet 2}$...	$\bar{x}_{\bullet j}$...	$\bar{x}_{\bullet v}$	\bar{x}

i -ური დონის A ფაქტორის გადაკვეთა j -ური დონის B ფაქტორთან ქმნის ij -ურ უჯრედს, სადაც ჩანერილია დაკვირვების შედეგი x_{ij} , მიღებული A და B ფაქტორების ერთდროული მოქმედების დროს. სიმარტივისთვის დავუშვათ, რომ უჯრედში გვაქვს მხოლოდ ერთი დაკვირვება და ფაქტორები ურთიერთდამოუკიდებლები არიან, ე.ი. ურთიერთქმედება გამორიცხულია. მაშინ ორფაქტორიანი დისპერსიული ანალიზის მოდელი შეიძლება ასე წარმოვადგინოთ:

$$x_{ij} = \mu + \gamma_i + g_j + e_{ij},$$

სადაც, μ – საერთო საშუალოა;

γ_i – ეფექტი, გამონვეული i -ური დონის A ფაქტორის მიერ;

g_j – ეფექტი, გამონვეული j -ური დონის B ფაქტორის მიერ;

e_{ij} – ij -ური უჯრედში შედეგების ვარიაცია.

თუ უჯრედში ერთი მნიშვნელობაა, მაშინ $e_{ij} = 0$. განვიხილოთ μ , γ და g შეფასებები. განვსაზღვროთ შემდეგი სიდიდეები:

$$\text{საერთო საშუალო: } \bar{x} = \frac{1}{r \cdot v} \sum_{i=1}^r \sum_{j=1}^v x_{ij};$$

საშუალოები დონეების მიხედვით:

$$\bar{x}_{i\bullet} = \frac{1}{v} \sum_{j=1}^v x_{ij} \quad , \quad \bar{x}_{\bullet j} = \frac{1}{r} \sum_{i=1}^r x_{ij} \quad .$$

დისპერსიების შეფასებისათვის განვიხილოთ შემდეგი გამოსახულება:

$$\begin{aligned} Q &= \sum_{i=1}^r \sum_{j=1}^v (x_{ij} - \bar{x})^2 = \sum_{i=1}^r \sum_{j=1}^v (x_{ij} - \bar{x}_{i\bullet} - \bar{x}_{\bullet j} + \bar{x} + \bar{x}_{i\bullet} - \bar{x} + \bar{x}_{\bullet j} - \bar{x})^2 = \\ &= v \underbrace{\sum_{i=1}^r (\bar{x}_{i\bullet} - \bar{x})^2}_{Q_A} + r \underbrace{\sum_{j=1}^v (\bar{x}_{\bullet j} - \bar{x})^2}_{Q_B} + \underbrace{\sum_{i=1}^r \sum_{j=1}^v (x_{ij} - \bar{x}_{i\bullet} - \bar{x}_{\bullet j} + \bar{x})^2}_{Q_e} \end{aligned}$$

ე.ი.

$$Q = Q_A + Q_B + Q_e,$$

სადაც, Q_A ახასიათებს პარამეტრის ცვლილებას, გამოწვეულს A ფაქტორის მიერ, Q_B – B ფაქტორის მიერ და Q_e – სხვა გაუთვალისწინებელი ფაქტორების მიერ. თუ ცნობილია Q , Q_A , Q_B და Q_e , მაშინ დისპერსიების შეფასებები იქნება:

$$\sigma^2 = \frac{Q}{r v - 1}; \quad \sigma_A^2 = \frac{Q_A}{r - 1}; \quad \sigma_B^2 = \frac{Q_B}{v - 1}; \quad \sigma_e^2 = \frac{Q_e}{(r - 1)(v - 1)} \quad .$$

A და B ფაქტორების ზეგავლენის დასადგენად საჭიროა დისპერსიების შედარება ფიშერის კრიტერიუმის გამოყენებით. ამისათვის უნდა განისაზღვროს

$$F_A = \frac{\sigma_A^2}{\sigma_e^2}$$

სიდიდე და იგი შედარდეს $F_{\alpha; v_1 v_3}$ კრიტიკულ მნიშვნელობას, სადაც, $v_1 = r - 1$ და $v_3 = (r - 1)(v - 1)$. თუ $F_A \geq F_{\alpha; v_1 v_3}$, მაშინ A ფაქტორის ზეგავლენა სარწმუნოა. წინააღმდეგ შემთხვევაში A ფაქტორის ზეგავლენა უმნიშვნელოა. ანალოგიურად ტარდება B ფაქტორის ზეგავლენის დადგენა.

დისპერსიული ანალიზი უმჯობესია წარმოვადგინოთ შემდეგი ცხრილების სახით:

დისპერსია	კვადრატების ჯამი	თავისუფლების ხარისხი
A ფაქტორი	$Q_A = v \sum_i (\bar{x}_{i\cdot} - \bar{x})^2$	$v_1 = r - 1$
B ფაქტორი	$Q_B = r \sum_j (\bar{x}_{\cdot j} - \bar{x})^2$	$v_2 = v - 1$
ნარჩენი	$Q_e = \sum_{i,j} (x_{ij} - \bar{x}_{i\cdot} - \bar{x}_{\cdot j} + \bar{x})^2$	$v_3 = (r-1)(v-1)$
საერთო	$Q = \sum_{i,j} (x_{ij} - \bar{x})^2$	$v = r v - 1$

დისპერსია	დისპერსიების შეფასება	F ფარდ.	F კრიტ.
A ფაქტორი	$\sigma_A^2 = \frac{Q_A}{v_1}$	$F_A = \frac{\sigma_A^2}{\sigma_e^2}$	$F_{\alpha; v_1 v_2}$
B ფაქტორი	$\sigma_B^2 = \frac{Q_B}{v_2}$	$F_B = \frac{\sigma_B^2}{\sigma_e^2}$	$F_{\alpha; v_2 v_3}$
ნარჩენი	$\sigma_e^2 = \frac{Q_e}{v_3}$		
საერთო	$\sigma^2 = \frac{Q}{v}$		

ჩვენ განვიხილეთ კერძო შემთხვევა, როცა უჯრედში იყო ერთი მონაცემი და ფაქტორებს შორის ურთიერთქმედება გამო-რიცხული იყო. ზოგადად, უჯრედში შეიძლება იყოს რამდენიმე, როგორც თანაბარი, ისე არათანაბარი რაოდენობის მონაცემები და ფაქტორებს შორის შეიძლება ადგილი ჰქონდეს ურ-თიერთქმედებას. სასურველია, რომ უჯრედებში მონაცემები იყ-ვნენ ერთი და იგივე რაოდენობის.

ამრიგად, ზოგადი შემთხვევისთვის – ორფაქტორიანი დის-პერსიული ანალიზის დროს – ერთი დაკვირვება შეიძლება წარ-მოვადგინოთ შემდეგნაირად

$$x_{ijk} = \mu + \gamma_i + g_j + v_{ij} + e_{ijk} ,$$

სადაც, μ – საერთო საშუალოა;

γ_i – ეფექტი, გამონვეული i -ური დონის A ფაქტორის ზეგავლენით;

g_j – ეფექტი, გამონვეული j -ური დონის B ფაქტორის ზეგავლენით;

v_{ij} – ეფექტი, გამონვეული A და B ფაქტორების ურთიერთქმედების შედეგად მიღებული ზეგავლენით;

e_{ij} – უჯრედშიგა ვარიაცია.

თუ უჯრედებში ერთნაირი რაოდენობის მონაცემებია, მაშინ დაკვირვების მატრიცა შეიძლება ასე წარმოვადგინოთ:

$A \backslash B$	B_1	B_2	...	B_v	$\bar{x}_{i\bullet\bullet}$
A_1	$\bar{x}_{11\bullet}$ $x_{111}, x_{112}, \dots, x_{11n}$	$\bar{x}_{12\bullet}$ $x_{121}, x_{122}, \dots, x_{12n}$...	$\bar{x}_{1v\bullet}$ $x_{1v1}, x_{1v2}, \dots, x_{1vn}$	$\bar{x}_{1\bullet\bullet}$
A_2	$\bar{x}_{21\bullet}$ $x_{211}, x_{212}, \dots, x_{21n}$	$\bar{x}_{22\bullet}$ $x_{221}, x_{222}, \dots, x_{22n}$...	$\bar{x}_{2v\bullet}$ $x_{2v1}, x_{2v2}, \dots, x_{2vn}$	$\bar{x}_{2\bullet\bullet}$
...
A_r	$\bar{x}_{r1\bullet}$ $x_{r11}, x_{r12}, \dots, x_{r1n}$	$\bar{x}_{r2\bullet}$ $x_{r21}, x_{r22}, \dots, x_{r2n}$...	$\bar{x}_{rv\bullet}$ $x_{rv1}, x_{rv2}, \dots, x_{rvn}$	$\bar{x}_{r\bullet\bullet}$
$\bar{x}_{\bullet j\bullet}$	$\bar{x}_{\bullet 1\bullet}$	$\bar{x}_{\bullet 2\bullet}$...	$\bar{x}_{\bullet v\bullet}$	\bar{x}

აქ, $x_{111}, x_{112}, \dots, x_{rvn}$ გამოსაკვლევი პარამეტრის დაკვირვებებია. დისპერსიული ანალიზის ჩატარებისათვის საჭიროა გამოვთვალოთ შემდეგი მნიშვნელობები:

– უჯრედის საშუალო მნიშვნელობა

$$\bar{x}_{ij\bullet} = \frac{1}{n} \sum_{k=1}^n x_{ijk} ;$$

– სტრიქონების (A ფაქტორი) საშუალო მნიშვნელობები

$$\bar{x}_{i\bullet\bullet} = \frac{1}{v} \sum_{j=1}^v \bar{x}_{ij\bullet} ;$$

- სვეტების (B ფაქტორი) საშუალო მნიშვნელობები

$$\bar{x}_{\bullet j \bullet} = \frac{1}{r} \sum_{i=1}^r \bar{x}_{ij \bullet};$$

- საერთო საშუალო

$$\bar{x} = \frac{1}{r v} \sum_{i=1}^r \sum_{j=1}^n \bar{x}_{ij \bullet},$$

სადაც, r, v - შესაბამისად A და B ფაქტორების დონეების რაოდენობებია. გარდა ამისა, განისაზღვრება შემდეგი კვადრატების ჯამი, ანუ დევიატები:

$$Q_A = v n \sum_{i=1}^r (\bar{x}_{i \bullet \bullet} - \bar{x})^2; \quad Q_B = m \sum_{j=1}^v (\bar{x}_{\bullet j \bullet} - \bar{x})^2; \quad Q_e = \sum_{i=1}^r \sum_{j=1}^v \sum_{k=1}^n (x_{ijk} - \bar{x}_{ij \bullet})^2;$$

$$Q_{AB} = n \sum_{i=1}^r \sum_{j=1}^v (\bar{x}_{ij \bullet} - \bar{x}_{i \bullet \bullet} - \bar{x}_{\bullet j \bullet} + \bar{x})^2; \quad Q = \sum_{i=1}^r \sum_{j=1}^v \sum_{k=1}^n (x_{ijk} - \bar{x})^2;$$

ე.ი.

$$Q = Q_A + Q_B + Q_{AB} + Q_e,$$

სადაც, Q_A -სა და Q_B -ს გააჩნია იგივე მნიშვნელობები, რაც წინა შემთხვევის დროს. Q_{AB} - არის კვადრატების ჯამი, რომელიც აფასებს A და B ფაქტორების ურთიერთქმედებას, Q_e - კვადრატების ჯამი, რომელიც აფასებს უჯრედშიგა ვარიაციას.

ამის შემდეგ უნდა განისაზღვროს დისპერსიების შეფასებები:

$$\sigma^2 = \frac{Q}{r v n - 1}; \quad \sigma_A^2 = \frac{Q_A}{r - 1}; \quad \sigma_B^2 = \frac{Q_B}{v - 1};$$

$$\sigma_{AB}^2 = \frac{Q_{AB}}{(v-1)(r-1)}; \quad \sigma_e^2 = \frac{Q_e}{r v (n-1)}$$

და სათანადო ფიშერის კრიტერიუმები:

$$F_A = \frac{\sigma_A^2}{\sigma_e^2}; \quad F_B = \frac{\sigma_B^2}{\sigma_e^2} \quad \text{და} \quad F_{AB} = \frac{\sigma_{AB}^2}{\sigma_e^2}.$$

შესაბამისი ნულოვანი ჰიპოთეზების შემოწმებით დგინდება A , B და AB ფაქტორების ზეგავლენის ან უმნიშვნელო ზეგავლენის ფაქტები.

ორფაქტორიანი დისპერსიული ანალიზი უმჯობესია წარმოვადგინოთ შემდეგი ცხრილების სახით:

დისპერსია	კვადრატების ჯამი	თავისუფლების ხარისხი
A ფაქტორი	$Q_A = v n \sum_i (\bar{x}_{i..} - \bar{x})^2$	$v_1 = r - 1$
B ფაქტორი	$Q_B = m \sum_j (\bar{x}_{.j.} - \bar{x})^2$	$v_2 = v - 1$
AB ფაქტორები	$Q_{AB} = n \sum_{i,j} (\bar{x}_{ij.} - \bar{x}_{i..} - \bar{x}_{.j.} + \bar{x})^2$	$v_3 = (v - 1)(r - 1)$
ნარჩენი	$Q_e = \sum_{i,j,k} (x_{ijk} - \bar{x}_{ij.})^2$	$v_4 = r v (n - 1)$
საერთო	$Q = \sum_{i,j,k} (x_{ijk} - \bar{x})^2$	$v = r v n - 1$

დისპერსია	დისპერსიების შეფასება	F ფარდობა	F კრიტიკ.
A ფაქტორი	$\sigma_A^2 = \frac{Q_A}{v_1}$	$F_A = \frac{\sigma_A^2}{\sigma_e^2}$	$F_{\alpha; v_1 v_4}$
B ფაქტორი	$\sigma_B^2 = \frac{Q_B}{v_2}$	$F_B = \frac{\sigma_B^2}{\sigma_e^2}$	$F_{\alpha; v_2 v_4}$
AB ფაქტორები	$\sigma_{AB}^2 = \frac{Q_{AB}}{v_3}$	$F_{AB} = \frac{\sigma_{AB}^2}{\sigma_e^2}$	$F_{\alpha; v_3 v_4}$
ნარჩენი	$\sigma_e^2 = \frac{Q_e}{v_4}$		
საერთო	$\sigma^2 = \frac{Q}{v}$		

შედეგობრივ პარამეტრზე ამა თუ იმ ფაქტორის ან ფაქტორების ერთობლივი ზემოქმედების სიძლიერე შეიძლება განისაზღვროს შემდეგი ფორმულებით:

$$h_A^2 = \frac{\hat{\sigma}_A^2}{\sigma_y^2}; \quad h_B^2 = \frac{\hat{\sigma}_B^2}{\sigma_y^2}; \quad h_{AB}^2 = \frac{\hat{\sigma}_{AB}^2}{\sigma_y^2},$$

სადაც,

$$\hat{\sigma}_A^2 = \frac{\sigma_A^2 - \sigma_e^2}{vn}; \quad \hat{\sigma}_B^2 = \frac{\sigma_B^2 - \sigma_e^2}{rn}; \quad \hat{\sigma}_{AB}^2 = \frac{\sigma_{AB}^2 - \sigma_e^2}{n};$$

$$\sigma_y^2 = \hat{\sigma}_A^2 + \hat{\sigma}_B^2 + \hat{\sigma}_{AB}^2,$$

σ_e^2 – ნარჩენი დისპერსია.

