2019

# THEORY PROBABILITY AND MATHEMATICAL STATISTICS

## textbook for students of economic specialties

Plaskon S., Eremenko V., Martyniuk O., Berezka K., Nemish V., Ruska R., Popina S., Seniv G., Homa-Mohylskaya S., Shinkarik M.

TERNOPIL

Plaskon S., Eremenko V., Martyniuk O., Berezka K., Nemish V., Ruska R., Popina S., Seniv G., Homa-Mohylskaya S., Shinkarik M.
**THEORY PROBABILITY AND MATHEMATICAL STATISTICS /** textbook for students of economic specialties/ - Тернопіль, ТНЕУ. – 2019. – 90 p

*Затверджено на засіданні кафедри прикладної математики протокол №2 від 15.10.2019р.*
*Рекомендовано вченою радою ФКІТ протокол №2 від 31.10.2019р.*

2

PART ONE
Probability theory

## § 1. Probability

1. During testing we mean the realization of the intended action and getting results in the performance of certain conditions set S. It is assumed that these conditions are fixed; they either objectively exist or are created artificially and can be played an unlimited number of times. The result of the test is an event. There are significant events, impossible and random.

Credible called an event that must occur during testing. Impossible to call an event that required when testing occurs. Random - this is the event that at the straighten-accommodated as can happen, did not happen. Credible event denoted by the letter $\Omega$, and impossible - $\varnothing$. Let us consider some properties of random events. Two events are called incompatible (compatible) if ye are dwelling serving one excludes (does not exclude) the departure of another. The set of random events form a complete group, if one of them when testing must occur, and any two events are mutually exclusive. Elementary will be called simple random events that can occur during testing.
Events "generated" one test, called elementary events with the same probability if there is reason to believe that none of them is more feasible than                                                                                     others.

2. Let consider the next notion.
An **experiment** is any process of observation that has an uncertain outcome. The process must be defined so that on any single repetition of the experi-

ment, one and only one of the possible outcomes will occur. The possible outcomes for an experiment are called **experimental outcomes**.

1. The probability assigned to each experimental outcome must be between 0 and 1. That is, If *A* represents an experimental outcome and if *P(A)* represents the probability of this outcome, then $0 \leq P(A) \leq 1$.
2. The probability of all of the experimental outcomes must sum to 1.

The **sample space** of an experiment is the set of all possible experimental outcomes. The experimental outcomes in the sample space are often called **sample space outcomes**.

An **event** is a set (or collection) of sample space outcomes.

The **probability of an event** is the **sum of the probabilities of the sample space outcomes** that correspond to the event.

**Definition.** Classical probability of event A is the ratio of the number of elementary events with the same probability contributing to the event A, and the total number of elementary events with the same probability that form a complete group. Each elementary random events, in fact, is one of outcomes tests. This approach is useful in analyzing problems. Given this observation analytical expression of the classical definitions come like this:

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

where *P (A)* - classical probability of event A;
*m* - number of elementary events with the same probability contributing to the event A (number of consequences test in which an event occurs A);
*n* - number of elementary events with the same probability that form a full group (the total number of effects (consequences) with the same probability of the tests).

This definition allows us to formulate the basic properties of the classical probabilities:
1) P (Ω) = 1 (probability of accurate event is 1);
2) P (∅) = 0 (the probability of the impossible event is 0);
3) if A - random event, then:

$$0 < P(A) < 1. \qquad (1.2)$$

3. In order to have some standard methods in the calculation scheme for classical probability, we present the basic formula of combinatorics, and consider the notion of combinations, arrangements and permutations.
Suppose there are $k$ groups of elements, the number of each of which is respectively $n_1, n_2, ..., n_k$.

Then the total number $N$ of ways that you can make a group, is given by

$$N = n_1 \cdot n_2 \cdot ... \cdot n_k, \tag{1.3}$$

called basic formula of combinatorics. Consider a set of different $n$ elements of arbitrary nature. We will form groups of $m$ $(m \leq n)$ different elements of this set.
Such groups in probability theory often referred samples.
Let $m < n$. Combinations are those groups that differ from each other by at least one element. The total number of combinations (read: it's from $n$ to $m$) is given by

$$C_n^m = \frac{n(n-1)\cdot(n-2)\cdot...\cdot(n-m+1)}{m!}, \tag{1.4}$$

де $m! = m(m-1) \cdot ... \cdot 2 \cdot 1$ (read: m factorial).

**Remarks**. In the numerator of (1.4) is the m factors. If the numerator and denominator multiplied by (n - m) !, then the next equation will be obtained:

$$C_n^m = \frac{n!}{(n-m)!m!}. \tag{1.4*}$$

Let $m \leq n$. Placement are those groups that are different from one another, or at least one element, or the order of these elements in the group. The number of placement (read: a from $n$ on $m$) is given by:

$$A_n^m = \underbrace{n(n-1)\cdot(n-2)\cdot...\cdot(n-m+1)}_{m \text{ співмножників}}. \tag{1.5}$$

If in the formula (1.5) $m = n$, then $A_n^m$ - the number of placements that differ only in the order of the elements, not the elements themselves. Such permutations are called accommodation. Their number $P_n$ by the formula (1.5)

$$A_n^n = n(n-1)(n-2)...2 \cdot 1 = n! = P_n,$$

that is,

$$P_n = A_n^n = n!. \tag{1.6}$$

The number n can take values not only natural, it can also be zero.

Empty set (the sample) is a subset of any set and naturally assume that it might be updated only one way. Therefore, it is considered that $0! = 1$.

The number of combinations has the following properties:

1) $C_n^0 = C_n^n = 1$;  2) $C_n^1 = C_n^{n-1} = n$;  3) $C_n^m = C_n^{n-m}$;

4) $C_n^0 + C_n^1 + C_n^2 + \ldots + C_n^n = 2^n$.

Note that the number of placements, permutations and combinations associated equality

$A_n^m = P_m C_n^m$.

When solving problems of combinatorics, the following rules apply:

**Rule amount**. If an object $\alpha$ can be selected from the set of objects $k$ ways, and the second object $\beta$ can be selected $s$ ways, then selection either $\alpha$, either $\beta$ can $k + s$ ways.

**Product rule**. If an item $\alpha$ can be selected from the set of objects $k$ means and after each such selection object $\beta$ can be selected $s$ ways, then pair of objects $(\alpha, \beta)$ in that order can be selected $k \cdot s$ ways.

4. Together with the basic concepts of probability to probability theory belongs to relative frequency.

Relative frequency of random events is the ratio of the number of tests in which the event occurred, the total number of tests actually performed. That is, the relative frequency of the event A is defined by the formula:

$$W(A) = \frac{M}{N}, \qquad\qquad (1.7)$$

where $M$ - number of occurrences of event $A$, $N$- total number of tests.

Comparison of definitions of probability and relative frequency suggests: the probability calculated before the test ( it is a priori value), and the relative frequency - after testing (posterior value).

From the definition (1.7) for the random event $A$ follows this double inequality (compare with (1.2)!):

$$0 \le W(A) \le 1.$$

In a small number of tests relative frequency of random events may be considerably changed from one test series to another. However, long-term observations have shown that when the same conditions is made sufficiently large number of tests, it shows the relative frequency stability property. This property lies in the fact that different series of tests relative frequency changes little (less, the greater the number of tests), fluctuating around some constant number. It turned out that this number is a constant probability of a random event. The mathematical formulation of the stability properties will be given in the topic "The law of large numbers" (J. Bernoulli's theorem).

Property relative frequency stability, and the ability, at least in principle, to hold an unlimited number of tests with respect to random events allow to formulate a statistical definition of probability, as statistical probability of random event of a relative frequency of the event.

**5.** One of the drawbacks of the classical definition of probability associated with the inability to use this definition for the case of an infinite number of trials testing the consequences can be eliminated by using the geometric probability - the probability of getting a point in the region (the segment of the plane, body parts, etc.). Assume that the test is carried out - a rectangle $\Omega$, which contains an arbitrary figure A, randomly throws point (Fig. 1.1).



Pic.1.1.

It is assumed that the following assumption: it can be in any part of the rectangle $\Omega$, the probability (possibility) of hit points on a figure proportional to the area $A$ of the figure and does not depend on the position of $A$ relative $\Omega$, or of the form $A$. Let random event $A$ - the point being deceived was in the shape $A$. Then the **geometric definition of the probability** of event $A$ given equality:

$$P(A) = \frac{S(A)}{S(\Omega)} \quad , \tag{1.8}$$

where S ($A$) - square shape $A$, S ($\Omega$) - the area of the rectangle $\Omega$.
Definition (1.8) is a particular case of the general definition of geometric probability. For $\Omega \subset \square^1$ *або* $\Omega \subset \square^3$.

$$P(A) = \frac{l(A)}{l(\Omega)} \quad a\text{бо } P(A) = \frac{V(A)}{V(\Omega)}. \tag{1.9}$$

# SOLVING COMMON PROBLEMS

Task 1.1. Bank may issue a month on credit loan five to its customers, while orders received on loan from 15 customers first region and 10 clients   second region. To maintain customer bank regards as a temporary forced action - playing randomly loans among five of whom acted order. Find the probability that the number of clients first district, which   poses   gets   on,   is:   a)   5;   b)   0;   c)   3.
○   Test - playing among clients. Outcome trials - five customers. The number of consequences equal to the number of groups with five  clients that can create from a set number of 25 items (bank customers). For these groups of elements characteristic is that they differ from each other by at least one element. In addition, as a basic set of elements (from 25 customers) and formed groups of 5 elements consist of different elements (no repeats). This leads to the conclusion that such groups are combinations. And their numbers for all three cases is:

$$n = C_{25}^5 = \frac{25 \cdot 24 \cdot 23 \cdot 22 \cdot 21}{5 \cdot 4 \cdot 3 \cdot 2 \cdot 1} = 53130.$$

a) $A$ - owners of the loan is five customers first district. The number of consequences tests for which there is an event $A$ is equal to the number groups of 5 elements that can create from a set number of 15 customers first region. Their total number

$$m = C_{15}^5 = \frac{15 \cdot 14 \cdot 13 \cdot 12 \cdot 11}{5 \cdot 4 \cdot 3 \cdot 2 \cdot 1} = 3003.$$

According to the classical definition of probability

$$P(A) = \frac{m}{n} = \frac{3003}{53130} \approx 0,057.$$

b) $B$ - owners of the loan is five clients of the second district. The number of consequences trial $m$, which occurs in the event, equal to the number groups of 5 elements that can create from a set number of 10 customers second region. Their total number

$$m = C_{10}^5 = \frac{10 \cdot 9 \cdot 8 \cdot 7 \cdot 6}{5 \cdot 4 \cdot 3 \cdot 2 \cdot 1} = 252,$$

and          $$P(B) = \frac{m}{n} = \frac{252}{53130} \approx 0,005.$$

c) C - owners of the loan is three clients first region and a second clients

6

second region. Then total number groups of three clients first region and a second clients second region using Product rule equals

$$m = C_{15}^3 \cdot C_{10}^2 = \frac{15 \cdot 14 \cdot 13}{3 \cdot 2 \cdot 1} \cdot \frac{10 \cdot 9}{2 \cdot 1} = 455 \cdot 45 = 20475,$$

where

$$P(C) = \frac{20475}{53130} \approx 0,385. \quad \bullet$$

$$P(A) = P(25 \leq x \leq 30) = \frac{30 - 25}{40 - 20} = \frac{5}{20} = 0,25. \quad \bullet$$

# § 2. Theorem multiplication and addition probability and their consequences

1. Performance of events (event algebra). Diagrams of Vienna.
2. The conditional probability. Multiplication theorem of probability.
3. Addition theorem of probability.
4. The main property of events that form a complete group. The probability of at least one event. The probability of serving only one event.
5. The algorithm for solving problems using addition and multiplication theorems of probability.
6. The formula of total probability. Bayes' Formula.
7. Algorithm for solving problems using formulas of total probability and Bayes Formula.

1. Enter into consideration actions on events. **Two events that form a complete group called opposite. Definition.** In the event $\overline{A}$ (read: not A) refers to an event that is the fact that the test event $A$ happens and an event occurs, the opposite of the event $A$. The action of an event is determined trait of this event, called the **negation** of (the event). **The sum of events $A_1, A_2, ..., A_k$ is called event $A$, which is that when testing occurs at least one of the events $A_1, A_2, ..., A_k$.** Symbolic Entry: $A = A_1 + A_2 + ... + A_k$. In particular, if $k = 2$, then the event $A = A_1 + A_2$ is the fact that the trial takes place or event $A_1$, or event $A_2$, or $A_1$ and $A_2$ (where $A_1$ and $A_2$ are compatible). In this connection, we obtain the following mnemonic rule: summation associated with the word "or» («+» ↔ «or»). If $A_1$ and $A_2$ are incompatible, then event $A_1 + A_2$ means occurring event $A_1$ or $A_2$ after testing. **The product of the events $A_1, A_2, ..., A_k$ is called event A, which is**

7

the fact that the trial take place all the events $A_1$, $A_2$, ..., $A_k$. Symbolically: $A = A_1 \cdot A_2 \cdot ... \cdot A_k$ (or $A = A_1 A_2 ... A_k$). Departures all $k$ events involves serving and event $A_1$, and event $A_2$, ..., and event $A_k$. Thus, the effect of multiplying the mnemonics associated with «and» («×» ↔ «and»).

**2.** Let the trial, in which can occur B, supplemented by the condition of the departure of a random event A. Then the **probability of event B, found under the condition that the event A has occurred, is called the conditional probability** of event $B$ and is denoted by $P_A(B)$ or $P(B/A)$. **Example**. There are 6 white and 4 black balls in the box, identical in size and feel. Two balls are randomly selected without returning. Find the probability that the second randomly selected ball is black, if the first ball appeared white.

We introduce into consideration events: $A$ - the first ball is white, B - second ball is black. Then accordingly the classical definition

$$P_A(B) = \frac{m}{n} = \frac{4}{9}.$$

In calculating the conditional probabilities are encouraged to include the contents of random events, beginning with the reading of conditional probability. For example, searching for the right part of the recorded above formula, the left should read: what is the probability that the second randomly selected ball would be black if the first ball appeared white. That is, when using the classical definition of probability, be sure to consider the departure event A (in this example $n = 9$, for the first ball selected, $m = 4$, since the selection did not change the number of balls black).

Two events are called **independent** if the probability of one of them does not change that took place or not is different. Analytical criterion of independence of events A and B are equality

$$P_A(B) = P_{\bar{A}}(B). \tag{2.1}$$

In the case of (2.1) we can assume that

$$P_A(B) = P_{\bar{A}}(B) = P(B), \tag{2.2}$$

where the last probability is called unconditional. Two events are called **dependent** if the probability of one of them is affected by the departure of another. That is, if events A and B are dependent, then equation (2.2) is violated:

$$P_A(B) \neq P_{\bar{A}}(B). \tag{2.3}$$

**Multiplication theorem of probability**

The probability of the product of two compatible events is the product of the

probability of one of them on the conditional probability of an event, calculated under the assumption that the first event occurred:

$$P(AB) = P(A)P_A(B)\big(= P(B)P_B(A)\big). \tag{2.4}$$

*Remarks*. Continued equality in parentheses indicates " equality" events $A$ and $B$ with regard to actions commutativity of multiplication ($AB = BA$).

**Corollary 1**. The probability of the product of $k$ events is the product of the probability of one of them on the conditional probabilities of all the rest, and the probability of each next event is the assumption that all previous events have taken place:

$$P(A_1 A_2 A_3 .. A_k) = P(A_1)P_{A_1}(A_2)P_{A_1 A_2}(A_3)...P_{A_1 A_2 .. A_{k-1}}(A_k). \tag{2.5}$$

**Corollary 2**. **If events A and B are independent**, then

$$P(AB) = P(A)\ P(B). \tag{2.6}$$

To summarize of Corollary 2 to the case of an arbitrary number of events, consider the following definition.

Several events are called pairwise independent if every two of them are independent. For example, the events $A_1$, $A_2$, $A_3$ are pairwise independent, if independent events $A_1$ and $A_2$, $A_1$ and $A_3$, $A_2$ and $A_3$.

**Several events are called independent in the aggregate** (or **independent**) if they are pairwise independent, and are independent of each and all possible products of others. For example, if events $A$, $B$, $C$ are independent in the aggregate, the independent events $A$ and $B$, $A$ and $C$, $B$ and $C$, $A$ and $BC$, $B$ and $AC$, $C$ and $AB$. It turns out that several independent pairs of events does not guarantee their independence together.

**Corollary 3**. The probability of the product of $k$ events, independent in aggregate, is equal to the product of the probabilities of these events:

$$P(A_1 A_2 ... A_k) = P(A_1)\ P(A_2)\ ...\ P(A_k). \tag{2.7}$$

**3.Theorem 1**. The probability of the sum of two disparate events is the sum of the probabilities of these events:

$$P(A + B) = P(A) + P(B). \tag{2.8}$$

**Corollary**. The probability of the sum of $k$ pairwise incompatible events is the sum of the probabilities of these events:

$$P(A_1 + A_2 + ... + A_k) = P(A_1) + P(A_2) + ... + P(A_k). \tag{2.9}$$

**Theorem 2**. The probability of the sum of two joint events is the sum of the probabilities of these events without the probability of their joint occurrence:

$$P(A + B) = P(A) + P(B) - P(AB). \tag{2.10}$$

**Theorem 3**. The sum of the probabilities of events that form a complete group is unity:

$$P(A_1) + P(A_2) + ... P(A_n) = 1. \tag{2.11}$$

For n = 2 events $A_1$ and $A_2$ are opposite, so the equality (2.11) takes the following form:

$$P(A) + P(\overline{A}) = 1 \quad \text{or} \quad p + q = 1, \tag{2.11*}$$

where p = P (A).

***The probability occurring of at least one event***
Let the test result may occur events $A_1, A_2, ..., A_k$, A1, the probability of occurrence of each of which are known and which are independent in the aggregate. We denote $A$ - the emergence of at least one of these events in the trial. Then according to the definition of the sum of events $A = A_1 + A_2 + ... + A_k$. Since the events $A_1, A_2, ..., A_k$ are compatible, the addition theorem of probability does not "work" for k ≥ 3. Using (2.11 *), (2.7) and $\overline{A} = \overline{A}_1 \overline{A}_2 ... \overline{A}_k$ obtain:

$$P(A) = 1 - q_1 q_2 ... q_k, \tag{2.12}$$

where $q_1 = P(\overline{A}_1), q_2 = P(\overline{A}_2), ..., q_k = P(\overline{A}_k)$.

In the particular case where $P(A_1) = P(A_2) = ... = P(A_k) = p$, obtained the following formula

$$P(A) = 1 - q^k, \tag{2.12*}$$

де $q = 1 - p$.

Finally, if the events $A_1, A_2, ..., A_k$ are dependent (not possess the independence property in total), then

$$P(A) = 1 - P(\overline{A}_1) P_{\overline{A}_1}(\overline{A}_2) ... P_{\overline{A}_1 \overline{A}_2 ... \overline{A}_{k-1}}(\overline{A}_k). \tag{2.13}$$

5. The algorithm for solving problems using addition and multiplication theorems of probability.
1) Introduce into consideration the event, the probability of which must be found, and as more simple events, probabilities it are known or can be found in the classical definition.
2) "The required" random event (the probability to find what you need), expressed in terms of simpler events using algebra of events, transactions that amount to the product, the negation (opposite event). It should be guided by the mnemonic rule «+» ↔ or, «×» ↔» and.
3) Depending on the type of the resulting expression using the addition theorem of probability either (i) multiplying the probability theorem and its consequences. In implementing this paragraph is necessary to find out the

properties of events (compatibility, incompatibility, dependence, independence, or completeness of opposite pairs or groups of events).
Remarks. Keep in mind that many problems of implementation of paragraph

## 6. The formula of total probability

Let event A can occur only if the appearance of one of the events $B_1$, $B_2$, ..., $B_n$, which form a complete group. Let the known probability of $P(B_k)$, $P_{B_k}(A)$, $k = (\overline{1,n})$. Answering the question: how to find $P(A)$? - gives

**Theorem. The probability of event $A$, which can happen only after the appearance of one of the events $B_1$, $B_2$, ..., $B_n$, which form a complete group is a so-called formula of total probability:**

$$P(A) = P(B_1)P_{B_1}(A) + P(B_2)P_{B_2}(A) + ... + P(B_n)P_{B_n}(A). \quad (2.14)$$

***Remarks***. Since before the test is unknown, after which events $B_1$, $B_2$, ..., $B_n$ event $A$ occurs, therefore this events are called hypotheses (assumptions).

### Bayes' Formula

Let the conditions of theorem respect to events $B_1$, $B_2$, ..., $B_n$ and $A$. Suppose that conducted tests in which the event $A$ occurred. The question arises: how "overestimated" the probability hypotheses $B_1$, $B_2$, ..., $B_n$ with taking into account the departure of event A, i.e. find the conditional probability $P_A(B_k)$, $(k = \overline{1,n})$? Answer given so-called Bayes' formula:

$$P_A(B_k) = \frac{P(B_k)P_{B_k}(A)}{\sum_{i=1}^{n} P(B_i)P_{B_i}(A)}, k = 1,2,...,n. \quad (2.15)$$

## 7. Algorithm for solving problems using formulas of total probability and Bayes.

We recommend a sequence of solving problems.
1) Formulate hypotheses $B_1$, $B_2$, ..., $B_n$ and event $A$. In this case, check the completeness of hypotheses, and the fact that the event $A$ can occur only after the appearance of one of the hypotheses.
2) The probability of hypotheses are calculated. The correctness of the calculation is controlled execution of equality
$P(B_1) + P(B_2) + ... + P(B_n) = 1$. Calculated conditional probabilities $P_{B_1}(A), P_{B_2}(A),..., P_{B_n}(A)$.

3) Select the formula of total probability and Bayes' formula. The latter are used when there is information about the departure of random events.

## § 3. Repeated independent trials
1. Bernoulli trial.
2. Likely number of an event.
3. Local Laplace formula.
4. Poisson.
5. Laplace integral formula.
6. The probability of deviation of the relative frequency of the event from its wait-term probability.
7. The algorithm for solving repeated independent trials-you.

In practice, often there are cases when multiple tests are possibly. These trials are called repeating , and their collection - the scheme repeating trials or Bernoulli scheme. The same random event $A$ can occur in every of these trials. If $P(A)$ remains the same for each test, these tests will be called independent (with respect to event A).

## 1. Bernoulli formula
A **binomial experiment** has the following characteristics:
1. The experiment consists of $n$ identical trials.
2. Each trial results in a success or a failure.
3. The probability of a success on any trial is $p$ and remains constant from trial to trial. This implies that the probability of failure, $q$, on any trial is $1 - p$ and remains constant from trial to trial.
4. The trials are independent (that is, the result of the trials have nothing to do with each other). Furthermore, if we define the random variable

**Theorem. The probability that an event $A$ happens exactly $m$ times in $n$ independent repeated trials, is calculated with Bernoulli formula:**
$$P_n(m) = C_n^m p^m q^{n-m}, \tag{3.1}$$

where $C_n^m$ - number of combinations, that is defined by formula (1.4) $p$ - probability of event A in one trial $(p = P(A))$, $q$ - the probability of disclosure event A in one trial $(q = P(\overline{A}))$.

## 2. Likely number of an event.
In the $n$ repeated independent trials event $A$ can occur number of times
$$0, 1, ..., m_0-1, m_0, m_0+1, ..., n \tag{3.2}$$

with corresponding probabilities (which can be found Bernoulli equation):

$$P_n(0), P_n(1), ..., P_n(m_0-1), P_n(m_0), P_n(m_0+1), ..., P_n(n). \qquad (3.3)$$

**Definition. Number $m_0$ in the sequence (3.2) is called most probabilitive number of an event in $n$ independent trials (or moda) if it meets the highest probability in sequence (3.3).**

Most probabilitive number can be found on this double inequality:

$$np - q \le m_0 \le np + p. \qquad (3.4)$$

Since the difference between the right and left side of this inequality is 1, it can be concluded that the maximum number of most likely event equal to two.

### 3. Local Laplace formula.

Using Bernoulli's equation for large values of $n$ is difficult. For example, for $p = 0{,}2$, $q = 0{,}8$ $P_{50}(30) = C_{50}^{30}(0{,}2)^{30}(0{,}8)^{20}$, where $C_{50}^{30} = 4712921 \cdot 10^7$.

One of the approximate realization of the right part of the formula (3.1) gives

**Local Laplace theorem (Moivre-Laplace).**

If the probability $p$ of occurrence the random event $A$ in each of the $n$ repeated trials remains unchanged (with $0 < p < 1$), and the number of trials is enough large, then probability that in $n$ trials event occurs $m$ times, is calculated with the approximate formula

$$P_n(m) \approx \frac{1}{\sqrt{npq}}\, \varphi(x), \qquad (3.5)$$

where $\varphi(x) = \dfrac{1}{\sqrt{2\pi}}\, e^{-x^2/2}$ - Gauss function, $x = \dfrac{m - np}{\sqrt{npq}}$, $q = 1 - p$.

***Note.*** Equation (3.5) is locally Laplace equation. Its accuracy increases with $n$.

Gauss function is table function and its value are shown in Table. 1 applications. To properly use this table should know Gauss function following properties: 1) $\varphi(x)$ is defined for all $x \in R$; 2) $\varphi(x)$ - even function $(\varphi(-x) = \varphi(x))$; 3) for positive $x$ $\varphi(x)$ very quickly tends to 0 as $x$ increases, particularly $\varphi(3{,}99) = 0{,}0001$. Based on these properties in the table shows the values of Gauss function for $x$ in the interval [0; 5]. **Practically we can assume that the local Laplace formula gives a good approximation if $npq > 9$.** If the requirements for accuracy value is higher,

it should require the inequality $npq \geq 25$.

**1.    Poisson formula.**

Inequality $npq > 9$ even for large n can not be performed (and hence the error when using local Laplace formula will do the very large) in the case of rare (unlikely) event $A$, that is, those events for which $p$ was significantly less than 0.1 ($p \ll 0.1$). In such cases, should be used an another approach right side Bernoulli formula. One of them is given this statement.

**Poisson                                                                    theorem.**
**If each of *n* repeated trials, the probability *p* of an event *A* is constant and small (p ≪ 0.1), and the number of tests *n* is rather significant, then the probability that event *A* will   come in this trials exactly *m* times, is given by**

$$P_n(m) \approx \frac{\lambda^m e^{-\lambda}}{m!}, \qquad\qquad (3.6)$$

where                                                                        $\lambda = np.$

*Note*. The error in the approximate equality (3.6), which is named Poisson formula, the lower, the higher the number of tests $n$.

The value function $P(m) = \dfrac{\lambda^m e^{-\lambda}}{m!}$ of two variables $\lambda$ and $m$ for some $m$

and    $\lambda$    are    shown    in    Table.    2    applications.
Note that the Poisson formula is also used to number appearance of disclosure event $A$ if $q \ll 0,1$, and $nq$ small.


## 5. Laplace integral formula.

**The integral theorem of Laplace (Moivre-Laplace).**
**If the probability *p* of an event *A* in each of the n repeated trials is constant (0 <*p* <1), and the number of trials is large enough, then the probability that event *A* in these tests will take place at least *m1* times and no more *m2* times, is the following approximate equation (Laplace integral formula):**

$$P_n(m_1 \leq m \leq m_2) \approx \Phi(x_2) - \Phi(x_1), \qquad\qquad (3.7)$$

where $\Phi(x) = \dfrac{1}{\sqrt{2\pi}} \int\limits_0^x e^{-t^2/2} dt$ - Laplace function,

$$x_1 = \frac{m_1 - np}{\sqrt{npq}}, x_2 = \frac{m_2 - np}{\sqrt{npq}}.$$

14

The function of the Laplace is table function and its values are given in Table. 3 applications. When using this table must take into account the properties of Laplace function:
1) $\Phi(x)$ is defined for all $x \in R$;

2) $\Phi(x)$ is an odd function $(\Phi(-x) = -\Phi(x))$;
3) $\Phi(x)$ increases monotonically for all $x \in R$, with $y = -0,5$ - left-sided asymptote and $y = 0,5$ – right-sided;
4) the rate of growth $\Phi(x)$ on the interval [0; 5] is very high ( $(\Phi 5) = 0.499997$), so for all $x > 5$ with negligible error $\Phi(x) \approx 0,5$.

Precision integrated Laplace formula is greater, the greater the number of trials $n$. As a local, integrated formula gives a good approximation if $npq > 9$. If the requirements are much higher accuracy, it is necessary to require the inequality $npq \geq 25$.

## 6. The probability of deviation relative frequency of events from her constant probability.
**Theorem. If the probability $p$ of occurrence, the random event $A$ in each of the $n$ repeated trials is constant, and the number of trials is large enough, then the probability that the absolute deviation of the relative frequency $m / n$ event $A$ from its probability $p$, is not exceed given number $\varepsilon > 0$ , is given by**
$$P\left(|m/n - p| \leq \varepsilon\right) \approx 2\Phi\left(\varepsilon\sqrt{n/pq}\right). \tag{3.8}$$

## SOLVING COMMON PROBLEMS

**Example**

Suppose that historical sales records indicate that 40 percent of all customers who enter a discount department store make a purchase. What is the probability that two of the next tree customers will make a purchase?

In order to find this probability, we first note that the experiment of observing three customers making a purchase decision has several distinguishing characteristics:
1.  The experiment consists of three identical trials; each trial consists of a customer making a purchase decision.
2.  Two outcomes are possible on each trial: the customer makes a purchase (which we call a success and denote as S), or the cus-

tomer does not make a purchase (which we call a failure and denote as F).

3.   Since 40 percent of all customers make a purchase, it is reasonable to assume that P(S), the probability that a customer makes a purchase, is .4 and is constant for all customers. This implies that P(F), the probability that a customer does not make a purchase, is .6 and is constant for all customers.

4.   We assume that customers make independent purchase decision. That is, we assume that the outcomes of the three trials are independent of each other.

It follows that the sample space of the experiment consists of the following eight sample space outcomes:

| | |
|---|---|
| SSS | FFS |
| SSF | SFS |
| FSF | SFF |
| FSS | FFF |

Here the sample space outcome SSS represents all three customers making purchases. On the other hand, the sample space outcome SFS represents the first customer making a purchase, the second customer not making a purchase, and the third customer making a purchase.

Two out of three customers make a purchase if one of the sample space outcomes SSF, SFS, or FSS occurs. Furthermore, since the trials (purchase decision) are independent, we can simply multiply the probabilities associated with the different trial outcomes (each of which is or F) to find the probability of a sequence of outcomes:

$$P(SSF) = P(S)P(S)P(F) = (.4)(.4)(.6) = (.4)^2(.6)$$
$$P(SFS) = P(S)P(F)P(S) = (.4)(.6)(.4) = (.4)^2(.6)$$
$$P(FSS) = P(F)P(S)P(S) = (.6)(.4)(.4) = (.4)^2(.6)$$

It follows that the probability that two out of the next three customers make a purchase is

$$P(SSF) + P(SFS) + P(FSS) = (.4)^2(.6) + (.4)^2(.6) + (.4)^2(.6) = 3(.4)^2(.6) = .288$$

We can now generalize the previous result and find the probability that $x$ of the next $n$ customers will make a purchase. Here we will assume that $p$ is the probability that a customer makes a purchase, $q = 1 - p$ is the probability that a customer does not make a purchase, and purchase decisions (tri-

als) are independent. To generalize the probability that two out of the next three customers make a purchase, which equals

$$3(.4)^2(.6)$$

we note that

2. The 3 in this expression is the number of sample space outcomes (SSF, SFS, and FSS) that correspond to the event "two out of the next three customers make a purchase". Note that this number equals of ways we can arrange two successes among the three trials.
3. The .4 is $p$, the probability that customer makes a purchase.
4. The .6 is $q = 1 - p$, the probability that a customer does not make a purchase.

Therefore, the probability that two of the next three customers make a purchase is

(The number of ways to arrange 2 successes among 3 trials) $p^2 q^1$

## § 4. DISCRETE RANDOM VARIABLES AND THEIR NUMERIC CHARACTERISTIC

1. *1. Random variables and their types. The law of probability distribution of a discrete random variable. 2. Basic discrete distributions (integer) random values (uniform, binomial, Poisson, geometric, hypergeometric). The simplest flow. 3. Acts of random variables. 4. Numerical characteristics of discrete random variables and their properties (expectation, variance, standard deviation). 5. Numerical characteristics of the basic laws of distribution.*

**1.** Random value is called value, which when tested shall take a single value from all possible with some probability, i.e. known in advance which specific possible value it enters, as it depends on accidental causes. Discrete (discontinuous) is a random variable, whose possible values are isolated numbers. The number of possible values of a discrete random variable can be finite or countable. In the latter case, you can set a one-to-one correspondence between the possible values and the natural numbers 1, 2, 3, ..., $n$, ... .

Continuous is a random variable, whose value may completely fill a finite or infinite interval. Obviously, the number of possible values for each continuous is infinite value.

*Note*. The definition of continuous random variable is tentative. Refinement is done in the next subsection.

Information on the set of possible values is not enough to describe random variable (different values may be have the same possible values). We need to know probability the possible values of a random variable.

**Law probability distribution of a discrete random variable** is called the correspondence between the possible values and probabilities with which they recruited a random variable.

Tabular format setting of the distribution has this type

$$\begin{array}{c|cccc} X & x_1 & x_2 & \dots & x_n \\ \hline P & p_1 & p_2 & \dots & p_n \end{array}, \tag{4.1}$$

where $p_i = P(X = x_i)$, $i = \overline{1,n}$. Because in one trial random variable gaining only one of its possible values, the random events

$(X = x_1)$, $(X = x_2)$, ..., $(X = x_n)$ form a complete group. Therefore, the sum of probabilities is unity:

$$p_1 + p_2 + \dots p_n = 1. \tag{4.2}$$

This equality is called normalization condition.

If the set of possible values of a discrete random value is countable:

$x_1, x_2, \dots, x_n$, then row $\sum\limits_{t=1}^{\infty} p_i$ converges and its amount is equal to one.

## 2. Basic discrete distributions (integer) random variables.

The law of distribution of a discrete random variable X can be set as analytically, that is, from the formula $p_i = P(X = x_i) = g(x_i)$, $i = \overline{1,n}$. All laws are given below are analytically.

• Integer random variable **distributed according to a uniform law (evenly distributed)** if the probability in the law of distribution have the following form: $p_k = P(X = k) = 1/n$, $k = \overline{1,n}$.