თუ რომელიმე რეგულირებადი ფაქტორის ან ფაქტორთა ერთობლივი ზეგავლენა შედეგობრივ პარამეტრზე არ დასტურდება, მაშინ მისი შესაბამისი კომპონენტი σ_y^2 გამოსახულებაში უნდა გამოირიცხოს.

მაგალითი. გამოკვლეულ იქნა სამი ტიპის მიკროელემენტის ზეგავლენა ძროხის რძის ცხიმოვანობაზე. ექსპერიმენტი ჩატარდა ერთნაირი ასაკის ოთხი სხვადასხვა ჯიშის ცხოველთა ჯგუფზე. მონაცემები მოყვანილია შემდეგ ცხრილში:

ძროხ. ჯიში	რძის ცხიმოვანობა %-ში								
	A ₁			A ₂			A ₃		
B ₁	2,1	2,0	3,4	2,8	2,6	3,0	2,4	2,1	2,8
B ₂	4,0	3,2	4,1	3,9	4,1	4,5	3,0	3,9	4,2
B ₃	3,0	2,8	2,7	3,5	4,0	2,8	4,8	3,1	2,9
B ₄	3,4	3,0	2,9	3,0	2,9	3,0	3,3	2,8	3,0

აქ *A*-თი აღნიშნულია მიკროელემენტები, ხოლო *B*-თი – სხვადასხვა ჯიშის ძროხები. *A* ფაქტორის გრადაციების რაოდენობაა $r = 3$, ხოლო *B* ფაქტორის – $v = 4$. დისპერსიული ანალიზის შედეგები მოყვანილია შემდეგ ცხრილში:

დისპერ- სია	კვადრატე- ბის ჯამი	თავისუ- ფლების ხარისხი	დისპერ- სიების შეფასება	F ფარ- დობა	F კრიტიკ.
A ფაქ- ტორი	$Q_A = 0,51$	$\nu_1 = 2$	$\sigma_A^2 = 0,26$	$F_A = 1,0$	$F_{0,05;2,24}=3,4$
B ფაქ- ტორი	$Q_B = 7,94$	$\nu_2 = 3$	$\sigma_B^2 = 2,65$	$F_B = 10,2$	$F_{0,05;3,24}=3,0$
AB ფაქ- ტორები	$Q_{AB} = 1,10$	$\nu_3 = 6$	$\sigma_{AB}^2 = 0,18$	$F_{AB} = 1,4$	$F_{0,05;6,24}=3,8$
ნარჩენი	$Q_e = 6,23$	$\nu_4 = 24$	$\sigma_e^2 = 0,26$		
საერთო	$Q = 15,78$	$\nu = 35$			

როგორც ამ ცხრილიდან ჩანს, მხოლოდ B ფაქტორის დროს ხდება ჰიპოთეზის უარყოფა. ეს იმას ნიშნავს, რომ ამ ჯიშის ცხოველებს გააჩნიათ რძის ცხიმოვანობაზე მიდრეკილება, რომელიც ალბათ შთამომავლობით უნდა აიხსნას და ამიტომ მიკროეულემენტების ზეგავლენას იგი არ ექვემდებარება. ალბათ, ამიტომაც, რომ A და B ფაქტორების ურთიერთქმედებაც ვერ ახდენს რძის ცხიმოვანობაზე ზეგავლენას.

რადგან განხილული ორი ფაქტორიდან მხოლოდ B ფაქტორი (ძროხის ჯიში) მოქმედებს რძის ცხიმოვანობაზე, ამიტომ შეგვიძლია რძის ცხიმოვანობაზე ძროხის ჯიშის ზემოქმედების სიძლიერის განსაზღვრა. ამისათვის გვაქვს: $\sigma_B^2 = 2,65$; $\sigma_e^2 = 0,26$; $n = 3$; $r = 3$. თუ გამოვიყენებთ ზემოთ მოყვანილ ფორმულებს, მაშინ მივიღებთ:

$$\hat{\sigma}_B^2 = \frac{2,65 - 0,26}{3,3} = 0,266, \quad h_B^2 = \frac{0,266}{0,266 + 0,26} = 0,50, \quad \text{ანუ } 50\%.$$

ამრიგად, რძის ცხიმოვანობაზე ძროხის ჯიშის ზემოქმედების სიძლიერე 50%-ის ტოლია.

ანალოგიურად ტარდება სამი, ოთხი და ზოგადად, მრავალფაქტორიანი დისპერსიული ანალიზი.

15.4. დისპერსიული ანალიზის არაპარამეტრული მეთოდები

თუ საწყის მონაცემებს არ გააჩნიათ ნორმალური განაწილება, მაშინ დისპერსიული ანალიზის ჩასატარებლად, საჭიროა გამოვიყენოთ არაპარამეტრული მეთოდები. თუ დონეების დაკვირვებათა რაოდენობები ტოლია, მაშინ შეიძლება გამოვიყენოთ **ფრიდმანის რანგული დისპერსიული ანალიზი**. ამისათვის საჭიროა ყოველი დონისათვის დაკვირვებების რანჟირება და რანგების დადგენა. მაშინ ფრიდმანის კრიტერიუმს გააჩნია შემდეგი სახე:

$$\chi_R^2 = \frac{12}{mn(m+1)} \sum_{j=1}^n \left(\sum_{i=1}^m R_{ij} \right)^2 - 3n(m+1),$$

სადაც, m – დონეების (სტრიქონების) რაოდენობაა, n – დაკვირვებათა (სვეტების) რაოდენობა, $\sum_i R_{ij}$ – i -ური სვეტის რანგების ჯამი.

უნდა შევამოწმოთ ჰიპოთეზა: არის თუ არა დონეებს შორის განსხვავება. რადგან χ_R^2 სიდიდეს გააჩნია χ^2 განაწილება, ამიტომ კრიტიკული წერტილი $\chi_{\alpha;v}^2$ მოიძებნება χ^2 განაწილების ცხრილიდან α და $v = m - 1$ სიდიდეების საშუალებით. თუ $\chi_R^2 \geq \chi_{\alpha;v}^2$, მაშინ ნულოვანი ჰიპოთეზა უარყოფილია, ე.ი. განსხვავება დონეებს შორის სარწმუნოა, წინააღმდეგ შემთხვევაში – იგი არ არის სარწმუნო.

მაგალითი. ჯანმრთელი კბილების მქონე ბავშვთა სამ ასაკობრივ ჯგუფში განისაზღვრა ჰიგიენური ინდექსი, რომელიც პირობით ერთეულებშია წარმოდგენილი

დონე დაკვ.	3 წლის		4 წლის		5 წლის	
	A_1	R_1	A_2	R_2	A_3	R_3
1	1	1	3,0	2,5	3,0	2,5
2	1	1	1,2	2	1,3	3
3	1	1	2,0	2	2,2	3
4	1	1,5	1,0	1,5	1,2	3
5	2,6	2	2,7	3	1,3	1
6	2,7	2	1,3	1	3,0	3
Σ		8,5		12,0		15,5

$m = 3; n = 6; \chi_R^2 = \frac{12}{3 \cdot 6 \cdot 4} (8,5^2 + 12,0^2 + 15,5^2) - 3 \cdot 6 \cdot 4 = 4,08; \alpha =$
 $= 0,05; v = m - 1 = 2; \chi_{0,05;2}^2 = 5,99.$ რადგან $4,08 < 5,99$, ამიტომ
 ბავშვთა ასაკობრივ ჯგუფებს შორის ჰიგიენური ინდექსით გან-
 სხვავება არ შეიმჩნევა.

როდესაც დონეებში დაკვირვებათა რაოდენობა სხვადა-
 სხვაა, მაშინ შეგვიძლია გამოვიყენოთ **კრასკელ-უოლისის**
კრიტერიუმი

$$H = \frac{12}{N(N+1)} \sum_{j=1}^n \frac{1}{n_i} \left[\sum_{i=1}^m R_{ij} \right]^2 - 3(N+1),$$

სადაც, N – დაკვირვებათა საერთო რაოდენობაა, n_i – i -ური დონის
 დაკვირვებათა რაოდენობა. $\sum_i R_{ij}$ – i -ური სვეტის რანგების ჯამი,
 რომლებიც მიიღება გაერთიანებული მწკრივის რანჟირებისას.

H სიდიდეს გააჩნია χ^2 განაწილება $v = m-1$ თავისუფლების
 ხარისხით. m არის დონეების რაოდენობა. როცა $H < \chi_{\alpha;v}^2$, მაშინ
 ნულოვანი ჰიპოთეზა მიიღება, ე.ი. დონეებს შორის განსხვავება არ
 შეიმჩნევა, წინააღმდეგ შემთხვევაში, როცა $H \geq \chi_{\alpha;v}^2$ – განსხვავე-
 ბა სარწმუნოა.

მაგალითი. მოცემულია ბავშვების 5 ჯგუფისათვის
 სისხლის ნაკადის სიჩქარე. ჯგუფები დაყოფილია ავადმყოფობის
 სიმძიმის მიხედვით (უძიმესი – A_1 ჯგუფი).

დონე \ დაკვ.	A_1		A_2		A_3		A_4		A_5	
	x	R_1	x	R_2	x	R_3	x	R_4	x	R_5
1	30	22	36	25	25	18	14	2,5	19	10
2	17	7,5	32	23,5	17	7,5	17	7,5	16	4,5
3	22	12,5	42	26	24	16,5			27	19
4	32	23,5	22	12,5	11	1			23	14,5
5	24	16,5			14	2,5			28	20,5
6					28	20,5			23	14,5
7					17	7,5			20	11
8					44	27			16	4,5
Σ		82,0		87,0		100,5		10,0		98,5

$$H = \frac{12}{27 \cdot 28} \left(\frac{82^2}{5} + \frac{87^2}{4} + \frac{100,5^2}{8} + \frac{10^2}{2} + \frac{98,5^2}{8} \right) - 3 \cdot 28 = 7,47.$$

$\alpha = 0,05$, $\nu = 5 - 1 = 4$, $\chi_{\alpha;\nu}^2 = 9,49$. რადგან $7,47 < 9,49$, ამიტომ დონეებს შორის განსხვავება არ შეიმჩნევა.

16. მთავარი კომპონენტების მეთოდი

პრაქტიკაში ძალიან ხშირად გვხვდება ისეთი სიტუაცია, როცა პარამეტრების რაოდენობა ძალზე დიდია და, მიუხედავად ამისა, საჭიროა საწყისი მონაცემების სტატისტიკური დამუშავება და გარკვეული გადაწყვეტილების მიღება. აქედან გამომდინარე, საჭიროა საწყისი ინფორმაციის შეკუმშული სახით წარმოდგენა, ანუ შესასწავლი ობიექტის აღწერა მცირე რაოდენობის განზოგადებული მაჩვენებლებით, მაგალითად, მთავარი კომპონენტებით ან ფაქტორებით. მთავარი კომპონენტები წარმოადგენენ მეტად მოსახერხებელ გამსხვილებულ მაჩვენებლებს, რომლებიც ასახავენ ობიექტის (პროცესის) იმ შინაგან კანონზომიერების აღწერას, რაც შეუძლებელია დაკვირვებების საშუალებით.

მთავარი კომპონენტების მეთოდით შესაძლებელია შემდეგი ამოცანების გადაწყვეტა.

1. შესასწავლ მოვლენაში ობიექტურად არსებული ფარული კანონზომიერების გამოვლენა;

2. შესასწავლი პროცესის აღწერა მცირე რაოდენობის მთავარი კომპონენტებით, რომელთა რიცხვი გაცილებით ნაკლებია საწყისი ცვლადების რაოდენობაზე. ამ შემთხვევაში, მთავარი კომპონენტები პროცესს ადეკვატურად ასახავენ უფრო კომპაქტური ფორმით და შეიცავენ საშუალოდ უფრო მეტ ინფორმაციას, ვიდრე უშუალოდ გაზომვადი ცვლადები;

3. ცვლადების მთავარ კომპონენტებთან სტატისტიკური კავშირის გამოვლენა და შესწავლა, რაც საშუალებას იძლევა უფრო აქტიურად ვიმოქმედოთ პროცესზე მისი ეფექტური ფუნქციონირებისთვის;

4. პროცესის განვითარების ტენდენციის პროგნოზირება რეგრესიის განტოლებით, რომელიც აგებულია მთავარი კომპონენტების საშუალებით. პროგნოზირების ასეთ მეთოდს გააჩნია გარკვეული უპირატესობა კლასიკურ რეგრესიულ ანალიზთან შედარებით, განსაკუთრებით იმ შემთხვევის დროს, როცა საქმე გვაქვს მულტიკოლინეარობის პრობლემასთან.

მთავარი კომპონენტების მოდელი. ვთქვათ, მოცემულია n შემთხვევითი ცვლადები X_1, X_2, \dots, X_n , რომლებსაც გააჩნიათ $\mu = (\mu_1, \mu_2, \dots, \mu_n)'$ საშუალოების ვექტორი და კოვარიაციული მატრიცა $S_{n,n}$. საჭიროა, განისაზღვროს ამ ცვლადების ურთიერთკავშირი, ანუ სტრუქტურული დამოკიდებულება. სტრუქტურული დამოკიდებულების ერთ-ერთ მეთოდს წარმოადგენს მთავარი კომპონენტების მეთოდი, რომლის ძირითადი ამოცანა შეიძლება ასე ჩამოვყალიბოთ: უნდა მოიძებნოს საწყისი X_1, X_2, \dots, X_n ცვლადების ისეთი წრფივი კომბინაცია

$$Y_i = \sum_{j=1}^n a_{ij} x_{ij} = \sum_{j=1}^n a_{ij} X_j, \quad i = 1, 2, \dots, m \quad (16.1)$$

სადაც, a_{ij} – წარმოადგენს წონით კოეფიციენტებს, როცა სრულდება შემდეგი პირობები:

$$\text{corr}(Y_i Y_j) = 0, \quad i, j = 1, 2, \dots, n, \quad i \neq j$$

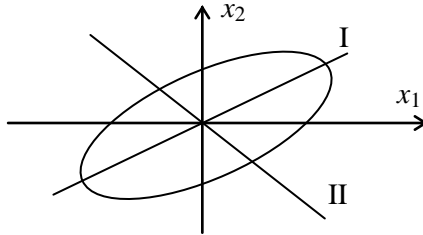
$$\text{var}(Y_1) \geq \text{var}(Y_2) \geq \dots \geq \text{var}(Y_m);$$

$$\sum_{i=1}^n \text{var}(Y_i) = \sum_{i=1}^n s_{ii}.$$

როგორც ამ ფორმულიდან ჩანს, ახალი ცვლადები Y_1, Y_2, \dots, Y_m , რომლებსაც მთავარ კომპონენტებს უწოდებენ, ერთმანეთის მიმართ არიან არაკორელირებული და რანჟირებული – დისპერსიის კლებაობის მიხედვით. უნდა აღინიშნოს, რომ ჯამური დისპერსია გარდაქმნის შემდეგ არ იცვლება. აქედან გამომდინარე, Y_i ცვლადების პირველ q ქვესიმრავლეზე მოდის საერთო დისპერსიის ძირითადი ნაწილი და ამიტომ შესაძლებელი ხდება საწყისი ცვლადების დამოკიდებულების სტრუქტურის შეკვეცილი აღწერა.