• Distribution law integer random variable, the probability of which are Bernoulli formula, called the **binomial**:

$$p_{m+1} = P(X = m) = C_n^m p^m q^{n-m}, \tag{4.3}$$

where $m = 0, 1, 2, \dots, n$; $q = 1 - p$.

• Distribution law integer random variable, the probability of which are for Poisson:

$$p_{m+1} = P(X = m) = \frac{\lambda^m e^{-\lambda}}{m!}, \quad m = 0, \ 1, \ 2, \ ..., \ n,$$

called Poisson ($n$ - great $p << 0,1, \ \lambda = np \le 9$).

• The law of probability distribution, which is determined by the sequence of probability:

$$p_1 = p, \quad p_2 = qp, \quad p_3 = q^2 p, \quad ..., \quad p_n = q^{n-1} p, \ ... \tag{4.4}$$

called geometric.

• The law of probability distribution, which is given by

$$\bullet \quad P(X = m) = C_M^m \cdot C_{N-M}^{n-m} / C_N^n, \quad m = 0, \ 1, \ 2, \ ..., \ n. \tag{4.5}$$

called hypergeometric, where $N$ - total number of objects, including $M$ have some feature ($M < N$), while randomly selected products $n$ ($n \le M$) without the return of each of them.

## The simplest flow of events

**Flow of events** called the sequence of events that occur at random times.
**The simplest** is called the **flow of events**, which has the properties of **stationary, no aftereffect and ordinariness.**
The property of stationarity is that the probability of $m$ events flow for any length of time interval $t$ depends only on m and $t$, and is independent of the timing; with different time intervals must not overlap.
Property no aftereffect determines that the probability of $m$ events at arbitrary time interval is independent of the occurrence or event of disclosure flow at time points prior to the beginning of the period.
**Intensity $\lambda$ flow** is called the average number of events that occur per unit time.
Mathematical model of simple flow of events is Poisson tormular
$$P_t(m) = (\lambda t)^m e^{-\lambda t} / m!, \tag{4.6}$$
through which it is possible to find the probability of $m$ simple flow of events during the time interval length $t$.

## 3. Acts of random variables.
We define the product **of a constant C in a discrete random variable X** given distribution law (4.1) as a discrete random variable $CX$, the law of distribution which has the form

$$\frac{CX \mid Cx_1 \quad Cx_2 \quad ... \quad Cx_n}{P \mid p_1 \quad p_2 \quad ... \quad p_n}. \tag{4.7}$$

**Square random variable** $X$, i.e. $X^2$ - a new discrete random variable which is described by distribution

$$\frac{X^2 \mid x_1^2 \quad x_n^2 \quad ... \quad x_n^2}{P \mid p_1 \quad p_2 \quad ... \quad p_n}. \tag{4.8}$$

Let the random variable $X$ given distribution (4.1), and $Y$ – law of

distribution $\dfrac{Y \mid y_1 \quad y_2 \quad ... \quad y_m}{P \mid g_1 \quad g_2 \quad ... \quad g_m}$. These values are called **independent** if

random events $(X = x_i)$, $(Y = y_j)$ are independen in arbitrary $i = \overline{1,n}$ and $j = \overline{1,m}$. In case zhnomu protyle-called dependent random variables.

Discrete random variables $X_1$, $X_2$, ..., $X_k$ are called **independent together** (**mutually independent**) if the law of distribution of each of them does not change if an arbitrary random variable or any group of these variables which will take anything from their possible values. In otherwise random variables called **dependent**.

Determine the **sum of random variables** $X$ and $Y$ as a random variable $X + Y$, the possible values which equal sum of each possible value of $X$ each possible value of $Y$, and likelihood of possible values of $X + Y$ for independent variables $X$ and $Y$ is the product of the probabilities of terms; for dependent variables - product of the probabilities of one term in the conditional probability of a second. If some sum $x_i + y_j$ are equal, then the probability of the possible value of the sum is the sum of the corresponding probabilities.

We define the **product of independent random variables** $X$ and $Y$ as a random variable $XY$, whose possible values equal to the product of each possible value of $X$ for each possible value of Y; probability of possible values of the product $XY$ is the product of the probabilities of possible values of the factors. In case of equality of products $x_i y_j$ probability $XY$ possible value equal to the sum of corresponding probabilities.

## 4. Numerical characteristics of discrete random variables and their properties.

The law of probability distribution gives complete information about discrete random variable. However, in many cases of practical importance economist-researcher is enough to know one or more numbers associated with a random variable, which provide less complete, but a visual

description of a random variable. These numbers are totally random variable describing, called its numerical characteristics.

## Mathematical expectation

**Mathematical expectation** of a discrete random variable $X$ is the sum of products of all possible values of the relevant probabilities:

$$M(X) = x_1p_1 + x_2p_2 + ... + x_np_n. \qquad (4.9)$$

If the value of $X$ may take countable set of values (such as in the case of geometric or Poisson distribution), then assuming that this series is absolutely convergent,

$$M(X) = \sum_{i=1}^{\infty} x_i p_i. \qquad (4.10)$$

Probabilistic content expectation: **$M(X)$ is approximately equal** (more precisely, the greater the number of tests) the **arithmetic mean of observed values of a random variable.**

*Properties expectation.*

1. **Expected a constant equal to this constant**:
$$M(C) = C.$$

2. **Constant multiplier can be taken before a sign expectation**:
$M(CX) = CM(X).$

3. **Expected sum of two random variables is the sum of the mathematical expectations of terms:**
$M(X + Y) = M(X) + M(Y).$

**Corollary. Expected sum of several random variables is the sum of the mathematical expectations of terms.**

4. **Expected product of two independent random variables is the product of their mathematical expectations:**
$$M(XY) = M(X) \cdot M(Y).$$

**Corollary. Expected product of several independent random variables in the aggregate equal to the product of their mathematical expectations.**

# Dispersion

To assess the spread of possible values of a random variable around average value (expectation) is used numerical characteristic, which is called **dispersion**.

**Dispersion** (spread) of a random variable $X$ is the expectation of the square of the deviation of $X$ from $M(X)$:

$$D(X) = M[X - M(X)]^2. \qquad (4.11)$$

According to this definition characterizes the dispersion of the average of the spread of possible values of the random variable around its expected value (mean) in square units.

If $X$ is distributed according to (4.1), then with regard to the law (4.8) random variable $[X - M(X)]^2$ has the next distribution law:

| $[X - M(X)]^2$ | $[x_1 - M(X)]^2$ | $[x_2 - M(X)]^2$ | ... | $[x_n - M(X)]^2$ |
|---|---|---|---|---|
| $P$ | $p_1$ | $p_2$ | ... | $p_n$ |

Using the definition of the expectation value and dispersion is achieved formula for calculating the dispersion $D(X)$:

$$D(X) = M[X - M(X)]^2 = [x_1 - M(X)]^2 p_1 + [x_2 - M(X)]^2 p_2 +$$
$$+ ... + [x_n - M(X)]^2 p_n = \sum_{i=1}^{n} [x_i - M(X)]^2 p_i. \tag{4.11*}$$

Reducing the amount of computation is achieved by using **calculation formula** for calculating the variance:

$$D(X) = M(X^2) - [M(X)]^2. \tag{4.12}$$

Applying the definition to the expectation of the distribution (4.8), we obtain the **numerical implementation of the calculation formula**:

$$D(X) = \sum_{i=1}^{n} x_i^2 p_i - [M(X)]^2. \tag{4.12*}$$

**Properties dispersion.**

1. **Dispersion constant is zero**: $D(C) = 0$.

2. **Constant multiplier can make a sign variance, having presented it to the square**: $D(CX) = C^2 D(X)$.

3. **Dispersion sum of two independent random variables is the sum of the variances of these values**: $D(X + Y) = D(X) + D(Y)$.

**Corollary. Dispersion sum of several independent random variables in the aggregate the sum of the variances of these quantities.**

4. **Dispersion difference of two independent random variables is the sum of their variances**: $D(X - Y) = D(X) + D(Y)$.

*Note*. Even if $X$ and $Y$ are independent, in the general case the inequality $D(XY) \neq D(X) D(Y)$ is true.


**Standard deviation**

The disadvantage of using dispersion in some cases due to the fact that it has a dimension equal to the dimension of the square of the random variable. For example, if the possible values of $X$ are measured in kg, $D(X)$ in $(kg)^2$.

**Standard deviation** of a random variable $X$ is the square root of the variance:

$$\sigma(X) = \sqrt{D(X)}. \qquad (4.13)$$

$\sigma(X)$ characterizes the average of the possible spread of values $X$ around $M(X)$ (middle) in linear units.

5. **Numerical characteristics of the basic laws of distribution** Numerical characteristics of random variable, distributed by **binomial law**: $M(X) = np, \quad D(X) = npq.$ (4.14)

Numerical characteristics of random variable, distributed by
**Poisson:** $\qquad M(X) = D(X) = np = \lambda.$ (4.15)

## 5. The continuous random variables And their numeric characteristics

1. The probability distribution function and its properties.
2. The density probability distribution and its properties.
3. Numerical characteristics of continuous random variables (expectation, variance and standard deviation, mode and mean random variable).

1. **The probability distribution function** is a function $F(x)$ determined (non-random) argument $x$, which is numerically equal to the probability that a test result of a random variable $X$ attains values less than $x$, i.e.

$$F(x) = P(X < x). \qquad (5.1)$$

Sometimes the distribution function is called integrated. Clarify the definition of continuous random variables: random variable is called continuous if its distribution function is continuous on its definition, derivative from function distribution is continuous at all points, except, perhaps, finite number of points on arbitrary finite interval.

*Properties probability distribution function*

**1. The domain of the distribution function -** $R^1 = (-\infty; \infty),$ **the range of values - the interval [0; 1].**

**2. $F(x)$ - nondecreasing function, that is, for any pair of numbers $x_1, x_2$ of inequality $x_2 > x_1$ implies inequality $F(x_2) \geq F(x_1)$.**

**Corollary 1.** The probability that a random variable when testing will take possible value in the interval [a; b), is growth the distribution function of in this interval:

$$P(a \leq X < b) = F(b) - F(a). \tag{5.2}$$

**Corollary 2.** The probability that a continuous random variable $X$ will come when tested one possible value is zero.

***Note.*** The equality $P(X = x_1) = 0$, where $x_1$ - specific possible value of $X$ does not mean that the event $X = x_1$ is not possible (as opposed to the classical definition of probability). Remind geometrical definition of probability.

Corollary 2 using formula (5.2) and the theorem of probability allows you to add a chain of equalities:

$$P(a < X < b) = \big(P(a \leq X \leq b) = P(a \leq X < b) = P(a < X \leq b)\big) = $$
$$= F(b) - F(a). \tag{5.3}$$

**3. If all possible values of the random variable belongs interval[a; b], then $F(x) = 0$ for all $x \leq a$ and $F(x) = 1$ for all $x > b$.**

**Corollary. If the possible values of a continuous random variable is all real numbers, then there are marginal value:**

$$\lim_{x \to -\infty} F(x) = 0, \quad \lim_{x \to \infty} F(x) = 1.$$

2. The density of the probability distribution of continuous random value $X$ is a function $f(x)$, which is the first derivative of the distribution function $F(x)$:

$$f(x) = F'(x). \tag{5.4}$$

From the definition (5.4) it follows that $F(x)$ is the initial for density. Below is the formula for finding $F(x)$, if known function $f(x)$. Thus, setting one of the functions $f(x)$ and $F(x)$ allows to find another. Therefore, in the literature, these functions are called the **laws of distribution of continuous random variables.**

*Properties of the density distribution*

1. **The domain of the function $f(x)$ — $R^1$, and the range of values interval [0, ∞).**

2. **Improper integral of the density distribution in the range of -∞ to ∞ is unity:**

24

$$\int\limits_{-\infty}^{\infty} f(x)dx = 1. \qquad\qquad (5.5)$$

Equality (5.5) is called the **normalization condition** of continuous random variables.

*Note*. If $[\alpha, \beta]$ - the minimum amount, which contains all the possible values of a continuous random variable, while normalization condition is

$$\int\limits_{\alpha}^{\beta} f(x)dx = 1. \qquad\qquad (5.5^*)$$

There        are        the        following        formula:

$$P\left(a \overset{(\leq)}{<} X \overset{(\leq)}{<} b\right) = \int\limits_{a}^{b} f(x)dx, \qquad\qquad (5.6)$$

$$F(x) = P(X < x) = P(-\infty < X < x) = \int\limits_{-\infty}^{x} f(x)dx. \qquad (5.7)$$

Let us determine the probability content of density probability distribution. Definition *f(x)* (5.4) and the derivative of the function, and equality (5.3) give the following equality

$$f(x) = \lim_{\Delta x \to 0} \frac{F(x + \Delta x) - F(x)}{\Delta x} = \lim_{\Delta x \to 0} \frac{P(x < X < x + \Delta x)}{\Delta x}, \quad (5.8)$$

that let you capture approximate equality

$$P(x < X < x + \Delta x) \approx f(x)\Delta x. \qquad\qquad (5.9)$$

up to a  infinitesimal of higher order relative $\Delta x$.

The right side of (5.8) explains the name of the function *f(x)*  by analogy with the contents density mass distributed on the interval, and the approximate equality (5.9) - "contribution" value *f(x)* in the value of the probability of events $(x < X < x + \Delta x)$  for each segment length $\Delta x$.

3. **Mathematical expectation of continuous random value** $X$ , all possible values which are finite interval [a; b], called the definite integral

$$M(X) = \int_a^b xf(x)dx. \tag{5.10}$$

If the possible values of $X$ belongs $R^1$, then

$$M(X) = \int_{-\infty}^{\infty} xf(x)dx, \tag{5.10*}$$

where the assumption of improper integrals coincide completely (absolute).

As for the case of discrete values, **dispersion continuous random variable** is called the expected value of the square of the deviation from the expectation. In view of the definition of $M(X)$:

$$D(X) = \int_a^b [x - M(X)]^2 f(x)dx, \tag{5.11}$$

if all possible values of $X$ are finite interval $[a; b]$, and

$$D(X) = \int_{-\infty}^{\infty} [x - M(X)]^2 f(x)dx, \tag{5.11*}$$

if possible value of $X$ is filled $R^1$.

**Formulas for calculating the variance**:

$$D(X) = \int_a^b x^2 f(x)dx - [M(X)]^2, \tag{5.12}$$

$$D(X) = \int_{-\infty}^{\infty} x^2 f(x)dx - [M(X)]^2. \tag{5.12*}$$

Standard deviation of continuous random value defined as for the discrete case, equality

$$\sigma(X) = \sqrt{D(X)}. \tag{5.13}$$

Mode $Mo(X)$ of a discrete random variable $X$ is that its possible value, which corresponds to the probability of its occurrence. For continuous random variable $X$ fashion mode $Mo(X)$ is a possible value of $X$, which corresponds to a local maximum density distribution. The random variable can have multiple modes. There are also distributions, with no maximum. Such distributions are called antymodal.

The median $Me(X)$ of continuous random variable $X$ is its possible value for which the equality $P(X < Me(X)) = P(X > Me(X))$ is true.

§ 6. THE BASIC LAW CONTINUOUS RANDOM VARIABELS

1. *Normal Law (Probabilistic content settings distribution, normal curve,*

26

*the probability of getting at a given interval; the probability of a given deviation).*
*2. Law uniform distribution.*
*3. Exponential law.*

1. The random variable is called **normally distributed** (distributed accordingly the normal law) if it has the density distribution of this type

$$f(x) = \frac{1}{\sigma\sqrt{2\pi}} e^{\frac{-(x-a)^2}{2\sigma^2}}. \qquad (6.1)$$

The density distribution (6.1) is completely determined by two parameters: $a$ and $\sigma$. Probabilistic content of these parameters determined by equalities:

$$a = M(X), \qquad \sigma = \sqrt{D(X)}. \qquad (6.2)$$

Normal distribution random variable with parameters $a=0$, $\sigma=1$ is called a **standard** or **normalized**.

Figure density of the normal distribution curve is called a normal (Gaussian curve) (pic.6.1).



Pic. 6.1.

The probability that the normal random variable during tests attains significance in the interval $(\alpha, \beta)$ is calculated by the formula

$$P(\alpha < X < \beta) = \Phi\left(\frac{\beta - a}{\sigma}\right) - \Phi\left(\frac{\alpha - a}{\sigma}\right). \qquad (6.3)$$

In practically important problems there is a need for calculation probability that absolute value deviations normally distributed value from $a$ is less than a given positive number $\varepsilon$, that is $P\left(|X - a| < \varepsilon\right)$.

This probability is calculated by the formula

$$P\left(|X - a| < \varepsilon\right) = 2\Phi\left(\frac{\varepsilon}{\sigma}\right), \qquad (6.4)$$

27

which, in particular, shows that the probability of error for the same $\varepsilon$ will be greater, the smaller the $\sigma$.

**2.** The random variable is called a **distributed accordingly the unit law** ev (unit distribution) if its density distribution has this type

$$f(x) = \begin{cases} C = \text{const}, & \text{якщо } x \in [a,b], \\ 0, & \text{якщо } x \notin [a,b]. \end{cases}$$

We find the constant C, using a second property density distribution. The geometric properties of the contents of this means equality unit area bounded by the curve of distribution, i.e. $C(b-a) = 1$, where $C = 1/(b-a)$. Therefore, the density of uniformly distributed random variable has the form:

$$f(x) = \begin{cases} 1/(b-a), & \text{при } x \in [a,b], \\ 0, & \text{при } x \notin [a,b]. \end{cases} \tag{6.5}$$

The interval [a, b] is called a **segment concentration** unit distribution.

3. The random variable X is distributed by the **exponential** law (exponential distributed) when its density distribution of probability has the

form $f(x) = \begin{cases} \lambda e^{-\lambda x}, & \text{якщо } x \geq 0, \\ 0, & \text{якщо } x < 0, \end{cases}$ \hfill (6.6)

which constant $\lambda > 0$ is called parameter exponential distribution.

## § 7. LAW OF LARGE NUMBER
*1. Lemma and Chebyshev inequality.*
*2. Chebyshev theorem (medium resistance).*
*3. Bernoulli's theorem (stability relative frequencies).*
*4. The Central Lyapunov Limit Theorem.*

**1. Lemma Chebyshev. If all possible values of the random variable $Y$ non-negative, then the probability that it is in triales attains value greater than from a positive number $b$, not much of a fraction, the numerator of which - mathematical vaiting of $Y$, and the denominator - the number of $b$ :**

$$P(Y > b) \le M(Y)/b. \qquad (7.1)$$

***Chebyshev inequality. The probability that the random deviation variable $X$ from its expectation in absolute value no greater than the number of positive $\varepsilon$ not less than $1 - D(X)/\varepsilon^2$ :***

$$P\big(|X - M(X)| \le \varepsilon\big) \ge 1 - D(X)/\varepsilon^2. \qquad (7.2)$$

***Note***. If the random variable X is distributed according Binominal law (*X* - number of an event in *n* repeated independent trials), then Chebyshev inequality come like this:

$$P\left(\left|m - np\right| \leq \varepsilon\right) \geq 1 - \frac{npq}{\varepsilon^2} \qquad (7.3)$$

or

$$P\left(\left|\frac{m}{n} - p\right| \leq \varepsilon\right) \geq 1 - \frac{pq}{n\varepsilon^2}, \qquad (7.4)$$

where $m$ - number of an event $A$ in $n$ repeated independent trials, $p = P(A)$, $m/n$ - relative frequency of an event A.

2. **Chebyshev theorem. If** $X_1, X_2, ..., X_n$ **pairwise independent random variables variance which is uniformly bounded** $\left(D(X_i) \le C, i = \overline{1, n}\right)$, **then for arbitrary ε> 0 and sufficiently high $n$ probability event**

$$\left|\frac{X_1 + X_2 + ... + X_n}{n} - \frac{M(X_1) + M(X_2) + ... + M(X_n)}{n}\right| \le \varepsilon \qquad (7.5)$$

**will be arbitrarily close to unity. When installed proof theorems is established correctness of inequality**

$$P\left(\left|\frac{X_1 + X_2 + ... + X_n}{n} - \frac{M(X_1) + M(X_2) + ... + M(X_n)}{n}\right| \le \varepsilon\right) \ge 1 - C/(n\varepsilon^2). \qquad (7.6)$$

**3. Bernoulli's theorem. If in each of the $n$ repeated trials the probability $p$ of an event $A$ is constant, then was then arbitrarily close to unity probability that the deviation relative frequency event $A$ from the probability $p$, in the absolute value is arbitrarily small if the number of trials is sufficiently large.**

From inequality (7.6) can be obtained inequality

$$P\left(\left|m/n - p\right| \le \varepsilon\right) \ge 1 - \frac{0{,}25}{n\varepsilon^2}. \qquad (7.7)$$

In Chebyshev theorem has not been used for information about laws distribution of random variables. However, for the problems of the theory and practice important question is: by what law the sum of sufficiently large number of random variables is distributed?

The random variable

$$Z_n = \left[\sum_{i=1}^{n} X_i - \sum_{i=1}^{n} M(X_i)\right] \Big/ \sqrt{\sum_{i=1}^{n} D(X_i)}. \qquad (7.8)$$

will be called **the normalized sum or a centered random variable.**

**4. The Central Lyapunov Limit Theorem. If** $X_1, X_2, ..., X_n$ **- independent random variables having mathematical mean** $m_i$**, variance** $\sigma_i^2$ **and**

**finite absolute central moments of the third order** $|\mu_3|$, **satisfying the conditions**

$$\lim_{n \to \infty} \left( \sum_{i=1}^{n} |\mu_3(X_i)| \right) \Big/ \left( \sum_{i=1}^{n} \sigma_i^2 \right)^{3/2} = 0, \tag{7.9}$$

**then in unlimited increase** $n$ **distribution law n normalized sum (7.8) coincides probability to the normal law with parameters** $a = 0$, $\sigma = 1$.

Note that $\mu_3(X_i) = M[|X_i - m|^3]$ is called the central third-order moment random variable $X_i$.

**Contents of condition (7.9) is that the variance of each random variable** $X_i, i = \overline{1, n}$, **is only a small part of total variance sum** $\sum\limits_{i=1}^{n} X_i$.

**Outcome Lyapunov theorem. If the conditions Lyapunov theorem are performed, then the random variable** $\sum\limits_{i=1}^{n} X_i$ **for large** $n$ **with sufficient accuracy is normally distributed with parameters**

$$a = \sum_{i=1}^{n} m_i, \quad \sigma^2 = \sum_{i=1}^{n} \sigma_i^2.$$

In the practical tasks of the central limit theorem of Lyapunov ofen used to calculate the probability that the sum of several random variables attains a value that belongs to the listed interval. The reason of this is the use of so-called **Partial result Lyapunov theorem**. If $X_1$, $X_2$, ..., $X_n$ - independent random variables, which have equal mathematical expectation $M(X_i)=m$, the dispersion $D(X_i)= \sigma^2$, and absolute central moments of the third order $M[|X_i - m|^3] = \mu_3 \quad i = \overline{1, n}$, then the law of the distribution sum $Y_n=X_1+X_2+...+ X_n$ if $n \to \infty$ infinitely close to normal law with parameters $a = \sum\limits_{i=1}^{n} m$, $\sigma^2 = \sum\limits_{i=1}^{n} \sigma_i^2$ .

# PART TWO
# Mathematical Statistics

§ 1. Introduction In mathematical statistics.
sampling
1. Objectives of mathematical statistics.
2. General and the sample. Ways formation of the sample.
3. The statistical distribution of the sample.
4. Empirical distribution function and its properties.
5. Graphic representation of statistical distributions (polygon and histogram).
6. Numerical characteristics of the sample.


1. The first task of mathematical statistics - specify the method of boron-grouping and statistics, as well as the number of non-bypass test (statistics). The second problem of mathematical statistics is to develop methods statistical analysis according to the research objectives. One of them - the estimation of the unknown:
- The probability of a random event;
- Probability distribution function (density distribution);
- Distribution parameters, the type of which is known;
- Depending random variable from one or more instances Closed variables.
The second goal - a statistical test hypotheses about the unknown-type my size distribution or the distribution parameters, which kind is known.

1.  General called the totality of the same objects subject to study. The set of objects denoted through $\Omega$.

    Objects set $\Omega$ can characterized by one or few grounds. These symptoms may be quantitative and qualitative.

    Sometimes conduct a continuous survey that examined each item set $\Omega$ relative signs that investigation. However, in practice a continuous survey used relatively rarely. This is because when testing facility partially or completely destroyed, or research facilities require large investments. Sometimes conduct a complete inspection is physically impossible. In such cases, of course away from the general population, not limited number of objects to explore their quantitative or qualitative features, and then draw conclusions about the entire general population.

    The sample is called the set randomly selected objects from the general population. The sample forms a subset $V$ of the set $\Omega$ ($V \subset \Omega$).

Volume population (general or selective) is the number of objects together. Volume of the general population is denoted letter N, and sampling - *n*.

For opinions on the legality of the object under study the population under study sample is needed to properly sample objects representing the general population, that is, the sample must have property representativeness. The randomness of selection of objects in the sample and use the law of large numbers can solve the question of the representativeness of the sample.

The accuracy of a random observation, in the final sub-bag will depend on the method of selection of objects, the degree of fluctuation studied traits in the general population and the volume of the sample.

In practice, using different methods of forming the sample are essentially divided into two types:

1) selection that does not require a partition the population into parts (simple (actually random) selection);

2) selection, in which the population is divided into parts (typical selection, mechanical selection, serial selection, the combinations selection).

A simple random is a selection, in which objects are selected one at random from all the population. Simple random sampling can be repeated or unrepeated. Repeated called sampling, the formation of which selected object (before the next selection) returns in the general population. Unrepeated called sampling in the formation of any selected object in the general population is not refundable.

Let investigated quantitative trait $X$ objects general set $\Omega$. To reduce suppose that it is a one-dimensional random variable. After processing the sample *n* objects are received *n* numbers $x_1, x_2, \ldots, x_n$, called variants and form a series of variants or a simple statistical series.

Initial processing of row option is to make something equal variants this series. A number of variant arrange in ascending order and responsible this renumber them. As a result, we obtain a sequence of numbers $x_1 \leq x_2 \leq x_3 \leq \ldots \leq x_n$, called variational row. If some of this sequence are the same options, we have renumber, keeping the same number same options.

Let $x_1$ version number of variations is repeated $n_1$ times, $x_2$ - $n_2$ times, ..., $x_k$ - $n_k$ times. Numbers $n_i$ are called **frequency** (**absolute frequency**) and their relation to sample size $n_i/n = w_i$ - **relative frequencies**. From these definitions derive the equation:

$$\sum_{i=1}^{k} n_i = n, \quad \sum_{i=1}^{k} w_i = 1. \tag{1.1}$$

3. The statistical distribution of the sample is called the correspondence between options and frequencies or relative frequencies:

$$
\begin{array}{c|cccc}
x_i & x_1 & x_2 & \ldots & x_k \\
\hline
n_i & n_1 & n_2 & \ldots & n_k
\end{array}, \tag{1.2}
$$

$$
\begin{array}{c|cccc}
x_i & x_1 & x_2 & \ldots & x_k \\
\hline
w_i = {n_i}/{n} & w_1 & w_2 & \ldots & w_k
\end{array}. \tag{1.3}
$$

Most of the statistical distribution of the sample in the form (1.2) or (1.3) is used when a number of variant is an implementation of a **discrete** random variable *X*. If *X* is a continuous random variable, then the statistical distribution of the sample set is correspondence between intervals and frequencies or relative  frequencies of options that fall into these intervals, i.e. in the form of tabel:

$$
\begin{array}{c|cccc}
[x_i,x_{i+1}) & [x_1,x_2) & [x_2,x_3) & \ldots & [x_k,x_{k+1}] \\
\hline
n_i & n_1 & n_2 & \ldots & n_k
\end{array}, \tag{1.4}
$$

$$
\begin{array}{c|cccc}
[x_i,x_{i+1}) & [x_1,x_2) & [x_2,x_3) & \ldots & [x_k,x_{k+1}] \\
\hline
w_i = {n_i}/{n} & w_1 & w_2 & \ldots & w_k
\end{array}. \tag{1.5}
$$

These tables are called interval statistical distribution of sample. In constructing the interval statistical distribution based on a number option is considered *k* intervals of equal length. Number intervals can be determined approximately by the formula Sterdzhesa: $k = 1+3,322 \lg n$. In this case $k \in [5;6]$ if $20 \leq n \leq 30; k \in [7;8]$ if $60 \leq n \leq 70;$ $k \in [8;9]$ if $70 \leq n \leq 200;$ $k \in [9;15]$ if $n>200$.

Statistical distributions of the sample (1.3) and (1.5) are called empirical (research) distribution of the random variable X (quantitative-term features of objects of the general population).

One of the problems of mathematical statistics - grade (approximate location) unknown distribution function *F (x)* probabilities quantitative trait *X* objects the population. By the definition of *F (x) = P (X <x)*. In it the researcher statistics are grouped in the statistical distribution of frequencies or relative frequencies. Therefore, considering the relative frequency stability properties, it is advisable probability events *(X <x)* to approximate the relative frequency of the same event.

4. **The empirical distribution function** (**distribution function of the sample**) is a function $F^{*}(x)$ determined  argument *x,* that equals the relative

frequency of an event $(X < x)$ for the sample values of the random variable $X$, i.e.

$$F^*(x) = W(X < x) = n_x/n, \qquad (1.6)$$

where $n_x$ - the sum of the frequencies of the variant are smaller than $x$,
$n$ - sample size.

In contrast to $F^*(x)$ the function $F(x)$ is called the theoretical distribution function.

From the definition of empirical functions implies the following properties of similar properties theoretical distribution function:

1) $D(F^*) = R$, $E(F^*) = [0, 1]$;
2) $F^*(x)$ — nondecreasing function;
3) if $x_1$ — the smallest version, then $F^*(x) = 0$ for $x \le x_1$; if $x_k$ — the largest option, then $F^*(x) = 1$ for $x > x_k$.

5. In the analysis of statistical data essential role played by geo-metric illustration of this data. For clarity, build various graphs of statistical distributions, including polygon and histogram.

Let statistical distribution of the sample is determined by tables (1.2) or (1.3).

Polygon frequency (frequency polygon) is broken, straight segments which connect adjacent points $(x_1, n_1)$, $(x_2, n_2)$, …, $(x_k, n_k)$. For construction site on the x-axis lay the options $x_i$, and the vertical axis - the corresponding frequency $n_j$.

Polygon relative frequencies called broken, straight segments which connect adjacent points $(x_1, w_1)$, $(x_2, w_2)$, …, $(x_k, w_k)$. When building polygon relative frequency on the horizontal axis lay the options $x_i$, and the vertical axis - corresponding relative frequencies $w_i$.

In the case of continuous attributes is better to build a histogram. Let distributions (1.4) and (1.5) such that the length of each partial interval is the same number $h$.

Histogram of frequencies called the stepped shape, consisting of rectangles, which are based on partial interval length $h$, and the height is equal to $n_i/h$ (**density frequency**). For the structure in the frequency histogram on the horizontal axis deposited partial interval, and above them are held straight segments parallel to the x-axis the distance $n_i/h$. Area $i$ partial rectangle is $hn_i/h = n_i$, i.e. the sum of the frequencies of the option, that enter the $i$ interval. Therefore, the frequency histogram area is the sum of all frequencies of the sample $n$.

Comparison of two histograms suggests that the severity of the histogram depends essentially on the election of length h-parts are intervals.

**Histogram relative frequencies** called the stepped shape consisting of rectangles, which are based on partial interval length h, and the height is equal to $w_i/h$ (**density relative frequency**). To construct a histogram of relative frequency on the horizontal axis deposited partial intervals, and above them are held segments parallel to the x-axis at a distance $w_i/h$. Area $i$ partial rectangle is $hw_i/h = w_i$, i.e. the relative frequency of the option, that hit in the $i$ interval. Thus, the area of the histogram of relative frequency equal to one.

6. Construction of the statistical distribution of the sample (1.2) or (1.4) and their graphical representation - is only the first step towards solving the problems of mathematical statistics. The next step of the pre-foresees numerical characteristics that express in compact form the most significant features of the statistical distribution of the sample and serve estimates (approximate value) of unknown options quantitative trait distribution of the general                                                                                      population.
**The average sample (mean variant)** statistical distribution (1.2) is determined by the formula

$$\bar{x}_\text{в} = \left( \sum_{i=1}^{k} x_i n_i \right) \Big/ n . \qquad (1.7)$$

If all $n$ different option, then (1.7) is gaining like this:

$$\bar{x}_\text{в} = \left( \sum_{i=1}^{n} x_i \right) \Big/ n . \qquad (1.7^*)$$

If the sample is given by the statistical distribution of the frequency interval (1.4), whereas in finding the need to go to the discrete distribution (1.2), the "new" versions which are midpoints of intervals, and then use the formula (1.7).

In the said medium, the statistics used also **structural medium** that does not depend on the values of option located at the edges of the distribution, and associated with a number of frequencies. The structural medium         are         median         and         mode.
**The median** $Me^*$ discrete statistical distribution of the sample (1.2) is a number that divides variational series, that "creates" this distribution into two equal parts by the number of option. If the number is odd variant, i.e. $n = 2m$, then $Me^* = x_{m+1}$. If the sample size is an even number, i.e. $n = 2m$, then the median is the arithmetic mean of the "average" (median) pairs option:

$$Me^* = (x_m + x_{m+1})/2.$$

**The median** for interval statistical distribution is called a number $Me^*$, for which the equality:

$$F^*(Me^*) = 0,5, \qquad (1.8)$$

where $F^*(x)$ - function of the empirical distribution. The formula for calculating the median has the following form:

$$Me^* = x_m + \frac{0,5 - F^*(x_m)}{F^*(x_{m+1}) - F^*(x_m)}(x_{m+1} - x_m), \qquad (1.9)$$

where $[x_m, x_{m+1})$ - the so-called partial median interval $(1 \le m \le k)$ for which inequality $F^*(x_m) < 0,5, \ F^*(x_{m+1}) > 0,5$ performed. The **mode** $Mo^*$ of discrete statistical distribution (1.2) is named variant, which corresponds to the highest frequency. The mode for interval statistical distribution is calculated as follows. First you define modal interval $[x_m, x_{m+1})$, i.e. an interval for which $n_m/h_m = \max\limits_{1 \le i \le k}\{n_i/h_i\}$, where $h_i$ - length of partial interval $[x_m, x_{m+1})$, $n_i$ - number of variants from this interval. The value $Mo^*$ belongs the modal interval and is calculated by interpolation formula

$$Mo^* = x_m + \frac{n_m - n_{m-1}}{2n_m - n_{m-1} - n_{m+1}}h_m. \qquad (1.10)$$

Consider some numerical characteristics of the scattering option to vkola-selective medium.