იმისათვის, რომ უკეთ გავერკვეთ მთავარი კომპონენტების მეთოდის არსში, განვიხილოთ მისი **გეომეტრიული ინტერპრეტაცია**. დავუშვათ, რომ ორი X_1 და X_2 შემთხვევითი ცვლადი ნორმა-

ლურად არის განანილებული $\mu = (\mu_1, \mu_2)$ საშუალოების ვექტორითა და S კოვარიაციული მატრიცით. ამ განანილების სიმკვრივის ელიფსოიდი, რომლის ცენტრი მოთავსებულია (μ_1, μ_2) წერტილში, წარმოდგენილია შემდეგ ნახაზზე:



პირველ მთავარ კომპონენტს $Y_1 = \alpha_{11}x_1 + \alpha_{12}x_2$ შეესაბამება ელიფსოიდის I ღერძი, რომლის სიგრძის ნახევარი ტოლია $\sqrt{\lambda_1}$ სიდიდისა, სადაც, λ_1 არის კოვარიაციული მატრიცის მაქსიმალური საკუთრივი მნიშვნელობა. რადგან X გააჩნია ორგანზომილებიანი ნორმალური განანილება, ამიტომ მას გააჩნია აგრეთვე II პატარა ღერძი, რომელიც I ღერძის პერპენდიკულარულია და იგი შეესაბამება მეორე მთავარ კომპონენტს $Y_2 = \alpha_{21}x_1 + \alpha_{22}x_2$, რომლის სიგრძე პროპორციულია $\sqrt{\lambda_2}$ სიდიდისა. ამრიგად, ელიფსოიდის, როგორც I, ასევე II ღერძი, განისაზღვრება $X\alpha_1$ და $X\alpha_2$ სიდიდეებით, სადაც, α_1 და α_2 საკუთრივი ვექტორებია, რომლებიც შეესაბამება S მატრიცის λ_1 და λ_2 საკუთრივ მნიშვნელობებს.

თუ X_1 და X_2 ერთმანეთის მიმართ დადებით კორელაციურ დამოკიდებულებაშია, მაშინ რაც უფრო იზრდება კორელაცია ცვლადებს შორის, მით უფრო უახლოვდება ელიფსოიდი I წრფეს და თუ X_1 და X_2 შორის დამოკიდებულება ფუნქციონალურია, მაშინ ელიფსოიდი გარდაიქმნება წრფედ.

მთავარი კომპონენტების განსაზღვრა. მთავარი კომპონენტის მეთოდი მდგომარეობს α_{ij} კოეფიციენტების მოძებნაში. (16.1) გამოსახულება მატრიცული სახით ასე ჩაიწერება $Y = X\alpha$. მოცემული α -ს დროს ამ გამოსახულების დისპერსია ტოლია:

$$\text{var}(Y) = \text{var}(X\alpha) = \alpha' \text{var}(X) \alpha = \alpha' S \alpha,$$

სადაც, S მოცემული სანყისი ცვლადების კოვარიაციული მატრიცაა. თუ სანყისი მონაცემები ნორმირებულია, მაშინ გვექნება კორელაციური მატრიცა

$$R = \frac{1}{n-1} XX'$$

მთავარი კომპონენტების მეთოდის პირველი ამოცანა მდგომარეობს Y_1 კომპონენტის მოძებნაში, რომელსაც გააჩნია უდიდესი დისპერსია. სათანადო შეზღუდვების შემოტანის გარეშე, ზოგადად, ამ ამოცანის გადაწყვეტა შეუძლებელია. მაგალითად, თუ ფიქსირებულ α -თვის რაღაც c მუდმივის დროს მივიღებთ $\alpha^* = c\alpha$ სიდიდეს, სადაც c -ს ზრდასთან ერთად უსასრულოდ იზრდება დისპერსიაც. ეს მოვლენა თავიდან რომ ავიცილოთ, საჭიროა α ვექტორის ნორმირება ისე, რომ

$$\alpha' \alpha = \alpha_1^2 + \alpha_2^2 + \dots + \alpha_n^2 = 1.$$

მაშინ ამოცანა ჩამოყალიბდება შემდეგნაირად: მოვახდინოთ $\alpha' S \alpha$ გამოსახულების მაქსიმიზაცია, როცა $\alpha' \alpha = 1$. დავუშვათ

$$\varphi = \alpha' S \alpha - \lambda(\alpha' \alpha - 1),$$

სადაც, λ – ლაგრანჟის მამრავლია. φ -ის კერძო წარმოებულის ნულთან გატოლების შემდეგ

$$\frac{\partial \varphi}{\partial \alpha} = 2S\alpha - 2\lambda\alpha = 0$$

მივიღებთ შემდეგ განტოლებას:

$$(S - \lambda I)\alpha = 0,$$

რომელიც წარმოადგენს კლასიკური ტიპის განტოლებას და მას გააჩნია ამონახსნი მხოლოდ იმ შემთხვევაში, როცა

$$\det |S - \lambda I| = 0.$$

ამრიგად, მიღებული განტოლების ამოხსნისთვის საჭიროა მოიძებნოს S მატრიცის $\lambda_1 \geq \lambda_2 \geq \dots \geq \lambda_n$ მახასიათებელი ფესვები.

იმისათვის, რომ განვსაზღვროთ, ამ ფესვებიდან რომელი გამოვიყენოთ მახასიათებელი (საკუთრივი) ვექტორის შესარჩევად, რომელიც მოახდენს $\alpha' S \alpha$ გამოსახულების მაქსიმიზაციას, საჭიროა, განტოლების მარცხენა მხარე გავამრავლოთ α' -ზე, მაშინ, თუ მხედველობაში მივიღებთ $\alpha' S \alpha = 1$, გვექნება:

$$\alpha'(S - \lambda I)\alpha = \alpha' S \alpha - \lambda = 0, \quad \text{ე.ი.} \quad \alpha' S \alpha = \lambda,$$

მაგრამ ჩვენ ვიცით, რომ $\alpha'S\alpha = \text{var}(Y)$. ე.ი. λ წარმოადგენს დისპერსიას.

ამრიგად, იმისათვის, რომ მოვახდინოთ Y კომპონენტის დისპერსიის მაქსიმიზაცია, უნდა ავიღოთ მახასიათებელი ფესვებიდან უდიდესი მნიშვნელობის ფესვი, კერძოდ λ_1 და მისი შესაბამისი საკუთრივი ვექტორი α_1 . მაშინ პირველი მთავარი კომპონენტი მიიღებს შემდეგ სახეს:

$$Y_1 = X\alpha_1,$$

რომლის დისპერსია იქნება λ_1 .

ზოგადად, როცა საქმე გვაქვს n ცვლადთან, პირველი მთავარი კომპონენტი Y_1 წარმოადგენს n ცვლადების წრფივ კომბინაციას, რომელთა კოეფიციენტები ტოლია ნორმირებული საკუთრივი ვექტორის კომპონენტებისა, რომელიც, თავის მხრივ, შეესაბამება R ან S მატრიცის უდიდეს საკუთრივ მნიშვნელობას.

მეორე მთავარი კომპონენტი Y_2 წარმოადგენს საწყისი n ცვლადების წრფივ კომბინაციას კოეფიციენტებით, რომლებიც ნორმირებული საკუთრივი ვექტორის კომპონენტების ტოლია და იგი შეესაბამება λ_2 მახასიათებელ მნიშვნელობას, რომელიც წარმოადგენს λ_1 -ის შემდეგ უდიდეს მნიშვნელობას. ანალოგიურად განისაზღვრება მესამე მთავარი კომპონენტი და ა.შ. n -ური კომპონენტის ჩათვლით. თითოეული კომპონენტის დისპერსია შესაბამისად ტოლია მახასიათებელი λ_i , $i=1,2,\dots,n$ ფესვებისა და თითოეული კომპონენტი არ არის ერთმანეთის მიმართ დამოკიდებული. ვაჩვენოთ ეს ბოლო დამოკიდებულება.

ცნობილია თეორემა, რომლის თანახმად, ნებისმიერი სიმეტრიული S მატრიცისთვის არსებობს ისეთი ორთოგონალური მატრიცა α , რომლისთვისაც სრულდება შემდეგი ტოლობა:

$$\alpha'S\alpha = \begin{bmatrix} \lambda_1 & 0 & 0 & \dots & 0 \\ 0 & \lambda_2 & 0 & \dots & 0 \\ \dots & \dots & \dots & \dots & \dots \\ 0 & 0 & 0 & \dots & \lambda_n \end{bmatrix}. \quad (16.2)$$

ამასთან, თუ S დადებითად განსაზღვრული მატრიცაა, მაშინ ყველა $\lambda_i > 0$ და $|S| > 0$. რადგან α წარმოადგენს S მატრიცის საკუთრივ ვექტორებს და როგორც მიღებული (18.2) გამოსახ-

ულებიდან ჩანს, $\text{cov}(\alpha_{ij}) = 0$, $i = 1, 2, \dots, n, j = 1, 2, \dots, n, i \neq j$, ამიტომ კომპონენტები ურთიერთდამოუკიდებელი არიან.

თითოეული კომპონენტის ჯამური დისპერსიის ფარდობითი წილი განისაზღვრება შემდეგი გამოსახულებით:

$$q_i = \frac{\lambda_i}{\sum_{i=1}^n \lambda_i}, \quad i = 1, 2, \dots, n. \quad (16.3)$$

პრაქტიკულად, თუ ჯამური დისპერსიის 80-85% მოდის პირველი k რაოდენობის კომპონენტებზე, მაშინ დანარჩენი კომპონენტები $k + 1, k + 2, \dots, n$ შეიძლება მხედველობაში არ მივიღოთ და ამით მოვახდინოთ სივრცის განზომილების შემცირება. აქედან გამომდინარე, საჭიროა ჩამოვყალიბოთ ისეთი კრიტერიუმი, რომელიც ნაკლებდისპერსიანი კომპონენტების ანალიზიდან გამორიცხვის საშუალებას მოგვცემს.

მთავარი კომპონენტის მეთოდის გამოყენება მიზანშეწონილია იმ შემთხვევაში, როცა საწყის ცვლადებს გააჩნიათ საერთო ფიზიკური ბუნება და იზომებიან ერთი და იმავე ფიზიკურ ერთეულებში. თუ ცვლადები სხვადასხვა ფიზიკური ბუნებისაა, მაშინ აუცილებელია მათი ნორმირება, მაგალითად, შემდეგნაირად:

$$z_{ij} = \frac{x_{ij} - \bar{x}_j}{\sigma_j}, \quad i = 1, 2, \dots, m, \quad j = 1, 2, \dots, n,$$

სადაც, \bar{x}_j საშუალო მნიშვნელობებია, ხოლო σ_j – საშუალო კვადრატული გადახრები. ასეთი ნორმირების შემდეგ კოვარიაციული მატრიცა გარდაიქმნება კორელაციურ მატრიცად.

კორელაცია X_i ცვლადსა და Y_j მთავარ კომპონენტს შორის განისაზღვრება შემდეგნაირად:

$$\text{corr}(Y_i X_j) = \frac{\alpha_{ij} \sqrt{\lambda_i}}{\sigma_j}, \quad (16.4)$$

სადაც, σ_i – სტანდარტული გადახრაა. ამრიგად, თუ გვინდა შევადაროთ თითოეული X_j ცვლადის წილი Y_i კომპონენტის ფორმირებაში, საჭიროა შევადაროთ $\frac{\alpha_{ij}}{\sigma_j}$ მნიშვნელობები. თუ კორელაციური მატრიცა ცნობილია, მაშინ საკმარისია α_{ij} კოეფიციენტების

შედარება. კერძოდ, იმ X_j ცვლადს, რომელსაც გააჩნია უდიდესი α_{ij} კოეფიციენტი, მიუძღვის უდიდესი წილი Y_i მთავარი კომპონენტის ფორმირებაში.

მთავარი კომპონენტების ინტერპრეტაციისთვის, კერძოდ, მასში არსებული ინფორმაციის გამოსავლენად, ხშირად (16.1) გამოსახულებაში a_{ij} წონითი კოეფიციენტების მაგივრად იყენებენ კორელაციის კოეფიციენტებს, რომლებიც განისაზღვრება (16.4) ფორმულით.

ჰიპოთეზების შემოწმება. დავუშვათ, რომ კოვარიაციული (კორელაციური) მატრიცის მახასიათებელი რიცხვები ერთმანეთის ტოლია, ანუ (16.3) გამოსახულებიდან მიღებული q_i მნიშვნელობები ერთმანეთის ტოლია. ორი X_1 და X_2 ცვლადების დროს ეს ნიშნავს ელიფსოიდის გარდაქმნას წრედ, ე.ი. ორივე I და II ღერძები ერთმანეთის ტოლია. მრავალგანზომილებიანი სისტემის დროს საქმე გვაქვს სფეროსთან. ამრიგად, თუ ნულოვანი ჰიპოთეზა

$$H_0: \lambda_1 = \lambda_2 = \dots = \lambda_n$$

სამართლიანია, მაშინ მუდმივი სიმკვრივის ელიფსოიდი გარდაქმნება მუდმივი სიმკვრივის სფეროდ და მაშინ ნულოვანი ჰიპოთეზის შემოწმება შეიძლება გავიხილოთ როგორც სფეროზე შემოწმების ჰიპოთეზად. ასეთი ჰიპოთეზის შესამოწმებლად განვიხილოთ შემდეგი სტატისტიკა:

$$\chi^2 = -(m-1) \left[\sum_{i=1}^n \ln \lambda_i - n \ln \left(\frac{1}{n} \sum_{i=1}^n \lambda_i \right) \right], \quad (16.5)$$

რომელსაც გააჩნია χ^2 განაწილება $v=0,5n(n+1)-1$ თავისუფლების ხარისხით. თუ $\chi^2 < \chi_{\alpha;v}^2$, მაშინ ნულოვანი ჰიპოთეზა მიიღება.

დავუშვათ, რომ პირველ k მთავარ კომპონენტზე მოდის ჯამური დისპერსიის უდიდესი ნაწილი. ჩვენ გვინტერესებს, დარჩენილი კომპონენტები განსხვავდებიან თუ არა ერთმანეთისგან. თუ ისინი არ განსხვავდებიან, მაშინ მათი გამორიცხვა შემდგომი ანალიზიდან მიზანშეწონილია. ამრიგად, ჩამოვაცალიბოთ შემდეგი ნულოვანი ჰიპოთეზა:

$$H_0: \lambda_{k+1} = \lambda_{k+2} = \dots = \lambda_n,$$

მაშინ (16.5) სტატისტიკას აქვს შემდეგი სახე:

$$\chi^2 = -(m-1) \left[\sum_{j=k+1}^n \ln \lambda_j - q \ln \left(\frac{1}{q} \sum_{j=K+1}^n \lambda_j \right) \right],$$

სადაც, $q = n - k$. თუ $\chi^2 < \chi_{\alpha;v}^2$, მაშინ ნულოვანი ჰიპოთეზა მიიღება და ბოლო q რაოდენობის მთავარი კომპონენტები შეიძლება გამოვრიცხოთ ანალიზიდან.

უნდა გვახსოვდეს, რომ ამ ჰიპოთეზების შემოწმების კრიტერიუმები მეტად მგრძნობიარეა ნორმალური განაწილების მიმართ. განსაკუთრებით საეჭვოა მათი გამოყენება დროითი მწკრივების (მაგალითად, ბიოსიგნალების) მიმართ, რადგან მონაცემების დამოუკიდებლობა ამ შემთხვევაში იშვიათობას წარმოადგენს.

მთავარი კომპონენტების მეთოდი გამოიყენება იმ შემთხვევაში, როდესაც საწყისი ცვლადები ურთიერთდამოკიდებული არიან. როცა ცვლადები ურთიერთდამოუკიდებელია, მაშინ ამ მეთოდის გამოყენებას აზრი არა აქვს, რადგან ამ დროს ფაქტიურად ხდება საწყისი ცვლადების რანჟირება დისპერსიების კლების მიხედვით. აქედან გამომდინარე, სასურველია შევამოწმოთ პარამეტრების დამოუკიდებლობის ჰიპოთეზა. ამისათვის განვიხილოთ სტატისტიკა:

$$\gamma = - \left(m - \frac{2n+1}{6} \right) \ln |R|,$$

სადაც, $|R|$ — კორელაციური მატრიცის დეტერმინანტია, რომელიც შეიძლება ასე განისაზღვროს: $|R| = \prod_{i=1}^n \lambda_i$, სადაც, λ_i კორელაციური მატრიცის საკუთრივი მნიშვნელობებია.

$$\gamma \text{ სტატისტიკას გააჩნია } \chi^2 \text{ განაწილება } v = \frac{m(m-1)}{2}$$

თავისუფლების ხარისხით. თუ აღმოჩნდება, რომ $\gamma < \chi_{\alpha;v}^2$, მაშინ პარამეტრები დამოუკიდებელია, წინააღმდეგ შემთხვევაში, როცა $\gamma \geq \chi_{\alpha;v}^2$, ისინი დამოკიდებული არიან.

მაგალითი. შესწავლილ იქნა მექანიკური ტრავმის 2300 შემთხვევის ავადმყოფობის ისტორია [7]. სტატისტიკური დამუშავებისთვის გამოყენებულ იქნა შემდეგი 11 მაჩვენებელი:

1. მდგომარეობის სიმძიმის სუბიექტური შეფასება (დამაკმაყოფილებელი, საშუალო სიმძიმის, მძიმე და უკიდურესად მძიმე);
2. ცნობიერების მდგომარეობა (ნათელი, არეული, არ ჰქონდა);
3. არტერიული წნევის სიდიდე;
4. ჰიპერტონიის ხანგრძლივობა (არტერიული წნევის სიდიდე 100მმ ვერცხლის სვეტის სიმაღლეზე ნაკლებია);
5. პულსის სიხშირე;
6. სისხლკარგვის სიდიდე;
7. სიცოცხლისთვის მნიშვნელოვანი დაზიანებული ორგანოების რაოდენობა;
8. გადასხმული სისხლის რაოდენობა;
9. გადასხმული სისხლის შემცველის რაოდენობა;
10. ოპერატიული ჩარევის ხანგრძლივობა;
11. ოპერატიული ჩარევის რაოდენობა.

ჩვენ შემთხვევაში კორელაციურ მატრიცას აქვს შემდეგი სახე:

$$R = \begin{matrix} & \begin{matrix} 1 & 2 & 3 & 4 & 5 & 6 & 7 & 8 & 9 & 10 & 11 \end{matrix} \\ \begin{matrix} 1 \\ 2 \\ 3 \\ 4 \\ 5 \\ 6 \\ 7 \\ 8 \\ 9 \\ 10 \\ 11 \end{matrix} & \begin{bmatrix} 1 & 0,714 & -0,624 & 0,426 & 0,525 & 0,748 & 0,741 & 0,503 & 0,588 & 0,434 & 0,507 \\ & 1 & -0,558 & 0,354 & 0,326 & 0,533 & 0,672 & 0,315 & 0,454 & 0,265 & 0,428 \\ & & 1 & -0,600 & -0,382 & 0,706 & 0,674 & 0,538 & 0,601 & 0,296 & 0,386 \\ & & & 1 & 0,304 & 0,467 & 0,378 & 0,510 & 0,436 & 0,364 & 0,387 \\ & & & & 1 & 0,482 & 0,406 & 0,351 & 0,341 & 0,273 & 0,293 \\ & & & & & 1 & 0,750 & 0,691 & 0,747 & 0,419 & 0,539 \\ & & & & & & 1 & 0,529 & 0,617 & 0,368 & 0,553 \\ & & & & & & & 1 & 0,737 & 0,430 & 0,515 \\ & & & & & & & & 1 & 0,421 & 0,556 \\ & & & & & & & & & 1 & 0,605 \\ & & & & & & & & & & 1 \end{bmatrix} \end{matrix}$$

მიღებული კორელაციური მატრიცისათვის განისაზღვრა საკუთრივი მნიშვნელობები და საკუთრივი ვექტორები, რომელთა საშუალებით მივიღეთ მთავარი კომპონენტების კოეფიციენტების მნიშვნელობები, რომლებიც მიღებულია (16.4) ფორმულით.

პირველი ოთხი მთავარი კომპონენტის საკუთრივი მნიშვნელობები და კოეფიციენტები მოცემულია შემდეგ ცხრილში:

ცვლადები	მთავარი კომპონენტები			
	1	2	3	4
1	0,84	-0,27	0,21	0,08
2	0,69	-0,44	0,24	-0,30
3	-0,79	0,23	0,32	0,09
4	0,62	0,15	-0,44	0,04
5	0,57	-0,20	0,20	0,76
6	0,86	0,06	-0,09	0,04
7	0,84	-0,25	0,12	-0,17
8	0,76	0,30	-0,13	0,11
9	0,81	0,15	-0,19	-0,04
10	0,59	0,58	0,36	-0,02
11	0,72	0,40	0,32	-0,16
დისპერსიის აბსოლუტური მნიშვნელობა (λ)	6,087	1,058	0,849	0,754
კომპონენტის წილი მთელ დისპერსიაში (%)	55,3	9,6	7,7	6,9
ჯამური წილი (%)	55,3	64,3	72,0	78,9

როგორც ამ ცხრილიდან ჩანს, პირველ ოთხ მთავარ კომპონენტზე მოდის მთელი დისპერსიის დაახლოებით 80% და აქედან მხოლოდ პირველ კომპონენტზე მოდის 55%.

განვიხილოთ უფრო დეტალურად პირველი მთავარი კომპონენტის $Z_1 = 0,84 X_1 + 0,69 X_2 - 0,79 X_3 + \dots + 0,72 X_{11}$ კოეფიციენტები, რომლებიც წარმოადგენენ კორელაციურ კოეფიციენტებს სანყის მაჩვენებლებთან.

1. პირველი მთავარი კომპონენტის დადებითი კორელაცია მდგომარეობის სიმძიმის სუბიექტურ შეფასებასთან გასაგებია, რადგან რაც უფრო მძიმეა ტრავმა, მით უფრო მაღალია სუბიექტური შეფასება.

2. რაც უფრო მძიმეა ტრავმა, მით უფრო გამოკვეთილია ცნობიერების დარღვევა, რომელიც ბალური სისტემით არის შეფასებული. ამასთან, უდიდესი ქულა ენიჭებოდა იმ შემთხვევაში,

როცა პაციენტს ცნობიერება არ გააჩნდა. სწორედ აქედან გამომდინარეობს მეორე მაჩვენებლის დადებითი კორელაცია პირველ კომპონენტთან.

3. რაც უფრო მძიმეა ტრავმა, მით უფრო მცირეა არტერიული წნევა, რაც ობიექტურად დასტურდება უარყოფითი კორელაციით მესამე მაჩვენებლისა პირველ კომპონენტთან.

4. ხანგრძლივი ჰიპოტონია მძიმე ტრავმების დამახასიათებელი ნიშანია, რასაც მიუთითებს მეოთხე მაჩვენებლის დადებითი კორელაცია პირველ მთავარ კომპონენტთან.

5-6. მძიმე ტრავმის დროს პულსის სიხშირე მატულობს, ხოლო სისხლკარგვა ხელს უწყობს ორგანიზმის მდგომარეობის გაართულებას. ამით აიხსნება პირველი კომპონენტის დადებითი კორელაცია მე-5 და მე-6 მაჩვენებლებთან.

7. რაც უფრო მეტი სიცოცხლისთვის მნიშვნელოვანი ორგანოებია დაზიანებული, მით უფრო დიდია ტრავმის სიმძიმის ხარისხი. აქედან გამომდინარეობს დადებითი კორელაცია მე-7 და პირველ მთავარ კომპონენტებს შორის.

8-11. ტრავმის სიმძიმის ზრდასთან ერთად იზრდება ჩარევის რაოდენობებიც. ამიტომაც ამ მაჩვენებლების დადებითი კორელაცია პირველ მთავარ კომპონენტთან.

17. ფაქტორული ანალიზი

ფაქტორული ანალიზის მიზანია მარტივი სტრუქტურის დადგენა, რომელიც გამოავლენს შესასწავლ მოვლენაში ობიექტურად არსებულ ფარულ კანონზომიერებას. გარდა ამისა, ფაქტორული ანალიზი საშუალებას იძლევა განისაზღვროს ახალი ცვლადები ანუ ფაქტორები და ისეთი სიდიდეების შეფასება, რომელთა უშუალო გაზომვა შეუძლებელია. ფაქტორული ანალიზით შეგვიძლია სანყისი მონაცემების გარდაქმნა, განზომილების შემცირება და სხვა სპეციფიკური ამოცანების გადაწყვეტა.

ფაქტორული ანალიზის ძირითადი მოდელი. ვთქვათ გვაქვს m ობიექტი, რომლებიც აღწერილი არიან X_1, X_2, \dots, X_n პარამეტრებით. ფაქტორული ანალიზის ჩასატარებლად ინფორმაცია წარმოდგენილი უნდა იყოს $m \times n$ განზომილებიანი მატრიცის სახით. იმისათვის, რომ გამოირიცხოს სხვადასხვა ფიზიკურ ერთეულებში

გაზომილი პარამეტრების ეფექტი, საჭიროა საწყისი მონაცემები წარმოვადგინოთ სტანდარტიზირებული სახით, ე.ი. გადავიღეთ უგანზომილებო ცვლადებზე

$$Z_{ij} = \frac{x_{ij} - \bar{x}_i}{\sigma_i}, \quad i = 1, 2, \dots, m, \quad j = 1, 2, \dots, n.$$

მაშინ ფაქტორული ანალიზის ძირითადი მოდელი შეიძლება ასე წარმოვადგინოთ:

$$Z_i = \sum_{j=1}^n \lambda_{ij} F_j + e_i, \quad i = 1, 2, \dots, m, \quad (17.1)$$

სადაც, Z_i – i -ური შემთხვევითი ცვლადია, F_1, F_2, \dots, F_n – ზოგადი ფაქტორებია, რომლებიც შემთხვევით სიდიდეებს წარმოადგენენ და გააჩნიათ ნორმალური განაწილება; e_i – სპეციფიური ანუ მახასიათებელი ფაქტორებია, რომლითაც ხასიათდებიან თითოეული საწყისი Z_i ცვლადები (იგულისხმება, რომ მახასიათებელი ფაქტორები არაკორელირებული არიან); λ_{ij} – ფაქტორული დატვირთვებია, რომლითაც ხასიათდება თითოეული ფაქტორი და რომლებიც უნდა იყვნენ განსაზღვრულნი.

ამრიგად, ფაქტორული ანალიზის ძირითადი ამოცანაა, განისაზღვროს ფაქტორული დატვირთვები. თუ ზოგადი და მახასიათებელი ფაქტორები ურთიერთარაკორელირებული არიან, მაშინ i -ური ფაქტორის დისპერსია შეიძლება ასე წარმოვადგინოთ:

$$\sigma_i^2 = 1 = \lambda_{i1}^2 + \lambda_{i2}^2 + \dots + \lambda_{in}^2 + \tau_i, \quad i = 1, 2, \dots, m,$$

სადაც, $\lambda_i^2 - Z_i$ პარამეტრის დისპერსიის წილია, რომელიც მოდის i -ურ ფაქტორზე, ხოლო მახასიათებელი ვექტორისათვის $\text{var}(e_i) = \tau_i, \quad i = 1, 2, \dots, n$, სადაც, τ_i -ს უწოდებენ სპეციფიურ დისპერსიას. შემოვიღოთ აღნიშვნა

$$h_i^2 = \lambda_{i1}^2 + \lambda_{i2}^2 + \dots + \lambda_{in}^2,$$

რომელიც წარმოადგენს საერთო დისპერსიის წილს, გამოწვეულს ზოგადი ფაქტორებით და მას უწოდებენ i -ური საწყისი პარამეტრის ერთიანობას. k -ური ფაქტორის სრული წილი საერთო დისპერსიაში იქნება:

$$V_k = \sum_{j=1}^m \lambda_{jk}^2, \quad k = 1, 2, \dots, n.$$

ფაქტორების დამოუკიდებლობის შემთხვევაში, ადვილია პარამეტრებს შორის კორელაციის კოეფიციენტის განსაზღვრა

$$r'_{ik} = \lambda_{i1}\lambda_{k1} + \lambda_{i2}\lambda_{k2} + \dots + \lambda_{in}\lambda_{kn}, \quad i \neq k, k = 1, 2, \dots, n.$$

შემოვიტანოთ ნარჩენი კორელაციისა და ნარჩენი კორელაციური მატრიცის ცნებები. (17.1) მოდელის ჩანერისათვის საწყის ინფორმაციას წარმოადგენს სპირმენის კორელაციური მატრიცა. თუ გამოვიყენებთ ფაქტორულ მოდელს და ხელახლა გამოვთვლით კორელაციურ მატრიცას, მაშინ მათ შორის სხვაობა იქნება ნარჩენი კორელაციის კოეფიციენტი, ე.ი.

$$\bar{r}_{jk} = r_{jk} - r'_{jk},$$

ხოლო ნარჩენი კორელაციის კოეფიციენტებისგან შედგენილი მატრიცა

$$\bar{R} = R - R'.$$

ფაქტორული ანალიზის ამოცანაა, შეაფასოს λ_{ij} ფაქტორული დატვირთვების τ_i სპეციფიური დისპერსიები და ფაქტორული მნიშვნელობები. როდესაც ფაქტორული დატვირთვები ცნობილი იქნება, შემდეგ რჩება კიდევ ერთი ამოცანა, კერძოდ, ფაქტორების ინტერპრეტაციის პრობლემა. ამისათვის გამოიყენება ფაქტორული ბრუნვა.

მთავარი ფაქტორების განსაზღვრა. ფაქტორული ანალიზის პირველი ამოცანაა R კორელაციური (ან S კოვარიაციული) მატრიცის საშუალებით განისაზღვროს λ_{ij} ფაქტორული დატვირთვების l_{ij} შეფასებები და τ_i სპეციფიურ დისპერსიის t_i შეფასებები. ამისათვის არსებობს მრავალი მეთოდი, მაგრამ ჩვენ განვიხილავთ მთავარი ვექტორის განსაზღვრის მეთოდს, რომელიც მდგომარეობს n მთავარი კომპონენტის განსაზღვრაში:

$$Y_i = \sum_{j=1}^n a_{ij} X_j, \quad i = 1, 2, \dots, n.$$

გავიხსენოთ, რომ n მთავარი კომპონენტები ერთმანეთის მიმართ არაკორელირებული არიან და i -ური კომპონენტის დისპერსია $\text{var}(Y_i)$ ტოლია კორელაციური მატრიცის i -ური საკუთრივი მნიშვნელობისა, ე.ი. $\text{var}(Y_i) = \lambda_i$.