**Sampling variance** statistical distribution (1.2) is called arithmetic mean of squared deviations of variants from the mean sample:

$$D_в = \frac{1}{n}\sum_{i=1}^{k}(x_i - \bar{x}_в)^2 n_i. \qquad (1.11)$$

$D_в$ characterizes the average value of the option spread around $\bar{x}_в$ in square units.

In practice, more convenient to use the so-called **calculating formula for calculating the variance:**

$$D_в = \overline{x^2} - (\bar{x}_в)^2 = \frac{1}{n}\sum_{i=1}^{k}x_i^2 n_i - (\bar{x}_в)^2. \qquad (1.12)$$

Disadvantage of $D_в$ is its dimension. To remedy this shortcoming uses different numerical description:
**standard deviation of the sample**

38

$$\sigma_{_в} = \sqrt{D_{_в}}. \qquad (1.13)$$

The spread relatively expectation separate option values characterize performance variation. The simplest of these is the **scale of variation rate** $R$, which is equal to the difference between the largest and smallest distribution options: $R = x_{\max} - x_{\min}.$ The range of variation is used in the statistical study of quality products.

If $\bar{x}_{_в}$ nonzero, then to compare two statistical distributions in terms of their dimensions relative to the intermediate sample is reproduced coefficient of variation, which is the ratio of standard deviation to the average sample and expressed as a percentage:

$$V = \frac{\sigma_{_в}}{\bar{x}_{_в}} \cdot 100\%. \qquad (1.14)$$

Selectively share (proportion) is the ratio of the number $m$ objects of sampling with feature $\alpha$ to sample size: $w = m/n.$ Feature $\alpha$ you can be product standardization, product grade, gender rights, etc.. In content, $w$ is the relative frequency of random events, which is randomly selected as the object of the population has a feature $\alpha$

## SOLVING COMMON PROBLEMS

Task 1.1. In studying the question of production quotas weavers observed the next frequency yarn breakages similar loom stalls at various intervals of equal length $t$:

7, 1, 2, 1, 2, 5, 4, 4, 3, 2, 2, 6, 0, 1, 6,
5, 3, 2, 0, 1, 4, 3, 2, 1, 5, 3, 0, 4, 2, 3.

1) To make the statistical distribution of frequencies and relative frequencies among the cliffs yarn stalls;
2) build a landfill frequencies and relative frequencies;
3) find empirical distribution function and build its schedule;
4) Calculate the sample, average, variance, standard deviation, mode, median, range of variation, coefficient of variation.

1) The sample size n = 30. This version number is written as a series of variations:

0, 0, 0, 1, 1, 1, 1, 1, 2, 2, 2, 2, 2, 2, 2, 3, 3, 3, 3, 3, 4, 4, 4, 4, 5, 5, 5, 6, 6, 7.

Variants: $x_1 = 0$, $x_2 = 1$, $x_3 = 2$, $x_4 = 3$, $x_5 = 4$, $x_6 = 5$, $x_7 = 6$, $x_8 = 7$.
Frequency: $n_1 = 3$, $n_2 = 5$, $n_3 = 7$, $n_4 = 5$, $n_5 = 4$, $n_6 = 3$, $n_7 = 2$, $n_8 = 1$.
As a result, we obtain the statistical distribution of frequencies:

| $x_i$ | 0 | 1 | 2 | 3 | 4 | 5 | 6 | 7 |
|-------|---|---|---|---|---|---|---|---|
| $n_i$ | 3 | 5 | 7 | 5 | 4 | 3 | 2 | 1 |

Checking: $\sum n_i = 3 + 5 + 7 + 5 + 4 + 3 + 2 + 1 = 30 = n$.

According to the formula $w_i = {}^{n_i}\!/_n$ consistently calculate the relative frequencies: $w_1 = 3/30$, $w_2 = 5/30$, $w_3 = 7/30$, $w_4 = 5/30$, $w_5 = 4/30$, $w_6 = 3/30$, $w_7 = 2/30$, $w_8 = 1/30$. Control: $\sum w_i = 3/30 + 5/30 + 7/30 + {} + 5/30 + 4/30 + 3/30 + 2/30 + 1/30 = 1$. Thus, the statistical distribution of relative frequencies among the cliffs yarn has the stalls type:

| $x_i$ | 0 | 1 | 2 | 3 | 4 | 5 | 6 | 7 |
|-------|------|------|------|------|------|------|------|------|
| $w_i$ | 3/30 | 5/30 | 7/30 | 5/30 | 4/30 | 3/30 | 2/30 | 1/30 |

2) Polygon of frequencies depicted in Fig.1.1 and polygon of relative frequencies — Fig.1.2.



Fig.1.1.



Fig.1.2.

3) The sample size $n = 30$. If $x \le 0$, then there is no options, lesser $x$, $n_x = 0$, and therefore $F^*(x) = n_x/30 = 0$.

Let $x \in (0; 1]$. Then variant $x = 0$ is less than $x$, so $n_x = 3$ and $F^*(x) = 3/30 = 0,1$.

If $x$ satisfies the double inequality $1 < x \le 2$, while less than $x$ are variants 0 and 1, the amount of which frequencies $n_x = 3 + 5 = 8$. So $F^*(x) = 8/30 = 4/15$ for $x \in (1; 2]$.

If $x$ is performed by the double inequality $2 < x \le 3$, then variants 0, 1, 2 are lesser $x$, sum frequency which $n_x = 3 + 5 + 7 = 15$. Therefore, for $x \in (2; 3]$ $F^*(x) = 15/30 = 0,5$.

Similarly, we find the value of $F^*(x)$ for intervals (3; 4], (4; 5], (5; 6], (6; 7], (7; ∞]. As a result, we obtain the required empirical distribution function:

$$F^*(x) = \begin{cases} 0, & \text{якщо } x \le 0, \\ 1/10, & \text{якщо } 0 < x \le 1, \\ 4/15, & \text{якщо } 1 < x \le 2, \\ 1/2, & \text{якщо } 2 < x \le 3, \\ 2/3, & \text{якщо } 3 < x \le 4, \\ 4/5, & \text{якщо } 4 < x \le 5, \\ 9/10, & \text{якщо } 5 < x \le 6, \\ 29/30, & \text{якщо } 6 < x \le 7, \\ 1, & \text{якщо } x > 7. \end{cases}$$

Graph this function is shown in Fig. 1.3.

Рис. 1.3.

2) To calculate $\bar{x}_в$ and $D_в$ use the formulas (1.7) і (1.12):

$$\bar{x}_в = \frac{\sum_{i=1}^{k} x_i n_i}{n} = \frac{0 \cdot 3 + 1 \cdot 5 + 2 \cdot 7 + 3 \cdot 5 + 4 \cdot 4 + 5 \cdot 3 + 6 \cdot 2 + 7 \cdot 1}{30} = 84/30 = 2,8;$$

$$D_в = \overline{x^2} - (\bar{x}_в)^2 = \frac{\sum_{i=1}^{k} x_i^2 n_i}{n} - (\bar{x}_в)^2 =$$

$$= \frac{0^2 \cdot 3 + 1^2 \cdot 5 + 2^2 \cdot 7 + 3^2 \cdot 5 + 4^2 \cdot 4 + 5^2 \cdot 3 + 6^2 \cdot 2 + 7^2 \cdot 1}{30} - (2,8)^2 =$$

$$= 338/30 - 7,84 = 3,4267;$$

$$\sigma_в = \sqrt{D_в} = \sqrt{3,4267} = 1,8511.$$

Conclusion: The average number of yarn breakages stalls during the time interval is 2.8, and the average spread of numbers breaks stalls average of 2.8 is 1.8511.

Mode $Mo^* = 2$, since variant 2 corresponds to the highest frequency of 7.

Median $Me^*$ can be find, using variation row obtained in 1), or directly from the statistical distribution. Sum frequency first three variants of this distribution is 15 (half-volume of sample) and the next five - also 15. So the median is between variants 2 and 3: $Me^* = (2+3)/2 = 2,5$. A mismatch $\bar{x}_в$, $Mo^*$ and $Me^*$ indicates the absence of strict symmetry of the distribution.

The coefficient of variation (according to the formula (1.14))

$$V = \frac{\sigma_в}{\bar{x}_в} \cdot 100\% = \frac{1,8511}{2,8} \cdot 100\% = 66,11\%.$$

Task 1.2. During the period between regular readjustment equipment made control measurements of thickness (in mm) 200 liners connecting rod bearings. The data are presented in Table. 1.2.
    1) To make the interval statistical distribution of frequencies and relative frequencies of the sample. Based on the statistical distribution of the resulting interval:
    2) build a histogram of frequencies and relative frequencies;

3) find empirical distribution function and build its schedule;
4) to calculate $\bar{x}_в$, $D_в$, $\sigma_в$, the median and mode.

| 1,754 | 1,739 | 1,743 | 1,764 | 1,733 | 1,736 | 1,743 | 1,742 |
|-------|-------|-------|-------|-------|-------|-------|-------|
| 1,728 | 1,732 | 1,731 | 1,752 | 1,737 | 1,747 | 1,758 | 1,737 |
| 1,724 | 1,737 | 1,733 | 1,713 | 1,740 | 1,740 | 1,729 | 1,740 |
| 1,751 | 1,739 | 1,747 | 1,748 | 1,730 | 1,750 | 1,740 | 1,732 |
| 1,740 | 1,730 | 1,748 | 1,757 | 1,741 | 1,733 | 1,743 | 1,745 |
| 1,748 | 1,723 | 1,737 | 1,748 | 1,741 | 1,751 | 1,714 | 1,750 |
| 1,744 | 1,748 | 1,758 | 1,756 | 1,727 | 1,731 | 1,738 | 1,753 |
| 1,735 | 1,738 | 1,743 | 1,729 | 1,743 | 1,737 | 1,731 | 1,734 |
| 1,741 | 1,742 | 1,744 | 1,756 | 1,744 | 1,752 | 1,739 | 1,740 |
| 1,729 | 1,745 | 1,742 | 1,753 | 1,743 | 1,734 | 1,731 | 1,734 |
| 1,732 | 1,732 | 1,746 | 1,748 | 1,755 | 1,738 | 1,742 | 1,729 |
| 1,731 | 1,725 | 1,729 | 1,745 | 1,739 | 1,754 | 1,752 | 1,720 |
| 1,750 | 1,734 | 1,749 | 1,738 | 1,747 | 1,757 | 1,751 | 1,746 |
| 1,723 | 1,736 | 1,746 | 1,744 | 1,759 | 1,728 | 1,751 | 1,750 |
| 1,746 | 1,759 | 1,748 | 1,740 | 1,735 | 1,745 | 1,740 | 1,746 |
| 1,737 | 1,726 | 1,743 | 1,755 | 1,740 | 1,726 | 1,745 | 1,744 |
| 1,735 | 1,746 | 1,739 | 1,732 | 1,758 | 1,744 | 1,754 | 1,724 |
| 1,742 | 1,750 | 1,761 | 1,758 | 1,753 | 1,757 | 1,720 | 1,733 |
| 1,738 | 1,728 | 1,758 | 1,732 | 1,763 | 1,733 | 1,745 | 1,766 |
| 1,745 | 1,743 | 1,734 | 1,733 | 1,755 | 1,756 | 1,769 | 1,750 |
| 1,740 | 1,762 | 1,738 | 1,742 | 1,740 | 1,740 | 1,760 | 1,752 |
| 1,746 | 1,728 | 1,743 | 1,718 | 1,738 | 1,762 | 1,728 | 1,734 |
| 1,753 | 1,751 | 1,748 | 1,735 | 1,739 | 1,729 | 1,754 | 1,736 |
| 1,762 | 1,748 | 1,738 | 1,726 | 1,757 | 1,738 | 1,726 | 1,720 |
| 1,751 | 1,734 | 1,724 | 1,741 | 1,752 | 1,732 | 1,738 | 1,739 |

○ 1). All variants are in the range [1,713; 1.769] length of 0,056 mm. We divide it into 8 equal the length of intervals:

[1,713; 1,720) [1,720; 1.727) [1.727; 1,734) [1,734; 1.741)
[1.741; 1.748) [1.748; 1.755) [1.755; 1.762) [1.762; 1.769].

Each version of the table. 1.2 assign one of these partial intervals and find sum of option, enter the first, second, ..., eighth intervals. The result:

$n_1 = 3$,  $n_2 = 13$,  $n_3 = 32$,  $n_4 = 49$,  $n_5 = 42$,  $n_6 = 35$,  $n_7 = 19$,  $n_8 = 7$,

$$n = \sum_{i=1}^{8} n_i = 200 .$$

The result is a statistical interval frequency distribution is shown in Table 1.3.

*Table 1.3*

| $[x_i; x_{i+1})$ | $n_i$ | $[x_i; x_{i+1})$ | $n_i$ |
|---|---|---|---|
| [1,713; 1,720) | 3 | [1,741; 1,748) | 42 |
| [1,720; 1,727) | 13 | [1,748; 1,755) | 35 |
| [1,727; 1,734) | 32 | [1,755; 1,762) | 19 |
| [1,734; 1,741) | 49 | [1,762; 1,769] | 7 |

In view of the fact that the sample size $n = 200$, we write the inter-abundant statistical distribution of relative frequencies of the sample (Table. 1.4):

*Table 1.4*

| $[x_i; x_{i+1})$ | $w_i$ | $[x_i; x_{i+1})$ | $w_i$ |
|---|---|---|---|
| [1,713; 1,720) | 3/200 | [1,741; 1,748) | 42/200 |
| [1,720; 1,727) | 13/200 | [1,748; 1,755) | 35/200 |
| [1,727; 1,734) | 32/200 | [1,755; 1,762) | 19/200 |
| [1,734; 1,741) | 49/200 | [1,762; 1,769] | 7/200 |

2). The histogram of the frequency distribution $h = 0,007$ depicted in rys.1.4 and relative frequency histogram for rys.1.5.



Rys.1.4.

Rys.1.5.

3). The study quantitative trait $X$ is a continuous random variable. Therefore, the probability distribution function $F(x)$, and empirical function $F^*(x)$ are continuous functions determined arhument x. Let $x \leq 1,713$. Then $n_x = 0$, as observed values quantitative features smaller than x, are absent. Consequently, $F^*(x) = 0$ for all $x \leq 1,713$.

Find the value of the empirical distribution function for $x = 1,720$ - left end of the second interval. Thus, we assume that we can not do this for each internal point of the first interval. When $x = 1,720$ $n_x = 3$ i $F^*(1,720) = 3/200 = 0,015$.

For $x = 1,727$ $n_x = 3 + 13 = 16$, $F^*(1,727) = 16/200 = 0,08$.
For $x = 1,734$ $n_x = 16 + 32 = 48$, $F^*(1,734) = 48/200 = 0,24$.
Similarly, we find:
$F^*(1,741) = (48 + 49)/200 = 97/200 = 0,485$;
$F^*(1,748) = (97 + 42)/200 = 139/200 = 0,695$;
$F^*(1,755) = (139 + 35)/200 = 174/200 = 0,87$;
$F^*(1,762) = (174 + 19)/200 = 193/200 = 0,965$.
Finally, for $x > 1,769$ all 200 observed quantitative trait values smaller than x, i.e. $n_x = 200$.. So $F^*(x) = 1$ for $x > 1,769$.
Construct a graph of the empirical distribution function: first for intervals $(-\infty; 1,713]$ and $(1,769; \infty)$, then in these locations. In order to show the

46

continuity of change $F^*(x)$, obtained neighboring points connect straight segments (Fig. 1.6).



Fig. 1.6.

4). To find the numerical characteristics of the sample interval specified statistical distribution of claim 1, move on to the discrete distribution:

| $x_i$ | 1,7165 | 1,7235 | 1,7305 | 1,7375 | 1,7445 | 1,7515 | 1,7585 | 1,7655 |
|-------|--------|--------|--------|--------|--------|--------|--------|--------|
| $n_i$ | 3 | 13 | 32 | 49 | 42 | 35 | 19 | 7 |

This "new" options are midpoints partial intervals.

$$\overline{x}_в = \frac{\sum_{i=1}^{k} x_i n_i}{n} =$$

$$= \frac{1,7165 \cdot 3 + 1,7235 \cdot 13 + 1,7305 \cdot 32 + 1,7375 \cdot 49 + 1,7445 \cdot 42 + 1,7515 \cdot 35 + 1,7585 \cdot 19 + 1,7655 \cdot 7}{200} =$$

$$= 1,7417;$$

$$D_в = \overline{x^2} - (\overline{x}_в)^2 = \frac{\sum_{i=1}^{k} x_i^2 n_i}{n} - (\overline{x}_в)^2 =$$

47

$$= \frac{(1,7165)^2 \cdot 3 + (1,7235)^2 \cdot 13 + (1,7305)^2 \cdot 32 + (1,7375)^2 \cdot 49}{200} +$$

$$+ \frac{(1,7445)^2 \cdot 42 + (1,7515)^2 \cdot 35 + (1,7585)^2 \cdot 19 + (1,7655)^2 \cdot 7}{200} - (1,7417)^2 =$$

$$= 0,0001283;$$

$$\sigma_{_в} = \sqrt{D_{_в}} = \sqrt{0,0001283} = 0,0113269.$$

Since all partial distribution of the intervals have the same length, it is a modal interval [1,734; 1.741), which corresponds to the highest frequency 49. The value $Mo^*$ contained within this interval is calculated by the formula (1.10).

$$Mo^* = x_m + \frac{n_m - n_{m-1}}{2n_m - n_{m-1} - n_{m+1}} h = 1,734 + \frac{49 - 32}{2 \cdot 49 - 32 - 42} \cdot 0,007 =$$

$$= 1,734 + \frac{17}{24} \cdot 0,007 \approx 1,73896.$$

To calculate the median find the first partial median interval $[x_m; x_{m+1})$, that performs inequality:

$$F^*(x_m) < 0,5, \qquad F^*(x_{m+1}) > 0,5.$$

According to 3) $F^*(1,741){=}0,485{<}0,5$, $F^*(1,748){=}0,695{>}0,5$.