მთავარი ფაქტორების მეთოდიდან გამომდინარე, ზოგად ფაქტორებად მიიღება m პირველი მთავარი კომპონენტი, რომლებიც განისაზღვრება შემდეგნაირად:

$$F_j = \frac{Y_j}{\sqrt{\text{var}(Y_j)}}, \quad j = 1, 2, \dots, m,$$

ხოლო ფაქტორული დატვირთვების შეფასებები ტოლია:

$$l_{ij} = a_{ji} \sqrt{\text{var}(Y_j)}, \quad i = 1, 2, \dots, n, \quad j = 1, 2, \dots, m.$$

რაც შეეხება მახასიათებელი ფაქტორების შეფასებებს, ისინი ასე განისაზღვრება:

$$e_i = \sum_{j=m+1}^n a_{ji} Y_j, \quad i = 1, 2, \dots, n.$$

ამრიგად, მივიღებთ ფაქტორული მოდელის შემდეგ შეფასებას:

$$Z_i = \sum_{j=1}^m e_{ij} F_j + e_i, \quad i = 1, 2, \dots, n.$$

აქ ყველა ფაქტორი ერთმანეთის მიმართ არაკორელირებულია და დისპერსიები ერთის ტოლია. სპეციფიური დისპერსიებისა და ერთიანობების შეფასებები იქნება:

$$h_i^2 = \sum_{j=1}^m a_{ji}^2 \text{var}(Y_j) = \sum_{j=1}^m e_{ji}^2;$$

$$t_i = \sum_{i=m+1}^n a_{ji} \text{var}(Y_j).$$

ფაქტორების ბრუნვა. ფაქტორული დატვირთვების განსაზღვრის შემდეგ საჭიროა თითოეული ფაქტორის ინტერპრეტაცია. ამისათვის საჭიროა ფაქტორების ბრუნვა ახალი ორთოგონალური ფაქტორების მისაღებად, რომლებიც ერთმანეთის მიმართ არაკორელირებულია და გააჩნიათ ერთეულოვანი დისპერსია. ბრუნვის შემდეგ ფაქტორული მოდელი ასე ჩაიწერება:

$$Z_i = \sum_{j=1}^m C_{ij} F^{(R)} + e_i, \quad i = 1, 2, \dots, n,$$

სადაც, C_{ij} წარმოადგენს ახალი ფაქტორების დატვირთვებს. აქვე უნდა აღვნიშნოთ, რომ ორთოგონალური ბრუნვის შედეგად, თითოეული სანყისი Z_i ცვლადის ერთიანობა უცვლელი რჩება, ე.ი.

$$h_i^2 = \sum_{j=1}^m C_{ij}^2 = \sum_{j=1}^m l_{ij}^2, \quad i = 1, 2, \dots, n.$$

C_{ij} მუდმივები ისე უნდა შეირჩეს, რომ დატვირთვები იყვნენ მარტივი სტრუქტურის. ზოგადად, ფაქტორული დატვირთვების სტრუქტურა ითვლება მარტივად, როცა C_{ij} კოეფიციენტების უმრავლესობა ახლოსაა ნულთან და მხოლოდ ერთ ან რამდენიმე მათგანს გააჩნია ნულისგან შედარებით დიდი მნიშვნელობები. ამრიგად, ბრუნვის მიზანია, თითოეული საწყისი ცვლადი წარმოადგინოს ერთი ან მცირე რაოდენობის ფაქტორებით, ხოლო სხვა დანარჩენი ფაქტორების დატვირთვა ნულთან უნდა იყოს ახლოს.

არსებობს ფაქტორების ბრუნვის მრავალი მეთოდი, როგორც გრაფიკული, ასევე ანალიზური. ანალიზური მეთოდებიდან გამოირჩევა ე.წ. მიზნობრივი ფუნქციის მინიმიზაციის მეთოდი, რომელიც დამოკიდებულია C_{ij} სიდიდეებზე. ორთოგონალური ბრუნვისთვის, ძირითადად, იყენებენ შემდეგ ფუნქციას:

$$G = \sum_{k=1}^m \sum_{\substack{j=1 \\ j \neq k}}^m \left[\sum_{i=1}^n C_{ij}^2 C_{ik}^2 - \frac{\gamma}{n} \left(\sum_{i=1}^n C_{ij}^2 \right) \left(\sum_{i=1}^n C_{ik}^2 \right) \right], \quad (17.2)$$

სადაც, $0 \leq \gamma \leq 1$.

როცა $\gamma = 0$, ბრუნვას, რომელიც მიიღება G ფუნქციის მინიმიზაციით, ეწოდება „კვარტიმაქსის“ მეთოდი. ამ შემთხვევაში, G ფუნქციის მინიმიზაცია ექვივალენტურია

$$\frac{1}{nm} \sum_{j=1}^m \sum_{i=1}^n (C_{ij}^2 - \bar{C}_{..}^2)^2 \quad (17.3)$$

გამოსახულების მაქსიმიზაციისა. აქ, $\bar{C}_{..}^2 = \frac{1}{nm} \sum_{j=1}^m \sum_{i=1}^n C_{ij}^2$.

როგორც (17.3) გამოსახულებიდან ჩანს, „კვარტიმაქსის“ მეთოდით ხდება ფაქტორული დატვირთვების კვადრატების დისპერსიის მაქსიმიზაცია. ამ შემთხვევაში, ის ფაქტორები, რომლებთაც აქვთ დიდი დატვირთვის მნიშვნელობები, კიდევ უფრო იზრდებიან, ხოლო მცირე მნიშვნელობები კიდევ უფრო მცირენი ხდებიან.

როცა $\gamma = 1$, ბრუნვის მეთოდს „ვარიმაქს“ უწოდებენ. ეს მეთოდი ყველაზე უფრო ხშირად გამოიყენება პრაქტიკაში. ამ შემთხვევაში G ფუნქციის მინიმიზაცია ექვივალენტურია

$$\frac{1}{n} \sum_{j=1}^m \sum_{i=1}^n \left(C_{ij}^2 - \bar{C}_{\bullet j}^2 \right)^2 \quad (17.4)$$

გამოსახულების მაქსიმიზაციისა. აქ,

$$\bar{C}_{\bullet j}^2 = \frac{1}{n} \sum_{i=1}^n C_{ij}^2, j = 1, 2, \dots, m.$$

(17.4) გამოსახულება წარმოადგენს ფაქტორული დატვირთვების კვადრატების დისპერსიების ჯამს სვეტების მიხედვით და იგი იწვევს თითოეული ფაქტორის დატვირთვების კვადრატების დისპერსიის მაქსიმიზაციას. ეს უკანასკნელი, თავის მხრივ, დატვირთვების დიდ მნიშვნელობებს კიდევ უფრო ზრდის, ხოლო დატვირთვების მცირე მნიშვნელობებს უფრო ამცირებს. ამრიგად, ამ შემთხვევაში უბრალო სტრუქტურა მიიღება თითოეული ფაქტორისათვის ცალ-ცალკე, ხოლო წინა „კვარტიმაქსის“ მეთოდში უბრალო სტრუქტურა განისაზღვრებოდა ყველა ფაქტორისათვის ერთდროულად.

აქამდე ჩვენ განვიხილეთ ფაქტორების ბრუნვის მხოლოდ ერთოგონალური მეთოდები. არსებობს მოსაზრება, რომ მნიშვნელოვანია მივიღოთ ფაქტორული დატვირთვების უბრალო სტრუქტურა, ვიდრე შევინარჩუნოთ მათი ერთოგონალურობა. ასეთი ფაქტორების მიღების მეთოდს ირიბკუთხა ბრუნვა ეწოდება. ალენიშნით $R = (r_{ik})$ ფაქტორების მეორადი „უბრალო“ სტრუქტურა $i = 1, 2, \dots, n$, $k = 1, 2, \dots, m$, სადაც, r_{ik} არის კორელაციის კოეფიციენტი i -ური სანყის ცვლადის k -ურ მეორად ფაქტორთან. ირიბკუთხა ბრუნვის დროს ხდება

$$G = \sum_{k=1}^m \sum_{j \neq k}^m \left[\sum_{i=1}^n r_{ij}^2 r_{ik}^2 - \frac{\gamma}{n} \left(\sum_{i=1}^n r_{ij}^2 \right) \left(\sum_{i=1}^n r_{ik}^2 \right) \right]$$

მიზნობრივი ფუნქციის მინიმიზაცია, სადაც, $r_{ik} = \text{corr}(X_i, G_j)$ და იგი იცვლება 0-დან 1-მდე. ანალიტიკურ მეთოდებს, სადაც იძებნება უბრალო მეორადი სტრუქტურა, ეწოდება არაპირდაპირი მეთოდი „ობლიმინი“. აქ, $0 \leq \gamma \leq 1$. როცა $\gamma = 0$, მაშინ საქმე გვაქვს მძლავრ ირიბკუთხა ბრუნვასთან, როცა $\gamma = 0,5$ – ნაკლებად ირიბკუთხა ბრუნვასთან, ხოლო როცა $\gamma = 1$ – ყველაზე მცირე ირიბკუთხა ბრუნვასთან.

10	0,150	99
11	0,084	99
12	0,023	100
13	0,019	100

როგორც ამ ცხრილიდან ჩანს, პირველ ხუთ ფაქტორზე მოდის მთელი დისპერსიის 80%, ამიტომ შევჩერდეთ ამ ხუთ ფაქტორზე, რომელთა ფაქტორული დატვირთვების შეფასებები მოცემულია შემდეგ ცხრილში:

ცვლა- დები	ფაქტორები				
	1	2	3	4	5
1	0,21	0,88	-0,22	0,15	-0,09
2	0,33	0,90	-0,13	0,14	-0,09
3	0,05	-0,08	0,59	0,48	0,33
4	0,47	0,83	-0,04	0,13	-0,07
5	-0,07	-0,18	-0,35	0,71	-0,03
6	-0,70	0,33	-0,10	-0,06	0,34
7	0,61	-0,44	-0,42	-0,12	-0,20
8	0,71	-0,48	-0,26	0,05	-0,21
9	-0,13	0,31	0,18	-0,59	-0,22
10	-0,61	-0,03	-0,52	-0,07	0,31
11	0,40	0,03	-0,32	-0,23	0,69
12	0,87	-0,00	0,15	-0,09	0,26
13	0,88	-0,01	0,13	-0,08	0,26

როგორც აღვნიშნეთ, ფაქტორული დატვირთვები – ესაა კორელაცია სანყის ცვლადსა და ფაქტორს შორის. მაგ. $I_{11} = 0,21$ ეს არის კორელაციის კოეფიციენტი სისტოლურ წნევასა და პირველ ფაქტორს შორის; $I_{12} = 0,88$ კი არის იგივე ცვლადისა და მეორე ფაქტორს შორის და ა.შ.

ფაქტორების ინტერპრეტაციისთვის განვიხილოთ ის დატვირთვები, რომელთა სიდიდე მეტია რაიმე ზღვრულ მნიშვნელობაზე, მაგ. $r = 0,4$. ცხრილში მოცემული პირველი ფაქტორისათვის ასეთი მაჩვენებლებია რვა, ე.ი. პირველი ფაქტორი, ძირითადად, დამოკიდებულია ამ რვა მაჩვენებელზე, მეორე ფაქტორი კი მხოლოდ ხუთ ცვლადზე და ა.შ. როგორც ვხედავთ, ფაქტორების ინტერპრეტაცია ძნელდება. ამიტომ მიზანშეწონილია ჩავატაროთ

ფაქტორების ბრუნვა „ვარიმაქსის“ მეთოდის გამოყენებით. მიღებული შედეგები მოყვანილია ცხრილში.

ცვლა- დები	ფ ა ქ ტ ო რ ე ბ ი				
	1	2	3	4	5
1	-0,11	0,94	-0,99	0,02	0,03
2	-0,01	0,98	-0,00	-0,04	0,04
3	-0,08	-0,10	0,81	0,17	0,04
4	0,10	0,95	0,09	-0,08	0,09
5	0,03	-0,00	-0,00	0,81	-0,14
6	-0,85	0,02	-0,14	0,01	0,07
7	0,78	-0,13	-0,36	0,14	0,19
8	0,88	-0,12	-0,14	0,21	0,14
9	-0,15	0,14	-0,19	-0,67	-0,14
10	-0,61	-0,21	-0,49	0,27	0,18
11	0,08	0,07	-0,09	0,02	0,88
12	0,64	0,21	0,32	-0,16	0,54
13	0,65	0,21	0,30	-0,15	0,54

მიღებული ფაქტორული დატვირთვებით უკვე შესაძლებელია ფაქტორების ინტერპრეტაცია. კერძოდ, F_1 ფაქტორი წარმოადგენს სისხლდენის ფაქტორს, F_2 – არტერიული წნევის, F_3 – გულის შეკუმშვის სიხშირისა და პლაზმის, F_4 – დიურეზის და F_5 – სისხლის შედგენილობის.

ამრიგად, საწყისი 13 ცვლადის მაგივრად ჩვენ შევჩერდით ხუთ ძირითად ფაქტორზე, რომლებიც ერთმანეთის მიმართ არაკორელირებული არიან და შესაძლებელია მათი ინტერპრეტაცია.

სტანდარტიზირებული ნორმალური განაწილების

ფუნქციის $F(z) = \frac{1}{\sqrt{2\pi}} \int_{-\infty}^z e^{-\frac{t^2}{2}} dt$ მნიშვნელობები

z	0	1	2	3	4	5	6	7	8	9
0,0	0,5000	5040	5080	5120	5160	5199	5239	5279	5319	5359
0,1	5398	5438	5478	5517	5557	5596	5636	5675	5714	5754
0,2	5793	5832	5871	5910	5948	5987	6026	6064	6103	6141
0,3	6179	6217	6255	6293	6331	6368	6406	6443	6480	6517
0,4	6554	6591	6628	6664	6700	6736	6772	6808	6844	6879
0,5	6915	6950	6985	7020	7054	7088	7123	7157	7190	7224
0,6	7258	7291	7324	7357	7389	7422	7454	7486	7518	7549
0,7	7580	7612	7642	7673	7704	7734	7764	7794	7823	7852
0,8	7881	7910	7939	7967	7996	8023	8051	8079	8106	8133
0,9	8159	8186	8212	8238	8264	8289	8315	8340	8365	8389
1,0	0,8413	8438	8461	8485	8508	8531	8554	8577	8599	8621
1,1	8643	8665	8686	8708	8729	8749	8770	8790	8810	8830
1,2	8849	8869	8888	8907	8925	8944	8962	8980	8997	9015
1,3	9032	9049	9066	9082	9099	9115	9131	9147	9162	9177
1,4	9192	9207	9222	9236	9251	9265	9279	9292	9306	9319
1,5	9332	9345	9357	9370	9382	9394	9407	9418	9430	9441
1,6	9452	9463	9474	9485	9495	9505	9515	9525	9535	9545
1,7	9554	9564	9573	9582	9591	9599	9608	9616	9625	9633
1,8	9641	9649	9656	9664	9671	9678	9686	9693	9700	9706
1,9	9713	9720	9726	9732	9738	9744	9750	9756	9762	9767
2,0	0,9773	9778	9783	9788	9793	9798	9803	9808	9812	9817
2,1	9821	9826	9830	9834	9838	9842	9846	9850	9854	9857
2,2	9861	9865	9868	9871	9875	9878	9881	9884	9887	9890
2,3	9893	9896	9898	9901	9904	9906	9909	9911	9913	9916
2,4	9918	9920	9922	9925	9927	9929	9931	9932	9934	9936
2,5	9938	9940	9941	9943	9945	9946	9948	9949	9951	9952
2,6	9953	9955	9956	9957	9959	9960	9961	9962	9963	9964
2,7	9965	9966	9967	9968	9969	9970	9971	9972	9973	9974
2,8	9974	9975	9976	9977	9977	9978	9979	9980	9980	9981
2,9	9981	9982	9983	9983	9984	9984	9985	9985	9986	9986
3,0	0,9987	9987	9987	9988	9988	9989	9989	9989	9990	9990
3,1	9990	9991	9991	9991	9992	9992	9992	9992	9993	9993
3,2	9993	9993	9994	9994	9994	9994	9994	9995	9995	9995
3,3	9995	9995	9996	9996	9996	9996	9996	9996	9996	9997
3,4	9997	9997	9997	9997	9997	9997	9997	9997	9998	9998
3,5	9998	9998	9998	9998	9998	9998	9998	9998	9998	9998
3,6	9998	9999	9999	9999	9999	9999	9999	9999	9999	9999

სტანდარტიზირებული ნორმალური განაწილების

სიმკვრივის ფუნქციის $f(z) = \frac{1}{\sqrt{2\pi}} e^{-\frac{z^2}{2}}$ მნიშვნელობები

z	0	1	2	3	4	5	6	7	8	9
0,0	0,3989	3989	3989	3988	3986	3984	3982	3980	3977	3973
0,1	3970	3965	3961	3956	3951	3945	3939	3932	3925	3918
0,2	3910	3902	3894	3885	3876	3867	3857	3847	3836	3825
0,3	3814	3802	3790	3778	3765	3752	3739	3726	3712	3697
0,4	3683	3668	3653	3637	3621	3605	3589	3572	3555	3538
0,5	3521	3503	3485	3467	3448	3429	3410	3391	3372	3352
0,6	3332	3312	3292	3271	3251	3230	3209	3187	3166	3144
0,7	3123	3101	3079	3056	3034	3011	2989	2966	2943	2920
0,8	2897	2874	2850	2827	2803	2780	2756	2732	2709	2685
0,9	2661	2637	2613	2589	2565	2541	2516	2492	2468	2444
1,0	0,2420	2396	2371	2347	2323	2299	2275	2251	2227	2203
1,1	2179	2155	2131	2107	2083	2059	2036	2012	1989	1965
1,2	1942	1919	1895	1872	1849	1826	1804	1781	1758	1736
1,3	1714	1691	1669	1647	1626	1604	1582	1561	1539	1518
1,4	1497	1476	1456	1435	1415	1394	1374	1354	1334	1315
1,5	1295	1276	1257	1238	1219	1200	1182	1163	1145	1127
1,6	1109	1092	1074	1057	1040	1023	1006	0989	0973	0957
1,7	0904	0925	0909	0893	0878	0863	0848	0833	0818	0804
1,8	0790	0775	0761	0748	0734	0721	0707	0694	0681	0669
1,9	0656	0644	0632	0620	0608	0596	0584	0573	0562	0551
2,0	0,0540	0529	0519	0508	0498	0488	0478	0468	0459	0449
2,1	0440	0431	0422	0413	0404	0396	0387	0379	0371	0363
2,2	0355	0347	0339	0332	0325	0317	0310	0303	0297	0290
2,3	0283	0277	0270	0264	0258	0252	0246	0241	0235	0229
2,4	0224	0219	0213	0208	0203	0198	0194	0189	0184	0180
2,5	0175	0171	0167	0163	0158	0154	0151	0147	0143	0139
2,6	0136	0132	0129	0126	0122	0119	0116	0113	0110	0107
2,7	0104	0101	0099	0096	0093	0091	0088	0086	0084	0081
2,8	0079	0077	0075	0073	0071	0069	0067	0065	0063	0061
2,9	0060	0058	0056	0055	0053	0051	0050	0048	0047	0046
3,0	0,0044	0043	0042	0040	0039	0038	0037	0036	0035	0034
3,1	0033	0032	0031	0030	0029	0028	0027	0026	0025	0025
3,2	0024	0023	0022	0022	0021	0020	0020	0019	0018	0018
3,3	0017	0017	0016	0016	0015	0015	0014	0014	0013	0013
3,4	0012	0012	0012	0011	0011	0010	0010	0010	0009	0009
3,5	0009	0008	0008	0008	0008	0007	0007	0007	0007	0006
3,6	0006	0006	0006	0005	0005	0005	0005	0005	0005	0004
3,7	0004	0004	0004	0004	0004	0004	0003	0003	0003	0003
3,8	0003	0003	0003	0003	0003	0002	0002	0002	0002	0002
3,9	0002	0002	0002	0002	0002	0002	0002	0002	0001	0001

ათობითი ლოგარითმის მნიშვნელობები

<i>n</i>	5	6	7	8	9	10	11	12	13
<i>lgn</i>	0,699	0,778	0,845	0,903	0,954	1,000	1,041	1,079	1,114

<i>n</i>	14	15	16	17	18	19	20	21	22
<i>lgn</i>	1,146	1,176	1,204	1,230	1,255	1,279	1,301	1,322	1,342

<i>n</i>	23	24	25	26	27	28	29	30	31
<i>lgn</i>	1,362	1,380	1,398	1,415	1,430	1,447	1,462	1,477	1,491

<i>n</i>	32	33	34	35	36	37	38	39	40
<i>lgn</i>	1,505	1,518	1,532	1,544	1,556	1,568	1,580	1,591	1,602

<i>n</i>	41	42	43	44	45	46	47	48	49
<i>lgn</i>	1,613	1,623	1,634	1,644	1,653	1,663	1,672	1,681	1,690

<i>n</i>	50	51	52	53	54	55	56	57	58
<i>lgn</i>	1,699	1,708	1,716	1,724	1,732	1,740	1,748	1,756	1,763

<i>n</i>	59	60	65	70	75	80	85	90	95
<i>lgn</i>	1,771	1,778	1,813	1,845	1,875	1,903	1,929	1,954	1,978

<i>n</i>	100
<i>lgn</i>	2,000

Q კრიტიკული მნიშვნელობები

<i>m</i>	2	3	4	5	6	7	8	9
$\alpha = 0,05$	1,96	2,39	2,64	2,81	2,94	3,04	3,12	3,20
$\alpha = 0,01$	2,58	2,94	3,14	3,29	3,40	3,49	3,57	3,64

<i>m</i>	10	11	12	13	14	15	16	17
$\alpha = 0,05$	3,26	3,32	3,37	3,41	3,46	3,49	3,53	3,56
$\alpha = 0,01$	3,69	3,74	3,79	3,83	3,87	3,90	3,94	3,97

<i>m</i>	18	19	20	21	22	23	24	25
$\alpha = 0,05$	3,59	3,62	3,65	3,68	3,70	3,72	3,74	3,77
$\alpha = 0,01$	3,99	4,02	4,04	4,07	4,09	4,11	4,13	4,15

ლაპლასის $\Phi(U_i) = \frac{2}{\sqrt{2\pi}} \int_0^{u_i} e^{-\frac{x^2}{2}} dx$ ფუნქციის მნიშვნელობები

U_i	0	1	2	3	4	5	6	7	8	9
0,0	0,0000	0080	0160	0239	0319	3999	0478	0558	0638	0717
0,1	0797	0878	0955	1034	1113	1192	1271	1350	1428	1507
0,2	1585	1663	1741	1819	1897	1974	2051	2128	2205	2282
0,3	2358	2434	2510	2586	2661	2737	2812	2886	2960	3035
0,4	3108	3182	3255	3328	3401	3473	3545	3616	3688	3759
0,5	3829	3899	3969	4039	4108	4177	4245	4313	4381	4448
0,6	4515	4581	4647	4713	4778	4843	4907	4971	5035	5098
0,7	5161	5223	5285	5346	5407	5467	5527	5587	5646	5705
0,8	5763	5821	5878	5935	5991	6047	6102	6157	6211	6265
0,9	6319	6372	6424	6476	6528	6579	6629	6679	6729	6778
1,0	0,6827	6875	6923	6970	7017	7063	7109	7154	7199	7243
1,1	7287	7330	7373	7415	7457	7499	7540	7580	7620	7660
1,2	7699	7737	7775	7813	7850	7887	7923	7959	7994	8029
1,3	8064	8098	8132	8165	8198	8230	8262	8293	8324	8355
1,4	8358	8415	8444	8473	8501	8529	8557	8584	8611	8638
1,5	8664	8690	8715	8740	8764	8789	8812	8836	8859	8882
1,6	8904	8926	8948	8969	8990	9011	9031	9051	9070	9090
1,7	9109	9127	9146	9164	9181	9199	9216	9233	9249	9265
1,8	9281	9297	9312	9327	9342	9357	9371	9385	9399	9412
1,9	9421	9439	9451	9464	9476	9488	9500	9512	9523	9534
2,0	0,9545	9556	9566	9576	9586	9596	9606	9616	9625	9634
2,1	9643	9651	9660	9668	9676	9684	9692	9700	9707	9715
2,2	9722	9729	9736	9743	9749	9756	9762	9768	9774	9780
2,3	9786	9791	9797	9802	9807	9812	9817	9822	9827	9832
2,4	9836	9841	9845	9849	9853	9857	9861	9865	9869	9872
2,5	9876	9879	9883	9886	9889	9892	9895	9898	9901	9904
2,6	9907	9910	9912	9915	9917	9920	9922	9924	9926	9928
2,7	9931	9933	9935	9937	9939	9940	9942	9944	9946	9947
2,8	9949	9951	9952	9953	9955	9956	9958	9959	9960	9961
2,9	9963	9964	9965	9966	9967	9968	9969	9970	9971	9972
3,0	0,9973	9974	9975	9976	9976	9977	9978	9979	9979	9980
3,1	9981	9981	9982	9983	9983	9984	9984	9985	9985	9986
3,2	9986	9987	9987	9988	9988	9989	9989	9989	9990	9990
3,3	9990	9991	9991	9991	9992	9992	9992	9992	9993	9993
3,4	9993	9994	9994	9994	9994	9994	9995	9995	9995	9995
3,5	9995	9996	9996	9996	9996	9996	9996	9996	9997	9997
3,6	9997	9997	9997	9997	9997	9997	9997	9998	9998	9998
3,7	9998	9998	9998	9998	9998	9998	9998	9998	9998	9998
3,8	9999	9999	9999	9999	9999	9999	9999	9999	9999	9999
3,9	9999	9999	9999	9999	9999	9999	9999	9999	9999	9999

სტიუდენტის განაწილება

α	ცალმხრივი კრიტერიუმი								
	v	0,30	0,20	0,10	0,05	0,025	0,01	0,005	0,001
1	0,727	1,376	3,078	6,314	12,71	31,82	63,66	318,3	
2	0,617	1,061	1,886	2,920	4,303	6,965	9,925	22,33	
3	0,584	0,978	1,638	2,353	3,182	4,541	5,841	10,22	
4	0,569	0,941	1,533	2,132	2,776	3,747	4,604	7,173	
5	0,559	0,906	1,440	1,943	2,447	3,143	3,707	5,208	
6	0,553	0,920	1,476	2,015	2,571	3,365	5,032	5,893	
7	0,549	0,896	1,415	1,895	2,365	2,998	3,499	4,785	
8	0,546	0,889	1,397	1,860	2,306	2,896	3,355	4,501	
9	0,543	0,883	1,383	1,833	2,262	2,821	3,250	4,297	
10	0,542	0,879	1,372	1,812	2,228	2,764	3,169	4,144	
11	0,540	0,876	1,363	1,796	2,201	2,718	3,106	4,025	
12	0,539	0,873	1,356	1,782	2,179	2,681	3,055	3,930	
13	0,538	0,870	1,350	1,771	2,160	2,650	3,012	3,852	
14	0,537	0,868	1,345	1,761	2,145	2,624	3,977	3,787	
15	0,536	0,866	1,341	1,753	2,131	2,602	2,947	3,733	
16	0,535	0,865	1,337	1,746	2,120	2,583	2,921	3,686	
17	0,534	0,863	1,333	1,740	2,110	2,567	2,898	3,646	
18	0,534	0,862	1,330	1,734	2,101	2,552	2,878	3,611	
19	0,533	0,861	1,328	1,729	2,093	2,539	2,861	3,579	
20	0,533	0,860	1,325	1,725	2,086	2,528	2,845	3,552	
21	0,532	0,859	1,323	1,721	2,080	2,518	2,831	3,527	
22	0,532	0,858	1,321	1,717	2,074	2,508	2,819	3,505	
23	0,532	0,858	1,319	1,714	2,069	2,500	2,807	3,485	
24	0,531	0,857	1,318	1,711	2,064	2,492	2,797	3,467	
v	0,60	0,40	0,20	0,10	0,05	0,02	0,01	0,002	
α	ორმხრივი კრიტერიუმი								

სტიუდენტის განაწილება (გაგრძელება)

v \ α	ცალმხრივი კრიტერიუმი							
	0,30	0,20	0,10	0,05	0,025	0,01	0,005	0,001
25	0,531	0,856	1,316	1,708	2,060	2,485	2,787	3,450
26	0,531	0,856	1,315	1,706	2,056	2,479	2,779	3,435
27	0,531	0,855	1,314	1,703	2,052	2,473	2,771	3,421
28	0,530	0,855	1,313	1,701	2,048	2,467	2,763	3,408
29	0,530	0,854	1,311	1,699	2,045	2,462	2,756	3,398
30	0,530	0,854	1,310	1,697	2,042	2,457	2,750	3,385
40	0,529	0,851	1,303	1,684	2,021	2,423	2,704	3,307
50	0,528	0,849	1,298	1,676	2,009	2,403	2,678	3,262
60	0,527	0,848	1,296	1,671	2,000	2,390	2,660	3,232
80	0,527	0,846	1,292	1,664	1,990	2,374	2,639	3,195
100	0,526	0,845	1,290	1,660	1,984	2,365	2,626	3,174
200	0,525	0,843	1,286	1,653	1,972	2,345	2,601	3,131
500	0,525	0,842	1,283	1,648	1,965	2,334	2,586	3,106
∞	0,524	0,842	1,282	1,645	1,960	2,326	2,576	3,090
v \ α	0,60	0,40	0,20	0,10	0,05	0,02	0,01	0,002
	ორმხრივი კრიტერიუმი							

კრიტიკული მნიშვნელობები $Z_{\alpha,m}$
(ტომპსონის წესი)

$m \backslash \alpha$	0,01	0,05	0,10
1	1,414	1,410	1,397
2	1,715	1,645	1,559
3	1,917	1,757	1,611
4	2,051	1,814	1,631
5	2,142	1,848	1,640
6	2,207	1,870	1,644
7	2,256	1,885	1,647
8	2,294	1,896	1,648
9	2,324	1,904	1,649
10	2,348	1,910	1,649
11	2,368	1,915	1,649
12	2,385	1,920	1,650
13	2,399	1,923	1,650
14	2,411	1,926	1,650
15	2,422	1,929	1,649
16	2,431	1,931	1,649
17	2,440	1,933	1,649
18	2,447	1,934	1,649
19	2,454	1,936	1,649
20	2,460	1,937	1,649
21	2,465	1,938	1,649
22	2,470	1,939	1,649
23	2,475	1,940	1,649
24	2,479	1,941	1,649
25	2,483	1,942	1,648
26	2,486	1,943	1,648
27	2,490	1,943	1,648
28	2,493	1,944	1,648
29	2,496	1,945	1,648
30	2,498	1,945	1,648

$m \backslash \alpha$	0,01	0,05	0,10
31	2,501	1,946	1,648
32	2,503	1,946	1,648
33	2,505	1,947	1,648
34	2,507	1,947	1,648
35	2,509	1,947	1,648
36	2,511	1,948	1,648
37	2,513	1,948	1,648
38	2,514	1,948	1,648
39	2,516	1,949	1,647
40	2,518	1,949	1,647
41	2,519	1,949	1,647
42	2,520	1,950	1,647
43	2,522	1,950	1,647
44	2,523	1,950	1,647
45	2,524	1,950	1,647
46	2,525	1,951	1,647
47	2,526	1,951	1,647
48	2,527	1,951	1,647
49	2,528	1,951	1,647
50	2,529	1,951	1,647
55	2,533	1,952	1,647
60	2,537	1,953	1,647
65	2,540	1,953	1,646
70	2,542	1,954	1,646
75	2,545	1,954	1,646
80	2,547	1,955	1,646
85	2,548	1,955	1,646
90	2,550	1,955	1,646
95	2,551	1,955	1,646
100	2,553	1,956	1,646