Therefore, $x_m = 1,741$, $x_{m+1} = 1,748$. According to the formula (1.9)

$$Mo^* = x_m + \frac{0,5 - F^*(x_m)}{F^*(x_{m+1}) - F^*(x_m)} (x_{m+1} - x_m) =$$

$$= 1,741 + \frac{0,5 - 0,485}{0,695 - 0,485} (1,748 - 1,741) =$$

$$= 1,741 + \frac{0,015}{0,21} \cdot 0,007 = 1,741 + 0,0005 = 1,7415.$$

The coefficient of variation calculated using the formula (1.14):

$$V = \frac{\sigma_{\hat{a}}}{\overline{x}_{\hat{a}}} \cdot 100\% = \frac{0,0113269}{1,7417} \cdot 100\% = 0,65\%.$$

# § 2. STATISTICAL EVALUATION

*1. Point statistical estimation of distribution parameters and their properties.*
*2. Assessment of the general average for easy sample (repeated and unrepeated).*
*3. Evaluation of the general proportion for easy sample (repeated and*

*unrepeated).*

*4. The mean square error (MSE) simple sample. Correct the county-sample variance.*

*5. Interval statistical estimation (Confidence intervals for $\bar{x}_\text{г}$ and was non-small samples.*

*6. Finding the minimum sample size.*

*7. Confidence intervals for the estimates $\bar{x}_\text{г} = a$ for small sample.*

*Confidence intervals for and in the case of small sample).*


## 1. Spot statistical estimation of distribution parameters and their properties.

Let investigated continuous quantitative trait $X$ objects of general set in order to find the unknown distribution law. In it the researchers are sample statistics. Suppose that for whatever reasons hypothesized normal distribution features of $X$ (research question of the correctness of this hypothesis or falsity will be conducted in the next section). Since the normal distribution is completely determined by two parameters $a$ and $\sigma$, there is a need to evaluate them is to find approximate value using only the observed variations $x_1, x_2, ..., x_n$ sample of n. Setting $a = M(X)$, the expectation characterizes the arithmetic mean of observed possible values of the random variable $X$. On the other hand $\bar{x}_\text{в}$ - is the arithmetic mean of the option.

Therefore (yet intuitive) should approximate $a$ as mean sample:

$$a \approx \bar{x}_\text{в} = \left( \sum_{i=1}^{k} x_i n_i \right) \Big/ n.$$

Similarly,

$$\sigma \approx \sigma_\text{в} = \sqrt{ \sum_{i=1}^{k} (x_i - \bar{x}_e)^2 n_i \Big/ n }. \qquad (2.2)$$

Again note that the right parts of equations (2.1), (2.2) are **random variables** since objects in the sample **fall randomly**.

Now let quantitative trait $X$ objects the population is discrete random variable. Suppose that statistical distribution general statistics and sample tables are described:

$$\begin{array}{c|cccc} x_i & x_1 & x_2 & ... & x_m \\ \hline N_i & N_1 & N_2 & ... & N_m \end{array}, \qquad N = \sum_{i=1}^{m} N_i, \qquad (2.3)$$

$$\begin{array}{c|cccc} x_i & x_1 & x_2 & ... & x_m \\ \hline n_i & n_1 & n_2 & ... & n_m \end{array}, \qquad n = \sum_{i=1}^{m} n_i, \qquad (2.4)$$

where $N$ and $n$ - general volume and sample respectively. The distribution (2.3) is hypothetical - it will always be unknown to us, because otherwise no longer be a need to study the sample. Finally, we note that some of the frequency distribution (2.4) can be zero, which corresponds to a situation where the value of quantitative trait objects not met the population of variant sample (compare the distribution (2.4) to (1.2)).

By analogy with the numerical characteristics of the sample following formula determines the number of general characteristics general set:

$$\bar{x}_{\text{г}} = \left( \sum_{i=1}^{m} x_i N_i \right) \Big/ N, \tag{2.5}$$

$$D_{\text{г}} = \left( \sum_{i=1}^{m} (x_i - \bar{x}_{\text{г}})^2 N_i \right) \Big/ N, \tag{2.6}$$

$$\sigma_{\text{г}} = \sqrt{D_{\text{г}}}. \tag{2.7}$$

These numbers are unknown, and evaluation (approximation) they give the following equation:

$$\bar{x}_{\text{г}} \approx \bar{x}_{\text{в}}, \qquad D_{\text{г}} \approx D_{\text{в}}, \qquad \sigma_{\text{г}} \approx \sigma_{\text{в}}. \tag{2.8}$$

Examples can make some conclusions. Left of approximate equalities (2.1), (2.2) and (2.8) are unknown parameters "known" of the distribution or numerical characteristics of the general population; they are unknown numbers for a particular general population. The right side of these equations are functions of random variables, which are fixed to the sample statistics shall enter numeric values that can be represented by points. This allows you to call their statistical point estimates of the relevant parameters and numerical characteristics of the population. Denote the generalized symbol $\Theta$ left of approximate equality (2.1) (2.2) (2.8) (as well as many others that are available to other distribution laws) and symbol $\Theta^*$ - right parts of equations.

**Point statistical estimation, sampling function or statistics** numerical parameter $\Theta$ is a function of sample values (option $\Theta^* = \Theta^*(x_1, x_2, ..., x_n)$, which is in a statistical sense is close to the true value of this parameter. **Unbiased estimation** is called point statistical evaluation $\Theta^*$, the expectation it is equal estimated parameters $\Theta$ in an arbitrary sample size, i.e.

$$M(\Theta^*) = \Theta. \tag{2.9}$$

**Shifted estimation** is called evaluation, which does not the equality (2.9). Suppose that for estimation parameter $\Theta$ can be used $\Theta$ unbiased point

estimates $\Theta_1^*, \Theta_2^*, ..., \Theta_k^*$. Estimation $\Theta_m^*$, $1 \le m \le k$, is called **effective** if a given sample volume $n$ her the equality

$$D(\Theta_m^*) = \min_{1 \le i \le k} D(\Theta_i^*).$$

Assessment $\Theta^*$ is called **able to score** estimation parameter $\Theta$, if $n \to \infty$ it coincides with the probability in to $\Theta$, i.e. for arbitrarily small $\varepsilon > 0$ there is a limit switch

$$\lim_{n \to \infty} P\left( \left| \Theta - \Theta^* \right| < \varepsilon \right) = 1.$$

2. **Assessment of the general average for easy selection. Theorem 2.1. To repeating sample with volume $n$ mean is unbiased and able to estimate an unknown general average $\bar{x}_\text{г}$. If n is large enough, then $\bar{x}_\text{в}$ with sufficient accuracy normally distributed with parameters:**

$$a = \bar{x}_\text{г,} \qquad \sigma = \sqrt{D_\text{г}/n}. \qquad\qquad (2.10)$$

**Theorem 2.2. To unrepeating sample with volume $n$ mean $\bar{x}_\text{в}$ is unbiased estimate of an unknown general average $\bar{x}_\text{г}$. For sufficiently large $n$ $\bar{x}_\text{в}$ with sufficient accuracy normally distributed with parameters:**

$$a = \bar{x}_\text{г,} \qquad \sigma = \sqrt{\frac{D_\text{г}}{n} \cdot \frac{N-n}{N-1}}, \qquad\qquad (2.11)$$

where $N$ - the amount of the population.
*Note.* The amount of the population $N$, tend to be very large. Therefore, replacing the denominator of the second equality (2.11) $N-1$ on $N$ distinguishable to $\sigma$. In this regard will continue to enjoy equality

$$\sigma = \sqrt{\frac{D_\text{г}}{n} \cdot \left( 1 - \frac{n}{N} \right)}. \qquad\qquad (2.12)$$

3. **Assessment of the general proportion for easy selection.**
Let qualitative attributes of the objects general population whose number is finite are studied. **General share** will be called $M$ ratio of the population of objects, that have sign $\alpha$, to the amount of $N$ the population: $p = M/N$.

**Theorem 2.3.** To repeating sample with volume *n* sample average share *w* is unbiased and capable point estimate unknown general share. If *n* is large enough, then *w* with reasonable accuracy normally distributed with parameters:

$$a = p, \quad \sigma = \sqrt{pq/n}, \quad (2.13)$$

where $q = 1 - p.$

**Theorem 2.4.** For unrepeating sample with volume *n* random fraction *w* is unbiased point estimate unknown general proportion *p*. If *n* is large enough, then *w* with reasonable accuracy normally distributed with parameters:

$$a = p, \quad \sigma = \sqrt{\frac{pq}{n} \cdot \frac{N-n}{N-1}}, \quad (2.14)$$

*Note*. Since the volume of the population *N* in many cases is very large, then replacement in the denominator second equality (2.14) $N - 1$ for *N* is distinguishable to $\sigma$. With this in mind will continue using equality

$$\sigma = \sqrt{\frac{pq}{n}\left(1 - \frac{n}{N}\right)}. \quad (2.15)$$

**4. The mean square error (MSE) simple sample. Fixed dispersion sample variance.**

Random variables $\bar{x}_{\text{B}}$ and *w* is an unbiased point estimates of unknown statistics and *p* respectively. Their implementation or possible values obtained on the basis of simple sampling (repeating and unrepeating) does not coincide with the estimated parameters. And each such mismatch or reject call natural bias estimates, due to the fact that not all studied general population, but only a part (sample set). For statistics is very important information about the average of these errors.

Mean square error (MSE) in estimating the unknown general $\bar{x}_{\text{r}}$ and general secondary particles was called standard deviation of the sample and the sample average share w respectively.

**Mean square error** (MSE) in estimating the unknown general $\bar{x}_{\text{r}}$ and general proportion *p* is cold the average standard deviation $\bar{x}_{\text{B}}$ and sample

share $w$ respectively.

In light of the results of Theorems 2.1-2.4 can specify the following formula to determine the mean square error:

Standard deviation of repeating sample(see. (2.10))

$$\overline{\sigma}_{\overline{x}} = \sqrt{D_{\text{г}}/n};\qquad(2.16)$$

Standard deviation of unrepeating sample(see. (2.12))

$$\overline{\sigma}'_{\overline{x}} = \sqrt{\frac{D_{\text{г}}}{n}\left(1 - \frac{n}{N}\right)};\qquad(2.17)$$

MSE sample particle re-sampling (see. (2.13))

$$\overline{\sigma}_{w} = \sqrt{pq/n};\qquad(2.18)$$

MSE sample particles of unrepeating sample (see. (2.15))

$$\overline{\sigma}'_{w} = \sqrt{\frac{pq}{n}\left(1 - \frac{n}{N}\right)}.\qquad(2.19)$$

In practice, using formulas (2.16) - (2.19) is not possible, because it is necessary to know the variance or general or general share (analyze how aggravating this condition is remembering the initial condition of the problem of assessment). This fundamental difficulty can be eliminated by replacing in the formulas general dispersion $D_{\text{г}}$ on dispersion sample $D_{\text{в}}$ and the general share $p$ was - at selective share $w$. However, if such replacement must first ascertain whether appears even more systematic error due to bias estimates $D_{\text{в}}$ in both estimation problems. If they were displaced, then the replacement should be introduced "fixed" estimation of the unknown variances general. Towards the implementation of the above approach is useful such statements.

**Theorem 2.5. Expected sample variance in the task on the evaluation of an unknown general average for repeating sample is determined by equality**

$$M(D_{\text{в}}) = \frac{n-1}{n}D_{\text{г}},\qquad(2.20)$$

**and for unrepeating sample** –

$$M(D_{\text{в}}) = \frac{n-1}{n}\cdot\frac{N}{N-1}D_{\text{г}},\qquad(2.21)$$

where $n$ and $N$ respectively volumes of sample and the population.

Theorem 2.6. **Expected sample variance in the task on the evaluation of an unknown general proportion average for repeating sample is determined by equality**

$$M(D_{\text{в}}) = \frac{n-1}{n} pq, \tag{2.22}$$

**and for unrepeating sample** –

$$M(D_{\text{в}}) = \frac{n-1}{n} \cdot \frac{N}{N-1} pq. \tag{2.23}$$

Contents of theorems 2.5 and 2.6 is that the variance is biased estimation general variance in task estimation unknown $\bar{x}_{\text{г}}$ and $p$ in the case of pure sample (repeating and unrepeating). This formula (2.20) - (2.23) indicate a general underestimation of the general variance values due to the presence in each of the multiplier $(n-1)/n$.

However, this bias is easy to "fix": enough sample variance multiplied by a fraction $n/(n-1)$.

**Fixed dispersion (improved дисперсією)** is called numerical characteristic $S^2$, which is defined by equality

$$S^2 = \frac{n}{n-1} D_{\text{в}}.$$

For non-small samples (n≥30), replacing in the formulas (2.16) - (2.19) dispersion general $D_г$ selective $D_в$, and the share of general $p$ - selective $w$, we get used to practice the formula:
Standard deviation of repeating sample

$$\bar{\sigma}_{\bar{x}} = \sqrt{D_{\text{в}}/n} = \sigma_{\text{в}}/\sqrt{n}; \tag{2.24}$$

Standard deviation of unrepeating sample

$$\bar{\sigma}'_{\bar{x}} = \sqrt{\frac{D_{\text{в}}}{n}\left(1 - \frac{n}{N}\right)}; \tag{2.25}$$

MSE sample particle re-sampling

$$\bar{\sigma}_w = \sqrt{\frac{w(1-w)}{n}}; \tag{2.26}$$

MSE sample particle nore-sampling

$$\bar{\sigma}'_w = \sqrt{\frac{w(1-w)}{n}\left(1 - \frac{n}{N}\right)}. \tag{2.27}$$

54

## 5. Interval statistical evaluation.

Point statistical evaluation $\Theta^*$ does not match with the real value of the unknown parameter $\Theta$. Therefore, there is always uncertainty when replacing the unknown parameters of the evaluation, i.e. $\left|\Theta - \Theta^*\right| > 0$.

Magnitude error while unknown, although you need to know which errors can result indicated above replacement.

**Interval** called statistical evaluation, which is defined by **two numbers** - the ends of the interval. The advantage of interval estimates is that they allow you to set the accuracy and reliability of the estimates.

The number $\delta > 0$, that appears in the inequality

$$\left|\Theta - \Theta^*\right| < \delta, \tag{2.28}$$

naturally called **misstatements** $\Theta^*$. But talk about the inequality (2.28) can only in probabilistic sense because $\Theta^*$ - random value. That is, for sufficiently small $\delta$ inequality is random event.

**The reliability (confidence probability)** estimation $\Theta$ on $\Theta^*$ is called on probability γ performance inequality (2.28):

$$P\left(\left|\Theta - \Theta^*\right| < \delta\right) = \gamma. \tag{2.29}$$

**Confidence** is called the interval $(\Theta^* - \delta;\ \Theta^* + \delta)$, with a given reliability γ cover unknown parameter $\Theta$.

**Confidence intervals for the estimates $\bar{x}_r$ and $p$ for non-small samples**

The highest deviation average of sample (or sample proportion) from average general (or general share), which is possible for a given confidence probability γ, called **marginal error** Δ. Limiting error is given by

$$\Delta = t\bar{\sigma}, \tag{2.30}$$

where t - the root of the equation $2\Phi(t) = \gamma$.

Substituting in equation (2.30) expressions (2.24) - (2.27), we obtain suitable for practical formulas limit of error: average selective re-sampling

$$\Delta = t\sqrt{D_{\text{в}}/n}; \tag{2.31}$$

average selective nore-sampling

$$\Delta = t\sqrt{\frac{D_{\text{в}}}{n}\left(1 - \frac{n}{N}\right)}; \qquad (2.32)$$

proportion of selective re-sampling

$$\Delta = t\sqrt{w(1-w)/n}; \qquad (2.33)$$

proportion of selective nore-sampling

$$\Delta = t\sqrt{\frac{w(1-w)}{n}\left(1 - \frac{n}{N}\right)}. \qquad (2.34)$$

The result $(\bar{x}_{\text{в}} - \Delta; \ \bar{x}_{\text{в}} + \Delta)$ is a confidence interval, which covers with reliability $\gamma$ unknown average general $\bar{x}_{\text{г}}$. Similarly $(w - \Delta; \ w + \Delta)$ - confidence interval with the same reliability covers general unknown proportion $p$.

***Note***. In some cases it may be known numerical value $\sigma_{\text{г}}$ (such as information process). Then $D_{\text{в}}$ in formulas (2.31), (2.32) should be replaced by $\sigma_{\text{г}}^2$.