χ^2 განაწილება

$\nu \backslash \sigma$	0,50	0,30	0,20	0,10	0,05	0,025	0,01	0,001
1	0,455	1,07	1,64	2,71	3,84	5,02	6,63	10,83
2	1,39	2,41	3,22	4,61	5,99	7,38	9,21	13,82
3	2,37	3,66	4,64	6,25	7,81	9,35	11,34	16,27
4	3,36	4,88	5,99	7,78	9,49	11,14	13,28	18,47
5	4,35	6,06	7,29	9,24	11,07	12,83	15,09	20,52
6	5,35	7,23	8,56	10,64	12,59	14,45	16,81	22,46
7	6,35	8,38	9,80	12,02	14,07	16,01	18,48	24,32
8	7,34	9,52	11,0	13,36	15,51	17,53	20,09	26,12
9	8,34	10,7	12,2	14,68	16,92	19,02	21,67	27,88
10	9,34	11,8	13,4	15,99	18,31	20,48	23,21	29,59
11	10,3	12,9	14,6	17,28	19,68	21,92	24,73	31,26
12	11,3	14,0	15,8	18,55	21,03	23,34	26,22	32,91
13	12,3	15,1	17,0	19,81	22,36	24,74	27,69	34,53
14	13,3	16,2	18,2	21,06	23,68	26,12	29,14	36,12
15	14,3	17,3	19,3	22,31	25,00	27,49	30,58	37,70
16	15,3	18,4	20,5	23,54	26,30	28,85	32,00	39,25
17	16,3	19,5	21,6	24,77	27,59	30,19	33,41	40,79
18	17,3	20,6	22,8	25,99	28,87	31,53	34,81	42,31
19	18,3	21,7	23,9	27,20	30,14	32,85	36,19	43,82
20	19,3	22,8	25,0	28,41	31,41	34,17	37,57	45,32
22	21,3	24,9	27,3	30,81	33,92	36,78	40,29	48,27
24	23,3	27,1	29,6	33,20	36,42	39,36	42,98	51,18
26	25,3	29,2	31,8	35,56	38,88	41,92	45,64	54,05
28	27,3	31,4	34,0	37,92	41,34	44,46	48,28	56,89
30	29,3	33,5	36,2	40,26	43,77	46,98	50,89	59,70
35	34,3	38,9	41,8	46,06	49,80	53,20	57,34	66,62
40	39,3	44,2	47,3	51,81	55,76	59,34	63,69	73,40
50	49,3	54,7	58,2	63,17	67,50	71,42	76,15	86,66
60	59,3	65,2	69,0	74,40	79,08	83,30	88,38	99,61
80	79,3	86,1	90,4	96,58	101,88	106,63	112,33	124,84
100	99,3	106,9	111,7	118,50	124,34	129,56	135,81	149,45
120	119,3	127,6	132,8	140,23	146,57	152,21	158,95	173,62
150	149,3	158,6	164,6	172,6	179,6	185,8	193,2	209,3
200	199,3	210,0	216,6	226,0	234,0	241,1	249,4	267,5

უილკოქსონი-მანნა-უიტნის U- კრიტერიუმი
(ცალმხრივი კრიტერიუმი $\alpha = 0,01$)

$n_1 \backslash n_2$	5	6	7	8	9	10	11	12	13	14	15	16	17	18	19	20
3	0	0	0	0	1	1	1	2	2	2	3	3	4	4	4	5
4	0	1	1	2	3	3	4	5	5	6	7	7	8	9	9	10
5	1	2	3	4	5	6	7	8	9	10	11	12	13	14	15	16
6		3	4	6	7	8	9	11	12	14	15	16	18	19	20	22
7			6	7	9	11	12	14	16	18	19	21	23	24	26	28
8				9	11	13	15	17	20	22	24	26	28	30	32	34
9					14	16	19	21	23	26	28	31	33	36	38	40
10						19	22	24	27	30	33	36	38	41	44	47
11							25	28	31	34	37	41	44	47	50	53
12								31	35	38	42	46	49	53	56	60
13									39	43	47	51	55	59	63	67
14										47	51	56	60	65	69	73
15											56	61	66	70	75	80
16												66	71	76	82	87
17													77	82	88	94
18														88	94	100
19															101	107
20																114

უილკოქსონი-მანნა-უიტნის U- კრიტერიუმი
(ორმხრივი კრიტერიუმი $\alpha = 0,01$)

$n_1 \backslash n_2$	5	6	7	8	9	10	11	12	13	14	15	16	17	18	19	20
5	0															
6	1	2														
7	1	3	4													
8	2	4	6	7												
9	3	5	7	9	11											
10	4	6	9	11	13	16										
11	5	7	10	13	16	19	21									
12	6	9	12	15	18	21	24	28								
13	7	10	13	17	20	24	27	31	34							
14	7	11	15	18	22	26	30	34	38	42						
15	8	12	16	20	25	29	33	37	42	46	51					
16	9	13	18	22	27	31	36	41	46	50	55	60				
17	10	15	19	24	29	34	39	44	49	54	60	65	70			
18	11	16	21	26	31	37	42	47	53	59	64	70	75	77	81	
19	12	17	22	28	34	39	45	51	57	63	69	75	81	87	93	
20	13	18	24	30	36	42	48	54	60	67	73	79	86	92	99	105
21	14	19	25	32	38	44	51	58	64	71	78	84	91	98	105	112
22	14	21	27	34	40	47	54	61	68	75	82	89	97	104	111	118
23	15	22	29	36	43	50	57	64	72	79	87	94	102	109	117	125
24	16	23	30	37	45	52	60	68	76	83	91	99	107	115	123	131
25	17	24	32	39	47	55	63	71	79	88	96	104	113	121	129	135

ფიქერის განაწილება (F განაწილება) $\alpha = 0,05$

$v_1 \backslash v_2$	1	2	3	4	5	6	8	12	20	24	∞
1	161	200	216	225	230	234	239	244	248	249	254
2	18,51	19,00	19,16	19,25	19,30	19,33	19,37	19,41	19,44	19,45	19,50
3	10,13	9,55	9,28	9,12	9,01	8,94	8,84	8,74	8,66	8,64	8,53
4	7,71	6,95	6,59	6,39	6,26	6,16	6,04	5,91	5,80	5,77	5,63
5	6,61	5,79	5,41	5,19	5,05	4,95	4,82	4,68	4,56	4,53	4,37
6	5,99	5,14	4,76	4,53	4,39	4,28	4,15	4,00	3,87	3,84	3,67
7	5,59	4,74	4,35	4,12	3,97	3,87	3,73	3,57	3,44	3,41	3,23
8	5,32	4,46	4,07	3,84	3,69	3,58	3,44	3,28	3,15	3,12	2,93
9	5,12	4,26	3,86	3,63	3,48	3,37	3,23	3,07	2,93	2,90	2,71
10	4,97	4,10	3,71	3,48	3,33	3,22	3,07	2,91	2,77	2,74	2,54
11	4,84	3,98	3,59	3,36	3,20	3,09	2,95	2,79	2,65	2,61	2,41
12	4,75	3,89	3,49	3,26	3,11	3,00	2,85	2,69	2,54	2,51	2,30
13	4,67	3,81	3,41	3,18	3,03	2,92	2,77	2,60	2,46	2,42	2,21
14	4,60	3,74	3,34	3,11	2,96	2,85	2,70	2,53	2,39	2,35	2,13
15	4,54	3,68	3,29	3,06	2,90	2,79	2,64	2,48	2,33	2,29	2,07
16	4,49	3,63	3,24	3,01	2,85	2,74	2,59	2,42	2,28	2,24	2,01
17	4,45	3,59	3,20	2,97	2,81	2,70	2,55	2,38	2,23	2,19	1,96
18	4,41	3,56	3,16	2,93	2,77	2,66	2,51	2,34	2,19	2,15	1,92
19	4,38	3,52	3,13	2,90	2,74	2,63	2,48	2,31	2,15	2,11	1,88
20	4,35	3,49	3,10	2,87	2,71	2,60	2,45	2,28	2,12	2,08	1,84
21	4,33	3,47	3,07	2,84	2,69	2,57	2,42	2,25	2,09	2,05	1,81
22	4,30	3,44	3,05	2,82	2,66	2,55	2,40	2,23	2,07	2,03	1,78
23	4,28	3,42	3,03	2,80	2,64	2,53	2,38	2,20	2,04	2,00	1,76
24	4,26	3,40	3,01	2,78	2,62	2,51	2,36	2,18	2,02	1,98	1,73
25	4,24	3,39	2,99	2,76	2,60	2,49	2,34	2,17	2,00	1,97	1,71
26	4,23	3,37	2,98	2,74	2,59	2,47	2,32	2,15	1,99	1,95	1,69
27	4,21	3,35	2,96	2,73	2,57	2,46	2,31	2,13	1,97	1,93	1,67
28	4,20	3,34	2,95	2,71	2,56	2,45	2,29	2,12	1,96	1,92	1,65
29	4,18	3,33	2,93	2,70	2,55	2,43	2,28	2,10	1,94	1,90	1,64
30	4,17	3,32	2,92	2,69	2,53	2,42	2,27	2,09	1,93	1,89	1,62
40	4,09	3,23	2,84	2,61	2,45	2,34	2,18	2,00	1,84	1,79	1,51
60	4,00	3,15	2,76	2,52	2,37	2,25	2,10	1,92	1,75	1,70	1,39
120	3,92	3,07	2,68	2,45	2,29	2,18	2,02	1,83	1,66	1,61	1,25
∞	3,84	3,00	2,61	2,37	2,21	2,10	1,94	1,75	1,57	1,52	1,00

ფიქურის განაწილება (F განაწილება) $\alpha = 0,01$

$v_1 \backslash v_2$	1	2	3	4	5	6	8	12	20	24	∞
1	4052	4999	5403	5625	5764	5859	5981	6106	6208	6234	6366
2	98,49	99,00	99,17	99,25	99,30	99,33	99,37	99,42	99,45	99,46	99,50
3	34,12	30,82	29,46	28,71	28,24	27,91	27,49	27,05	26,69	26,60	26,12
4	21,20	18,00	16,69	15,98	15,52	15,21	14,80	14,37	14,02	13,93	13,42
5	16,26	13,27	12,06	11,39	10,97	10,67	10,29	9,89	9,55	9,47	9,02
6	13,74	10,92	9,78	9,15	8,75	8,47	8,10	7,72	7,39	7,31	6,88
7	12,25	9,55	8,45	7,85	7,46	7,19	6,84	6,47	6,15	6,07	5,65
8	11,26	8,65	7,59	7,01	6,63	6,37	6,03	5,67	5,36	5,28	4,86
9	10,56	8,02	6,99	6,42	6,06	5,80	5,47	5,11	4,80	4,73	4,31
10	10,04	7,56	6,55	5,99	5,64	5,39	5,06	4,71	4,41	4,33	3,91
11	9,65	7,20	6,22	5,67	5,32	5,07	4,74	4,40	4,10	4,02	3,60
12	9,33	6,93	5,95	5,41	5,06	4,82	4,50	4,16	3,86	3,78	3,36
13	9,07	6,70	5,74	5,20	4,86	4,62	4,30	3,96	3,67	3,59	3,16
14	8,86	6,51	5,56	5,03	4,69	4,46	4,14	3,80	3,51	3,43	3,00
15	8,68	6,36	5,42	4,89	4,56	4,32	4,00	3,67	3,36	3,29	2,87
16	8,53	6,23	5,29	4,77	4,44	4,20	3,89	3,55	3,25	3,18	2,75
17	8,40	6,11	5,18	4,67	4,34	4,10	3,79	3,45	3,16	3,08	2,65
18	8,28	6,01	5,09	4,58	4,25	4,01	3,71	3,37	3,07	3,00	2,57
19	8,18	5,93	5,01	4,50	4,17	3,94	3,63	3,30	3,00	2,92	2,49
20	8,10	5,85	4,94	4,43	4,10	3,87	3,56	3,23	2,94	2,86	2,42
21	8,02	5,78	4,87	4,37	4,04	3,81	3,51	3,17	2,88	2,80	2,36
22	7,94	5,72	4,82	4,31	3,99	3,76	3,45	3,12	2,83	2,75	2,31
23	7,88	5,66	4,76	4,26	3,94	3,71	3,41	3,07	2,78	2,70	2,26
24	7,82	5,61	4,72	4,22	3,90	3,67	3,36	3,03	2,74	2,66	2,21
25	7,77	5,57	4,68	4,18	3,86	3,63	3,32	2,99	2,70	2,62	2,17
26	7,72	5,53	4,64	4,14	3,82	3,59	3,29	2,96	2,66	2,58	2,13
27	7,68	5,49	4,60	4,11	3,79	3,56	3,26	2,93	2,63	2,55	2,10
28	7,64	5,45	4,57	4,07	3,76	3,53	3,23	2,90	2,60	2,52	2,06
29	7,60	5,42	4,54	4,04	3,73	3,50	3,20	2,87	2,57	2,49	2,03
30	7,56	5,39	4,51	4,02	3,70	3,47	3,17	2,84	2,55	2,47	2,01
40	7,31	5,18	4,31	3,83	3,51	3,29	2,99	2,66	2,37	2,29	1,61
60	7,08	4,98	4,13	3,65	3,34	3,12	2,82	2,36	2,20	2,12	1,60
120	6,85	4,79	3,95	3,48	3,17	2,96	2,66	2,34	2,03	1,95	1,38
∞	6,63	4,60	3,78	3,32	3,02	2,80	2,51	2,18	1,87	1,79	1,00

უილკოქსონის T-კრიტერიუმი

α v	ორმხრივი კრიტერიუმი		ცალმხრივი კრიტერიუმი		α v	ორმხრივი კრიტერიუმი		ცალმხრივი კრიტერიუმი	
	0,05	0,01	0,05	0,01		0,05	0,01	0,05	0,01
6	0		2		36	208	171	227	185
7	2		3	0	37	221	182	241	198
8	3	0	5	1	38	235	194	256	211
9	5	1	8	3	39	249	207	271	224
10	8	3	10	5	40	264	220	286	238
11	10	5	13	7	41	279	233	302	252
12	13	7	17	9	42	294	247	319	266
13	17	9	21	12	43	310	261	336	281
14	21	12	25	15	44	327	276	353	296
15	25	15	30	19	45	343	291	371	312
16	29	19	35	23	46	361	307	389	328
17	34	23	41	27	47	378	322	407	345
18	40	27	47	32	48	396	339	426	362
19	46	32	53	37	49	415	355	446	379
20	52	37	60	43	50	434	373	466	397
21	58	42	67	49	51	453	390	486	416
22	65	48	75	55	52	473	408	507	434
23	73	54	83	62	53	494	427	529	454
24	81	61	91	69	54	514	445	550	473
25	89	68	100	76	55	536	465	573	493
26	98	75	110	84	56	557	484	595	514
27	107	83	119	92	57	579	504	618	535
28	116	91	130	101	58	602	525	642	556
29	126	100	140	110	59	625	546	666	578
30	137	109	151	120	60	648	567	690	600
31	147	118	163	130	61	672	589	715	623
32	159	128	175	140	62	697	611	741	646
33	170	138	187	151	63	721	634	767	669
34	182	148	200	162	64	747	657	793	693
35	195	159	213	173	65	772	681	820	718