## 6. Finding the minimum sample size.

Before the formation of the sample to find out what should be its volume. The following formulas determine the minimum sample sizes in assessing unknown:

a) average general

$$n = \frac{t^2 D_{\text{г}}}{\Delta^2} \qquad (2.35)$$

to re-sample

$$n' = \frac{nN}{n + N} \qquad (2.36)$$

for nore-sample
($n$ determined by (2.35));

b) the general share

$$n = \frac{t^2 pq}{\Delta^2} \qquad (2.37)$$

to re-sample

$$n' = \frac{nN}{n + N} \qquad (2.38)$$

(2.38)
for no re-sample (*n* determined by (2.37));
where *N* - the amount of the population.

## 7. Confidence intervals for the estimates $\bar{x}_\text{г} = a$ for small sample

Let investigated on quantitative traits *X* knows just what she normally distributed. The problem:find confidence interval for the estimation of unknown parameters $a = M(X)$ (or $\bar{x}_\text{г}$ in the case of finite volume of the population) according to the re-sampling a small amount of *n*, given confidence probability γ and if unknown parameter $\sigma = \sigma_\text{г}$.

In [8] it is proved that the confidence interval for an unknown parameter *a* ( $\bar{x}_\text{г}$ ) if unknown $\sigma(\sigma_\text{г})$ with reliability γ has the following form:

$$x_\text{в} - t(\gamma, n-1)\frac{S}{\sqrt{n}} < a < x_\text{в} + t(\gamma, n-1)\frac{S}{\sqrt{n}}, \qquad (2.39)$$

where

$$S^2 = \sum_{i=1}^{n}(x_i - x_\text{в})^2 \Big/ (n-1), \qquad (2.40)$$

parameter $t=t(\gamma, n\text{-}1)$ - the root of the equation

$$P\left(\left|\frac{x_\text{в} - a}{S/\sqrt{n}}\right| < t\right) = P\big(|T| < t\big) = 2\int_0^t g_k(t)dt = \gamma, \qquad (2.41)$$

which can be found in Table 4, depending on the applications given confidential probability γ and number of degrees of freedom $k = n - 1$; $g_k(t)$ - density St'yudent distribution.

## Confidence intervals for $D_\text{г}$ and $\sigma_\text{г}$ in the case of small sample.

When finding the minimum sample size required is general information about the variance (standard deviation general). Often at the disposal of the researcher is only a small sample.

We can prove (see. [8]) that the confidence interval for estimation an unknown variance general $D_\text{г}=\sigma^2$ with reliability γ has the following form:

$$\frac{(n-1)S^2}{\chi_2^2} < \sigma^2 < \frac{(n-1)S^2}{\chi_1^2}, \qquad (2.42)$$

where $S$ is defined by equality (2.40), and values $\chi_1^2$ і $\chi_2^2$ are at Table 6 applications using equations

$$P\left(\chi^2(k) > \chi_1^2(p;k)\right) = p, \ \ p = \frac{1+\gamma}{2}; \qquad (2.43)$$

$$P\left(\chi^2(k) > \chi_2^2(p;k)\right) = p, \ \ p = \frac{1-\gamma}{2}; \qquad (2.44)$$

$k = n - 1$ - number of degrees of freedom of the distribution $\chi^2$.

With double inequality (2.42) is equivalent to double inequality (2.45)

$$\frac{S\sqrt{n-1}}{\chi_2} < \sigma < \frac{S\sqrt{n-1}}{\chi_1}, \qquad (2.45)$$

which defines the confidence interval for an unknown evaluation $\sigma(\sigma_r)$.

***Note.*** In the Table. 6 applications are values $\chi^2(p;k)$, that satisfy the equation $P\left(\chi^2(k) > \chi^2(p;k)\right) = p$ only for $k = n - 1 \le 30$, and for the $k = 40, 50, 100$. This due the next condition: when $k$ growth then the law of distribution of $\chi^2$ close to normal.

$$\frac{(n-1)S^2}{\chi_2^2} < \sigma^2 < \frac{(n-1)S^2}{\chi_1^2}, \qquad (2.42)$$

Із подвійної нерівності (2.42) отримаємо рівносильну подвійну нерівність

$$\frac{S\sqrt{n-1}}{\chi_2} < \sigma < \frac{S\sqrt{n-1}}{\chi_1}, \qquad (2.45)$$

яка визначає довірчий інтервал для оцінки невідомої $\sigma(\sigma_r)$.

***Зауваження***. В табл. 6 додатків наведені значення $\chi^2(p;k)$, що задовольняють рівняння $P\left(\chi^2(k) > \chi^2(p;k)\right) = p$ тільки для $k = n - 1 \le 30$, а також для $k = 40, 50, 100$. Це зумовлено тим, що при зростанні $k$ закон розподілу $\chi^2$ наближається до нормального.

SOLVING                    PROBLEMS                    COMMON
**Task 2.1.** The test results on the strength of the steel wires of equal length are shown in Table. 2.1 interval distribution.

58

| Breaking strength (kg / mm | [30; 32) | [32; 34) | [34; 36) | [36; 38) | [38; 40) | [40; 42) | [42; 44] |
|---|---|---|---|---|---|---|---|
| Number of wires | 6 | 12 | 17 | 29 | 19 | 13 | 4 |

1) Find a confidence interval that cover with reliability $\gamma = 0,95$ mean breaking the whole party 4000 wires.
2) Find a confidence probability that the sample mean depart trom general the average in absolute **value no more than 0.5 kg / mm2.**

○

**1)** Sample is nosmal ($n = 100 > 30$), sample is unrepeating, because after checking facility (wire) can not be restore a bearing in the general population (wire terminated due to checking its strength). To find the limits of the confidence interval limit error $\Delta$ find the formula (2.32), previously know sample numeric characteristics $\overline{x}_{\text{в}}$ and $D_{\text{в}}$.

$$\begin{array}{c|ccccccc} x_i = \dfrac{x_k + x_{k+1}}{2} & 31 & 33 & 35 & 37 & 39 & 41 & 43 \\ \hline n_i & 6 & 12 & 17 & 29 & 19 & 13 & 4 \end{array}$$

$$\overline{x}_{в} = \frac{\sum\limits_{i=1}^{k} x_i n_i}{n} =$$

$$= \frac{31 \cdot 6 + 33 \cdot 12 + 35 \cdot 17 + 37 \cdot 29 + 39 \cdot 19 + 41 \cdot 13 + 43 \cdot 4}{100} = 36,96;$$

$$D_{в} = \overline{x^2} - (\overline{x}_{в})^2 = \frac{\sum\limits_{i=1}^{k} x_i^2 n_i}{n} - (\overline{x}_{в})^2 =$$

$$= \frac{31^2 \cdot 6 + 33^2 \cdot 12 + 35^2 \cdot 17 + 37^2 \cdot 29 + 39^2 \cdot 19 + 41^2 \cdot 13 + 43^2 \cdot 4}{100} - (36,96)^2 =$$

$$= 9,0384$$

For the function tables of Laplace (tab. 3 applications) find the value of $t$ using the equation $\Phi(t) = \gamma/2 = 0{,}95/2 = 0{,}475$, $t = 1{,}96$. Formula (2.32), where $N = 4000$,

$$\Delta = 1{,}96\sqrt{\frac{9{,}0384}{100}\left(1 - \frac{100}{4000}\right)} \approx 0{,}5818,$$

then left limit of the confidence interval

$$\bar{x}_{\text{в}} - \Delta = 36{,}96 - 0{,}5818 = 36{,}3782,$$

права —

$$\bar{x}_{\text{в}} + \Delta = 36{,}96 + 0{,}5818 = 37{,}5418.$$

law -

$= 9.0384$

За таблицями значень функції Лапласа (табл. 3 додатків) знайдемо значення $t$ з рівняння $\Phi(t) = \gamma/2 = 0{,}95/2 = 0{,}475$, $t = 1{,}96$. За формулою (2.32), де $N = 4000$,

$$\Delta = 1{,}96\sqrt{\frac{9{,}0384}{100}\left(1 - \frac{100}{4000}\right)} \approx 0{,}5818,$$

тоді ліва межа довірчого інтервалу

$$\bar{x}_{\text{в}} - \Delta = 36{,}96 - 0{,}5818 = 36{,}3782,$$

right —

$$\bar{x}_{\text{в}} + \Delta = 36{,}96 + 0{,}5818 = 37{,}5418.$$

Finally, the desired confidence interval has the form (36.3782, 37.5418).

2) Quest is the probability $P\left(\left|\bar{x}_{\text{в}} - \bar{x}_{\text{г}}\right| \le 0{,}5\right)$, to be found by the formula

$$P\left(\left|\bar{x}_{\text{в}} - \bar{x}_{\text{г}}\right| < \varepsilon\right) = 2\Phi\left(\frac{\varepsilon}{\sigma}\right),$$

as $\bar{x}_{\text{в}}$ is normally distributed random variable (according to Theorem 2.2), $M(\bar{x}_{\text{в}}) = \bar{x}_{\text{г}}$. Accordingly formula (2.25)

$$\overline{\sigma} = \overline{\sigma}'_{\bar{x}} = \sqrt{\frac{D_{\text{в}}}{n}\left(1 - \frac{n}{N}\right)} = \sqrt{\frac{9,0384}{100}\left(1 - \frac{100}{4000}\right)} \approx 0,2968.$$

Then $\varepsilon/\overline{\sigma} = 0,5/0,2968 = 1,69$ and accordingly table. 3 application

$$P\left(\left|\bar{x}_{\text{в}} - \bar{x}_{\text{г}}\right| \le 0,5\right) = 2\Phi\left(0,5/0,2968\right) =$$

$$= 2\Phi\left(1,69\right) = 2 \cdot 0,45449 = 0,90898.$$

*Table 3.1*

| Inter-val $(x_i; x_{i+1})$ | $n_i$ | $x_{i+1} - \overline{x}_e$ | $x_i - \overline{x}_e$ | $\dfrac{x_{i+1} - \overline{x}_e}{\sigma_e}$ | $\dfrac{x_i - \overline{x}_e}{\sigma_e}$ | $\Phi\left(\dfrac{x_{i+1} - \overline{x}_e}{\sigma_e}\right)$ | $\Phi\left(\dfrac{x_i - \overline{x}_e}{\sigma_e}\right)$ | $p_i = \dots - \Phi$ |
|---|---|---|---|---|---|---|---|---|
| $(-\infty;\ 1{,}720)$ | 3 | –0,0217 | $-\infty$ | –1,916 | $-\infty$ | –0,47231 | –0,5 | |
| $[1{,}720;\ 1{,}727)$ | 13 | –0,0147 | –0,0217 | –1,298 | –1,916 | –0,40285 | –0,47231 | |
| $[1{,}727;\ 1{,}734)$ | 32 | –0,0077 | –0,0147 | –0,680 | –1,298 | –0,25175 | –0,40285 | |
| $[1{,}734;\ 1{,}741)$ | 49 | –0,0007 | –0,0077 | –0,062 | –0,680 | –0,02472 | –0,25175 | |
| $[1{,}741;\ 1{,}748)$ | 42 | 0,0063 | –0,0007 | 0,556 | –0,062 | 0,21089 | –0,02472 | |
| $[1{,}748;\ 1{,}755)$ | 35 | 0,0133 | 0,0063 | 1,174 | 0,556 | 0,3798 | 0,21089 | |
| $[1{,}755;\ 1{,}762)$ | 19 | 0,0203 | 0,0133 | 1,792 | 1,174 | 0,46336 | 0,3798 | |
| $(1{,}762;\ \infty)$ | 7 | $\infty$ | 0,0203 | $\infty$ | 1,792 | 0,5 | 0,46336 | |
| Усього | 200 | — | — | — | — | — | — | |

*Table 3.2*

| $i$ | $n_i$ | $n_i^0$ | $n_i - n_i^0$ | $(n_i - n_i^0)^2$ | $(n_i - n_i^0)^2 / n_i^0$ |
|---|---|---|---|---|---|
| 1 | 3 | 5,54 | –2,54 | 6,4516 | 1,1646 |
| 2 | 13 | 13,89 | –0,89 | 0,7921 | 0,0570 |
| 3 | 32 | 30,22 | 1,78 | 3,1684 | 0,1048 |
| 4 | 49 | 45,41 | 3,59 | 12,8881 | 0,2838 |
| 5 | 42 | 47,12 | –5,12 | 26,2144 | 0,5563 |
| 6 | 35 | 33,78 | 1,22 | 1,4884 | 0,0441 |
| 7 | 19 | 16,71 | 2,29 | 5,2441 | 0,3138 |
| 8 | 7 | 7,33 | –0,33 | 0,1089 | 0,0149 |
| $\sum$ | 200 | 200 | | | $\chi_{cn}^2 = 2{,}5393$ |

# Applications

Value Gaus function $\varphi(x) = \dfrac{1}{\sqrt{2\pi}}\, e^{-\frac{x^2}{2}}$

**Table *1***

| $x$ | Соті долі $x$ | | | | | | | | | |
|-----|---------|-------|-------|-------|-------|-------|-------|-------|-------|-------|
|     | 0 | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 | 9 |
| 0,0 | 0,39894 | 39892 | 39886 | 39876 | 39862 | 39844 | 39822 | 39797 | 39767 | 39733 |
| 0,1 | 39695 | 39654 | 39608 | 39559 | 39505 | 39448 | 39387 | 39322 | 39253 | 39181 |
| 0,2 | 39104 | 39024 | 38940 | 38853 | 38762 | 38667 | 38568 | 38466 | 38361 | 38251 |
| 0,3 | 38139 | 38023 | 37903 | 37780 | 37654 | 37524 | 37391 | 37255 | 37115 | 36973 |
| 0,4 | 36827 | 36678 | 36526 | 36371 | 36213 | 36053 | 35889 | 35723 | 35553 | 35381 |
| 0,5 | 35207 | 35029 | 34849 | 34667 | 34482 | 34294 | 34105 | 33912 | 33718 | 33521 |
| 0,6 | 33322 | 33121 | 32918 | 32713 | 32506 | 32297 | 32086 | 31874 | 31659 | 31443 |
| 0,7 | 31225 | 31006 | 30785 | 30563 | 30339 | 30114 | 29887 | 29659 | 29431 | 29200 |
| 0,8 | 28969 | 28737 | 28504 | 28269 | 28034 | 27798 | 27562 | 27324 | 27086 | 26848 |
| 0,9 | 26609 | 26369 | 26129 | 25888 | 25647 | 25406 | 25164 | 24923 | 24681 | 24439 |
| 1,0 | 24197 | 23955 | 23713 | 23471 | 23230 | 22988 | 22747 | 22506 | 22265 | 22025 |
| 1,1 | 21785 | 21546 | 21307 | 21069 | 20831 | 20594 | 20357 | 20121 | 19886 | 19652 |
| 1,2 | 19419 | 19186 | 18954 | 18724 | 18494 | 18265 | 18037 | 17810 | 17585 | 17360 |
| 1,3 | 17137 | 16915 | 16694 | 16474 | 16256 | 16038 | 15822 | 15608 | 15395 | 15183 |
| 1,4 | 14973 | 14764 | 14556 | 14350 | 14146 | 13943 | 13742 | 13542 | 13344 | 13147 |
| 1,5 | 12952 | 12758 | 12566 | 12376 | 12188 | 12001 | 11816 | 11632 | 11450 | 11270 |
| 1,6 | 11092 | 10915 | 10741 | 10567 | 10396 | 10226 | 10059 | 09893 | 09728 | 09566 |
| 1,7 | 09405 | 09246 | 09089 | 08933 | 08780 | 08628 | 08478 | 08329 | 08183 | 08038 |
| 1,8 | 07895 | 07754 | 07614 | 07477 | 07341 | 07206 | 07074 | 06943 | 06814 | 06687 |
| 1,9 | 06562 | 06438 | 06316 | 06195 | 06077 | 05959 | 05844 | 05730 | 05618 | 05508 |
| 2,0 | 05399 | 05292 | 05186 | 05082 | 04980 | 04879 | 04780 | 04682 | 04586 | 04491 |
| 2,1 | 04398 | 04307 | 04217 | 04128 | 04041 | 03955 | 03871 | 03788 | 03706 | 03626 |
| 2,2 | 03547 | 03470 | 03394 | 03319 | 03246 | 03174 | 03103 | 03034 | 02965 | 02898 |
| 2,3 | 02833 | 02768 | 02705 | 02643 | 02582 | 02522 | 02463 | 02406 | 02349 | 02294 |
| 2,4 | 02239 | 02186 | 02134 | 02083 | 02033 | 01984 | 01936 | 01888 | 01842 | 01797 |
| 2,5 | 01753 | 01709 | 01667 | 01625 | 01585 | 01545 | 01506 | 01468 | 01431 | 01394 |
| 2,6 | 01358 | 01323 | 01289 | 01256 | 01223 | 01191 | 01160 | 01130 | 01100 | 01071 |
| 2,7 | 01042 | 01014 | 00987 | 00961 | 00935 | 00909 | 00885 | 00861 | 00837 | 00814 |
| 2,8 | 00792 | 00770 | 00748 | 00727 | 00707 | 00687 | 00668 | 00649 | 00631 | 00613 |
| 2,9 | 00595 | 00578 | 00562 | 00545 | 00530 | 00514 | 00499 | 00485 | 00470 | 00457 |
| 3,0 | 00443 | 00430 | 00417 | 00405 | 00393 | 00381 | 00370 | 00358 | 00348 | 00337 |
| 3,1 | 00327 | 00317 | 00307 | 00298 | 00288 | 00279 | 00271 | 00262 | 00254 | 00246 |
| 3,2 | 00238 | 00231 | 00224 | 00216 | 00210 | 00203 | 00196 | 00190 | 00184 | 00178 |
| 3,3 | 00172 | 00167 | 00161 | 00156 | 00151 | 00146 | 00141 | 00136 | 00132 | 00127 |