σ – ის ნდობის ინტერვალები
ნდობის ინტერვალის γ_1 – ქვედა და γ_2 – ზედა საზღვრები

$$\gamma_1 \sigma < \sigma < \gamma_2 \sigma \left(\sigma = \sqrt{\frac{1}{n-1} \sum_{i=1}^n (x_i - \bar{x})^2} \right)$$

α $v = n - 1$	0,01		0,02		0,05		0,10	
	γ_1	γ_2	γ_1	γ_2	γ_1	γ_2	γ_1	γ_2
1	0,36	159	0,39	79,8	0,45	31,9	0,51	15,9
2	0,43	14,1	0,47	9,97	0,52	6,28	0,58	4,40
3	0,48	6,47	0,51	5,11	0,57	3,73	0,62	2,92
4	0,52	4,39	0,55	3,67	0,60	2,87	0,65	2,37
5	0,55	3,48	0,58	3,00	0,62	2,45	0,67	2,09
6	0,57	2,98	0,60	2,62	0,64	2,20	0,68	1,92
7	0,59	2,66	0,62	2,38	0,66	2,04	0,71	1,80
8	0,60	2,44	0,63	2,21	0,68	1,92	0,72	1,71
9	0,62	2,28	0,64	2,08	0,69	1,83	0,73	1,65
10	0,63	2,15	0,66	1,98	0,70	1,76	0,74	1,59
11	0,64	2,06	0,67	1,90	0,71	1,90	0,75	1,55
12	0,65	1,98	0,68	1,83	0,72	1,65	0,76	1,52
13	0,66	1,91	0,69	1,78	0,73	1,61	0,76	1,49
14	0,67	1,85	0,69	1,73	0,73	1,58	0,77	1,46
15	0,68	1,81	0,70	1,69	0,74	1,55	0,78	1,44
16	0,68	1,78	0,71	1,66	0,75	1,52	0,78	1,42
17	0,69	1,73	0,71	1,63	0,75	1,50	0,79	1,40
18	0,70	1,70	0,72	1,60	0,76	1,48	0,79	1,39
19	0,70	1,67	0,73	1,58	0,76	1,46	0,79	1,37
20	0,71	1,64	0,73	1,56	0,77	1,44	0,80	1,36
21	0,71	1,62	0,73	1,54	0,77	1,43	0,80	1,35
22	0,72	1,60	0,74	1,52	0,77	1,42	0,81	1,34
23	0,72	1,58	0,74	1,50	0,78	1,40	0,81	1,33
24	0,73	1,56	0,75	1,49	0,78	1,39	0,81	1,32
25	0,73	1,54	0,75	1,47	0,78	1,38	0,82	1,31
26	0,73	1,53	0,76	1,46	0,79	1,37	0,82	1,30
27	0,74	1,51	0,76	1,45	0,79	1,36	0,82	1,29
28	0,74	1,50	0,76	1,44	0,79	1,35	0,83	1,29
29	0,74	1,49	0,77	1,43	0,80	1,34	0,83	1,28
30	0,75	1,48	0,77	1,42	0,80	1,34	0,83	1,27
40	0,77	1,39	0,79	1,34	0,82	1,28	0,85	1,23
50	0,79	1,34	0,81	1,30	0,84	1,24	0,86	1,20
60	0,81	1,30	0,82	1,27	0,85	1,22	0,87	1,18
70	0,82	1,27	0,84	1,24	0,86	1,20	0,88	1,16
80	0,83	1,25	0,84	1,22	0,87	1,18	0,89	1,15
90	0,84	1,23	0,85	1,21	0,87	1,17	0,89	1,14
100	0,85	1,22	0,86	1,20	0,88	1,16	0,90	1,13
200	0,89	1,15	0,90	1,13	0,91	1,11	0,93	1,09

q კრიტიკული მნიშვნელობები ($\alpha' = 0,05$)

$l \backslash v$	2	3	4	5	6	7	8	9	10
1	17,97	26,98	32,82	37,08	40,41	43,40	45,40	47,36	49,07
2	6,09	8,33	9,80	10,88	11,74	12,44	13,03	13,54	13,99
3	4,50	5,91	6,83	7,50	8,04	8,48	8,85	9,18	9,46
4	3,93	5,04	5,76	6,29	6,71	7,05	7,35	7,60	7,83
5	3,64	4,60	5,22	5,67	6,03	6,33	6,58	6,80	7,00
6	3,46	4,34	4,90	5,31	5,63	5,90	6,12	6,32	6,49
7	3,34	4,17	4,68	5,06	5,36	5,61	5,82	6,00	6,16
8	3,26	4,04	4,53	4,89	5,17	5,40	5,60	5,77	5,92
9	3,20	3,95	4,42	4,76	5,02	5,24	5,43	5,60	5,74
10	3,15	3,88	3,33	4,65	4,91	5,12	5,31	5,46	5,60
11	3,11	3,82	4,26	4,57	4,82	5,03	5,20	5,35	5,49
12	3,08	3,77	4,20	4,51	4,75	4,95	5,12	5,27	5,40
13	3,06	3,74	4,15	4,46	4,69	4,89	5,05	5,19	5,32
14	3,03	3,70	4,11	4,41	4,64	4,83	5,00	5,13	5,25
15	3,01	3,67	4,08	4,37	4,60	4,78	4,94	5,08	5,20
16	3,00	3,65	4,05	4,33	4,56	4,74	4,90	5,03	5,15
17	2,98	3,63	4,02	4,30	4,52	4,71	4,86	4,99	5,11
18	2,97	3,61	4,00	4,28	4,50	4,67	4,82	4,96	5,07
19	2,96	3,59	3,98	4,25	4,47	4,65	4,79	4,92	5,04
20	2,95	3,58	3,96	4,23	4,45	4,62	4,77	4,90	5,01
24	2,92	3,53	3,90	4,17	4,37	4,54	4,68	4,81	4,92
30	2,89	3,49	3,85	4,10	4,30	4,46	4,60	4,72	4,82
40	2,86	3,44	3,79	4,04	4,23	4,39	4,52	4,64	4,74
60	2,83	3,40	3,74	3,98	4,16	4,31	4,44	4,55	4,65
120	2,80	3,36	3,69	3,92	4,10	4,24	4,36	4,47	4,57
∞	2,77	3,31	3,63	3,86	4,03	4,17	4,29	4,39	4,47

q კრიტიკული მნიშვნელობები ($\alpha' = 0,01$)

$v \backslash l$	2	3	4	5	6	7	8	9	10
1	90,03	135,0	164,3	185,6	202,2	215,8	227,2	237,0	245,5
2	14,04	90,02	22,29	24,72	26,63	28,20	29,53	30,68	31,69
3	8,26	10,62	12,17	13,33	14,24	15,00	15,64	16,20	16,69
4	6,51	8,12	9,17	9,96	10,58	11,10	11,55	11,93	12,27
5	5,70	6,98	7,80	8,42	8,91	9,32	9,70	9,97	10,24
6	5,24	6,33	7,03	7,56	7,97	8,32	8,61	8,87	9,10
7	4,95	5,92	6,54	7,01	7,37	7,68	7,94	8,17	8,37
8	4,75	5,64	6,20	6,63	6,96	7,24	7,47	7,68	7,86
9	4,60	5,43	5,96	6,35	6,66	6,92	7,13	7,33	7,50
10	4,48	5,27	5,77	6,14	6,43	6,67	6,88	7,06	7,21
11	4,39	5,15	5,62	5,97	6,25	6,48	6,67	6,84	6,99
12	4,32	5,05	5,50	5,84	6,10	6,32	6,51	6,67	6,81
13	4,26	4,96	5,40	5,73	5,98	6,19	6,37	6,53	6,67
14	4,21	4,90	5,32	5,63	5,88	6,09	6,26	6,41	6,54
15	4,17	4,84	5,25	5,56	5,80	5,99	6,16	6,31	6,44
16	4,13	4,79	5,19	5,49	5,72	5,92	6,08	6,22	6,35
17	4,10	4,74	5,14	5,43	5,66	5,85	6,01	6,15	6,27
18	4,07	4,70	5,09	5,38	5,60	5,79	5,94	6,08	6,20
19	4,05	4,67	5,05	5,33	5,55	5,74	5,89	6,02	6,14
20	4,02	4,64	5,02	5,29	5,51	5,69	5,84	5,97	6,09
24	3,96	4,55	4,91	5,17	5,37	5,54	5,69	5,81	5,92
30	3,89	4,46	4,80	5,05	5,24	5,40	5,54	5,65	5,76
40	3,83	4,37	4,70	4,93	5,11	5,27	5,39	5,50	5,56
60	3,76	4,28	4,60	4,82	4,99	5,13	5,25	5,36	5,45
120	3,70	4,20	4,50	4,71	4,87	5,01	5,12	5,21	5,30
∞	3,64	4,12	4,40	4,60	4,76	4,88	4,99	5,08	5,16

ლიტერატურა

1. Лакин Г.Ф. Биометрия, М., 1984.
2. Гланц С. Медико-биологическая статистика, М., Практика, 1999.
3. Реброва О.Ю. Статистический анализ медицинских данных. М., Медио Сфера, 2002.
4. Гелман В.Я. Медицинская информатика. СПб. 2001.
5. Компьютерная биометрика. Под ред. В.Н. Носова, МГУ, 1990.
6. Максимов Г.К., Синицин А.Н. Статистическое моделирование многомерных систем в медицине. М., “Медицина”, 1983.
7. Афифи А., Эйзен С. Статистический анализ. Подход с использованием ЭВМ. М., “Статистика” . 1982.

სარჩევი

შესავალი

ალბათობის თეორიის საფუძვლები

1. ხდომილობა და მისი ალბათობა-----	5
1.1 ხდომილობის კლასიფიკაცია-----	5
1.2 მოქმედებები ხდომილობებზე -----	6
1.3 ხდომილობის ალბათობა -----	7
1.4 კომბინატორიკის ძირითადი ფორმულები -----	10
1.5 ალბათობის თეორიის ძირითადი თეორემები -----	13
1.6 სრული ალბათობის ფორმულა. ბაიესის ფორმულა ----	17
1.7 ცდების გამეორება. ბერნულის ფორმულა -----	20
2. შამთხვევითი სიდიდეები -----	22
2.1 შემთხვევითი სიდიდის ცნება და მისი განაწილების კანონი -----	22
2.2 შემთხვევით სიდიდეთა სისტემა -----	27
2.3 პირობითი განაწილების სიმკვრივის ფუნქცია -----	30
2.4 შემთხვევითი სიდიდის რიცხვითი მახასიათებლები ---	33
2.5 ორგანზომილებიანი შემთხვევითი სისტემის რიცხვითი მახასიათებლები -----	38
3. შამთხვევითი სიდიდეების ძირითადი განაწილების კანონები -----	40
3.1 ბინომური განაწილების კანონი -----	40
3.2 პუასონის განაწილება -----	41
3.3 თანაბარი განაწილება -----	42
3.4 მაჩვენებლიანი განაწილება -----	44
3.5 ნორმალური განაწილება -----	45
3.6 ნორმალური განაწილება სიბრტყეზე -----	49
3.7 მრავალგანზომილებიანი სისტემის ნორმალური განაწილების კანონი -----	51
4. ალბათობის თეორიის ზღვარითი თეორემები-----	52
4.1 დიდ რიცხვთა კანონი -----	53
4.2 ცენტრალური ზღვარითი თეორემები -----	55
ბიოსტატისტიკის მეთოდები	
5. ბიოსტატისტიკის არსი -----	57
6. მონაცემების კლასიფიკაცია და მათი ნარმოდენის მეთოდები. სიხშირული ანალიზი -----	59

7.	მათემატიკურ სტატისტიკაში გამოყენებული ძირითადი განაწილების კანონები -----	68
8.	ძირითადი სტატისტიკური მახასიათებლები -----	72
8.1	განაწილები მდებარეობის მახასიათებლები -----	72
8.2	განაწილების ცვალებადობის მახასიათებლები -----	78
8.3	განაწილების ფორმის მახასიათებლები -----	81
8.4	თვისებრივი მაჩვენებლების სტატისტიკურ მახასიათებლები -----	86
9.	უცნობი პარამეტრების სტატისტიკური შეფასება	
9.1	პარამეტრების შეფასების ცნება -----	87
9.2	პარამეტრების შეფასების წერტილობანი მეთოდები --	90
9.3	პარამეტრების შეფასების ინტერვალური მეთოდები --	96
10.	ჰიპოთეზების სტატისტიკური შემოწმების პარამეტრული მეთოდები -----	103
10.1	სტატისტიკური ჰიპოთეზის ცნება -----	103
10.2	შემთხვევითი სიდიდის საშუალოსა და დისპერსიაზე სტატისტიკური ჰიპოთეზის შემოწმება -----	109
10.3	დისპერსიების ტოლობის ჰიპოთეზის შემოწმება. ფიშერის კრიტერიუმი -----	111
10.4	საშუალოების ტოლობის ჰიპოთეზის შემოწმება. სტიუდენტის კრიტერიუმი -----	112
10.5	საშუალოების მრავლობითი შედარება -----	115
10.6	ორი დამოკიდებული ამონარჩევის შედარება -----	118
10.7	ამონარჩევში უხეში შეცდომების გამოვლენის მეთოდები -----	120
10.8	ორზე მეტი ამონარჩევის ერთდროული შედარება --	124
10.9	მრავალგანზომილებიანი სისტემის ზოგიერთი ჰიპოთეზების შემოწმება -----	129
11.	ჰიპოთეზების სტატისტიკური შემოწმების არაპარამეტრული მეთოდები -----	131
11.1	განაწილების კანონის შესახებ ჰიპოთეზის -----	131
11.2	ორი ამონარჩევის შედარების ჰიპოთეზა -----	136
11.3	ამონარჩევების მრავლობითი შედარება -----	140
11.4	ორი დამოკიდებული ამონარჩევის შედარება -----	142
11.5	ფარდობითი სიხშირეების ტოლობის ჰიპოთეზის შემოწმება -----	144
12.	კორელაციური ანალიზი -----	146
12.1	კორელაციური ანალიზის არსი -----	146
12.2	კორელაციის კოეფიციენტის განსაზღვრის პარამეტრული მეთოდები -----	148

12.3	კორელაციის კოეფიციენტის განსაზღვრის არაპარამეტრული მეთოდები -----	155
12.4	კონკორდაციის კოეფიციენტი -----	160
13.	რებრესიული ანალიზის საფუძვლები -----	163
13.1	რეგრესიული ანალიზის არსი -----	163
13.2	უმცირეს კვადრატთა მეთოდი -----	167
13.3	წრფივი რეგრესია -----	168
13.4	წრფივი რეგრესიის განტოლების ნდობის ინტერვალი -----	172
13.5	არაწრფივი რეგრესია -----	174
13.6	ნაშთთა ანალიზი -----	180
14.	მრავლობითი რებრესიული ანალიზი -----	181
14.1	მრავლობითი რეგრესიის მოდელი -----	181
14.2	მრავლობითი წრფივი რეგრესიის განტოლების ნდობის ინტერვალი -----	186
14.3	ცვლადების შერჩევა -----	188
14.4	ნარჩენების ავტოკორელაციისა და მულტიკოლინეარობის პრობლემები -----	190
15.	დისპერსიული ანალიზის საფუძვლები -----	193
15.1	მეთოდის არსი -----	193
15.2	ერთფაქტორიანი დისპერსიული ანალიზი -----	195
15.3	ორფაქტორიანი დისპერსიული ანალიზი -----	201
15.4	დისპერსიული ანალიზის არაპარამეტრული მეთოდები -----	209
16.	მთავარი კომპონენტების მეთოდი -----	211
17.	ფაქტორული ანალიზი -----	221
	დანართი -----	230
	ლიტერატურა -----	246