*Table 1*

| x | Соті долі x | | | | | | | | | |
|---|---|---|---|---|---|---|---|---|---|---|
|  | 0 | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 | 9 |
| 3,4 | 00123 | 00119 | 00115 | 00111 | 00107 | 00104 | 00100 | 00097 | 00094 | 00090 |
| 3,5 | 00087 | 00084 | 00081 | 00079 | 00076 | 00073 | 00071 | 00068 | 00066 | 00063 |
| 3,6 | 00061 | 00059 | 00057 | 00055 | 00053 | 00051 | 00049 | 00047 | 00046 | 00044 |
| 3,7 | 00042 | 00041 | 00039 | 00038 | 0037 | 00035 | 00034 | 00033 | 00031 | 00030 |
| 3,8 | 00029 | 00028 | 00027 | 00026 | 00025 | 00024 | 00023 | 00022 | 00021 | 00021 |
| 3,9 | 00020 | 00019 | 00018 | 00018 | 00017 | 00016 | 00016 | 00015 | 00014 | 00014 |
| x | Десяті долі x | | | | | | | | | |
|  | 0 | | 2 | | 4 | | 6 | | 8 | |
| 4, | 0,0001338 | | 0000589 | | 0000249 | | 0000101 | | 0000040 | |
| 5, | 0000015 | | | | | | | | | |

Value function $P(m;\lambda) = \dfrac{\lambda^m e^{-\lambda}}{m!}$

**Table *2***

| m | λ | | | | | | | | | |
|---|---|---|---|---|---|---|---|---|---|---|
|  | 0,1 | 0,2 | 0,3 | 0,4 | 0,5 | 0,6 | 0,7 | 0,8 | 0,9 | 1,0 |
| 0 | 0,90484 | 81873 | 74082 | 67032 | 60653 | 54881 | 49659 | 44933 | 40657 | 36788 |
| 1 | 09048 | 16375 | 22225 | 26813 | 30327 | 32929 | 34761 | 35946 | 36591 | 36788 |
| 2 | 00452 | 01637 | 03334 | 05363 | 07582 | 09879 | 12166 | 14379 | 16466 | 18394 |
| 3 | 00015 | 00109 | 00333 | 00715 | 01264 | 01976 | 02839 | 03834 | 04940 | 06131 |
| 4 |  | 00005 | 00025 | 00072 | 00158 | 00296 | 00497 | 00767 | 01111 | 01533 |
| 5 |  |  | 00002 | 00006 | 00016 | 00036 | 00070 | 00123 | 00200 | 00307 |
| 6 |  |  |  |  | 00001 | 00004 | 00008 | 00016 | 00030 | 00051 |
| 7 |  |  |  |  |  |  | 00001 | 00002 | 00004 | 00007 |
| 8 |  |  |  |  |  |  |  |  |  | 00001 |

**Table 2**

| m | λ | | | | | | | | | |
|---|---|---|---|---|---|---|---|---|---|---|
|  | 1,5 | 2,0 | 2,5 | 3,0 | 3,5 | 4,0 | 4,5 | 5,0 | 5,5 | 6,0 |
| 0 | 0,22313 | 13534 | 08208 | 04979 | 03020 | 01832 | 01111 | 00674 | 00409 | 00248 |
| 1 | 33470 | 27067 | 20521 | 14936 | 10569 | 07326 | 04999 | 03369 | 02248 | 01487 |
| 2 | 25102 | 27067 | 25652 | 22404 | 18496 | 14653 | 11248 | 08422 | 06181 | 04462 |
| 3 | 12551 | 18045 | 21376 | 22404 | 21579 | 19537 | 16872 | 14037 | 11332 | 08924 |
| 4 | 04707 | 09022 | 13360 | 16803 | 18881 | 19537 | 18981 | 17547 | 15582 | 13385 |
| 5 | 01412 | 03609 | 06680 | 10082 | 13217 | 15629 | 17083 | 17547 | 17140 | 16062 |
| 6 | 00353 | 01203 | 02783 | 05041 | 07710 | 10420 | 12812 | 14622 | 15712 | 16062 |
| 7 | 00076 | 00344 | 00994 | 02160 | 03855 | 05954 | 08236 | 10444 | 12345 | 13768 |
| 8 | 00014 | 00086 | 00311 | 00810 | 01687 | 02977 | 04633 | 06528 | 08487 | 10326 |
| 9 | 00002 | 00019 | 00086 | 00270 | 00656 | 01323 | 02316 | 03627 | 05187 | 06884 |
| 10 |  | 00004 | 00022 | 00081 | 00230 | 00529 | 01042 | 01813 | 02853 | 04130 |
| 11 |  | 00001 | 00005 | 00022 | 00073 | 00192 | 00426 | 00824 | 01426 | 02253 |
| 12 |  |  | 00001 | 00006 | 00021 | 00064 | 00160 | 00343 | 00654 | 01126 |
| 13 |  |  |  | 00001 | 00006 | 00020 | 00055 | 00132 | 00277 | 00520 |
| 14 |  |  |  |  | 00001 | 00006 | 00018 | 00047 | 00109 | 00223 |
| 15 |  |  |  |  |  | 00002 | 00005 | 00016 | 00040 | 00089 |
| 16 |  |  |  |  |  |  | 00002 | 00005 | 00014 | 00033 |
| 17 |  |  |  |  |  |  |  | 00001 | 00004 | 00012 |
| 18 |  |  |  |  |  |  |  |  | 00001 | 00004 |
| 19 |  |  |  |  |  |  |  |  |  | 00001 |

Value Laplas function $\Phi(x) = \dfrac{1}{\sqrt{2\pi}} \int\limits_0^x e^{\frac{t^2}{2}} dt$

**Table 3**

| x | Соті долі x | | | | | | | | | |
|---|---|---|---|---|---|---|---|---|---|---|
|  | 0 | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 | 9 |
| 0,0 | 0,00000 | 00399 | 00798 | 01197 | 01595 | 01994 | 02392 | 02790 | 03188 | 03586 |
| 0,1 | 03983 | 04380 | 04776 | 05172 | 05567 | 05962 | 06356 | 06749 | 07142 | 07535 |
| 0,2 | 07926 | 08317 | 08706 | 09095 | 09483 | 09871 | 10257 | 10642 | 11026 | 11409 |
| 0,3 | 11791 | 12172 | 12552 | 12930 | 13307 | 13683 | 14058 | 14431 | 14803 | 15173 |
| 0,4 | 15542 | 15910 | 16276 | 16640 | 17003 | 17364 | 17724 | 18082 | 18439 | 18793 |
| 0,5 | 19146 | 19497 | 19847 | 20194 | 20540 | 20884 | 21226 | 21566 | 21904 | 22240 |
| 0,6 | 22575 | 22907 | 23237 | 23565 | 23891 | 24215 | 24537 | 24857 | 25175 | 25490 |
| 0,7 | 25804 | 26115 | 26424 | 26730 | 27035 | 27337 | 27637 | 27935 | 28230 | 28524 |
| 0,8 | 28814 | 29103 | 29389 | 29673 | 29955 | 30234 | 30511 | 30785 | 31057 | 31327 |
| 0,9 | 31594 | 31859 | 32121 | 32381 | 32639 | 32894 | 33147 | 33398 | 33646 | 33891 |

**Table *3***

| x | Соті долі x | | | | | | | | | |
|---|---|---|---|---|---|---|---|---|---|---|
| | 0 | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 | 9 |
| 1,0 | 34134 | 34375 | 34614 | 34850 | 35083 | 35314 | 35543 | 35769 | 35993 | 36214 |
| 1,1 | 36433 | 36650 | 36864 | 37076 | 37286 | 37493 | 37698 | 37900 | 38100 | 38298 |
| 1,2 | 38493 | 38686 | 38877 | 39065 | 39251 | 39435 | 39617 | 39796 | 39973 | 40147 |
| 1,3 | 40320 | 40490 | 40658 | 40824 | 40988 | 41149 | 41308 | 41466 | 41621 | 41774 |
| 1,4 | 41924 | 42073 | 42220 | 42364 | 42507 | 42647 | 42786 | 42922 | 43056 | 43189 |
| 1,5 | 43319 | 43448 | 43574 | 43699 | 43822 | 43943 | 44062 | 44179 | 44295 | 44408 |
| 1,6 | 44520 | 44630 | 44738 | 44845 | 44950 | 45053 | 45154 | 45254 | 45352 | 45449 |
| 1,7 | 45543 | 45637 | 45728 | 45818 | 45907 | 45994 | 46080 | 46164 | 46246 | 46327 |
| 1,8 | 46407 | 46485 | 46562 | 46638 | 46712 | 46784 | 46856 | 46926 | 46995 | 47062 |
| 1,9 | 47128 | 47193 | 47257 | 47320 | 47381 | 47441 | 47500 | 47558 | 47615 | 47670 |
| 2,0 | 47725 | 47778 | 47831 | 47882 | 47932 | 47982 | 48030 | 48077 | 48124 | 48169 |
| 2,1 | 48214 | 48257 | 48300 | 48341 | 48382 | 48422 | 48461 | 48500 | 48537 | 48574 |
| 2,2 | 48610 | 48645 | 48679 | 48713 | 48745 | 48778 | 48809 | 48840 | 48870 | 48899 |
| 2,3 | 48928 | 48956 | 48983 | 49010 | 49036 | 49061 | 49086 | 49111 | 49134 | 49158 |
| 2,4 | 49180 | 49202 | 49224 | 49245 | 49266 | 49286 | 49305 | 49324 | 49343 | 49361 |
| 2,5 | 49379 | 49396 | 49413 | 49430 | 49446 | 49461 | 49477 | 49492 | 49506 | 49520 |
| 2,6 | 49534 | 49547 | 49560 | 49573 | 49585 | 49598 | 49609 | 49621 | 49632 | 49643 |
| 2,7 | 49653 | 49664 | 49674 | 49683 | 49693 | 49702 | 49711 | 49720 | 49728 | 49736 |
| 2,8 | 49744 | 49752 | 49760 | 49767 | 49774 | 49781 | 49788 | 49795 | 49801 | 49807 |
| 2,9 | 49813 | 49819 | 49825 | 49831 | 49836 | 49841 | 49846 | 49851 | 49856 | 49861 |
| 3,0 | 49865 | 49869 | 49874 | 49878 | 49882 | 49886 | 49889 | 49893 | 49897 | 49900 |
| 3,1 | 49903 | 49906 | 49910 | 49913 | 49916 | 49918 | 49921 | 49924 | 49926 | 49929 |
| 3,2 | 49931 | 49934 | 49936 | 49938 | 49940 | 49942 | 49944 | 49946 | 49948 | 49950 |
| 3,3 | 49952 | 49953 | 49955 | 49957 | 49958 | 49960 | 49961 | 49962 | 49964 | 49965 |
| 3,4 | 49966 | 49968 | 49969 | 49970 | 49971 | 49972 | 49973 | 49974 | 49975 | 49976 |
| 3,5 | 49977 | 49978 | 49978 | 49979 | 49980 | 49981 | 49981 | 49982 | 49983 | 49983 |
| 3,6 | 49984 | 49985 | 49985 | 49986 | 49986 | 49987 | 49987 | 49988 | 49988 | 49989 |
| 3,7 | 49989 | 49990 | 49990 | 49990 | 49991 | 49991 | 49992 | 49992 | 49992 | 49992 |
| 3,8 | 49993 | 49993 | 49993 | 49994 | 49994 | 49994 | 49994 | 49995 | 49995 | 49995 |
| 3,9 | 49995 | 49995 | 49996 | 49996 | 49996 | 49996 | 49996 | 49996 | 49997 | 49997 |
| x | Десяті долі x | | | | | | | | | |
| | 0 | | 2 | | 4 | | 6 | | 8 | |
| 4, | 0,4999683 | | 4999867 | | 4999946 | | 4999979 | | 4999992 | |
| 5, | 4999997 | | | | | | | | | |

*Table **4***

$$P\left(\,|T| < t\right) = 2\int_0^t g_k(t)dt = \gamma,$$

$g_k(t)$ — Student dencity function

**Ст'юдента**, $k = n-1$



| $k = n-1$ | $\gamma$ | | | | |
|:---:|:---:|:---:|:---:|:---:|:---:|
| | 0,9 | 0,95 | 0,98 | 0,99 | 0,999 |
| 1 | 6,314 | 12,706 | 31,821 | 63,657 | 636,619 |
| 2 | 2,920 | 4,303 | 6,965 | 9,925 | 31,599 |
| 3 | 2,353 | 3,182 | 4,541 | 5,841 | 12,924 |
| 4 | 2, 132 | 2,776 | 3,747 | 4,604 | 8,610 |
| 5 | 2,015 | 2,571 | 3,365 | 4,032 | 6,869 |
| 6 | 1,943 | 2,447 | 3,143 | 3,707 | 5,969 |
| 7 | 1,895 | 2,365 | 2,998 | 3,499 | 5,408 |
| 8 | 1,860 | 2,306 | 2,896 | 3,355 | 5,041 |
| 9 | 1,833 | 2,262 | 2,821 | 3,250 | 4,781 |
| 10 | 1,812 | 2,228 | 2,764 | 3,169 | 4,587 |
| 11 | 1,796 | 2,201 | 2,718 | 3,106 | 4,437 |
| 12 | 1,782 | 2,179 | 2,681 | 3,055 | 4,318 |
| 13 | 1,771 | 2,160 | 2,650 | 3,012 | 4,221 |
| 14 | 1,761 | 2,145 | 2,624 | 2,977 | 4,140 |
| 15 | 1,753 | 2,131 | 2,602 | 2,947 | 4,073 |
| 16 | 1,746 | 2,120 | 2,583 | 2,921 | 4,015 |
| 17 | 1,740 | 2,110 | 2,567 | 2,898 | 3,965 |
| 18 | 1,734 | 2,101 | 2,552 | 2,878 | 3,922 |
| 19 | 1,729 | 2,093 | 2,539 | 2,861 | 3,883 |
| 20 | 1,725 | 2,086 | 2,528 | 2,845 | 3,850 |
| 25 | 1,708 | 2,060 | 2,485 | 2,785 | 3,725 |
| 30 | 1,697 | 2,042 | 2,457 | 2,750 | 3,646 |
| 40 | 1,684 | 2,021 | 2,423 | 2,704 | 3,551 |
| 50 | 1,676 | 2,009 | 2,403 | 2,678 | 3,496 |
| 60 | 1,671 | 2,000 | 2,390 | 2,660 | 3,460 |
| 70 | 1,667 | 1,994 | 2,381 | 2,648 | 3,435 |

**Table 4**

| 80 | 1,664 | 1,990 | 2,374 | 2,639 | 3,416 |
|----|-------|-------|-------|-------|-------|
| 90 | 1,662 | 1,987 | 2,368 | 2,632 | 3,402 |
| 100 | 1,660 | 1,984 | 2,364 | 2,626 | 3,390 |
| 120 | 1,658 | 1,980 | 2,358 | 2,617 | 3,373 |
| $\infty$ | 1,645 | 1,960 | 2,326 | 2,576 | 3,291 |

*Table 5*

Value $P\left( \chi^2(k) > \chi^2(p;k) \right) = p,$

$k$ — number of freedom



| $k$ | $p$ | | | | | | | |
|-----|-------|-------|-------|-------|-------|-------|-------|-------|
|     | 0,999 | 0,99 | 0,95 | 0,90 | 0,10 | 0,05 | 0,01 | 0,001 |
| 1 | $0,157 \cdot 10^{-5}$ | 0,0002 | 0,004 | 0,02 | 2,71 | 3,84 | 6,63 | 10,83 |
| 2 | 0,002 | 0,02 | 0,10 | 0,21 | 4,61 | 5,99 | 9,21 | 13,82 |
| 3 | 0,02 | 0,12 | 0,35 | 0,58 | 6,25 | 7,82 | 11,34 | 16,27 |
| 4 | 0,09 | 0,30 | 0,71 | 1,06 | 7,78 | 9,49 | 13,28 | 18,47 |
| 5 | 0,21 | 0,55 | 1,15 | 1,61 | 9,24 | 11,07 | 15,08 | 20,51 |
| 6 | 0,38 | 0,87 | 1,64 | 2,20 | 10,65 | 12,59 | 16,81 | 22,46 |
| 7 | 0,60 | 1,24 | 2,17 | 2,83 | 12,02 | 14,06 | 18,48 | 24,32 |
| 8 | 0,86 | 1,65 | 2,73 | 3,49 | 13,36 | 15,51 | 20,09 | 26,12 |
| 9 | 1,15 | 2,09 | 3,33 | 4,17 | 14,68 | 16,92 | 21,67 | 27,88 |
| 10 | 1,48 | 2,56 | 3,94 | 4,87 | 15,99 | 18,31 | 23,21 | 29,59 |
| 11 | 1,83 | 3,05 | 4,58 | 5,58 | 17,28 | 19,68 | 24,72 | 31,26 |
| 12 | 2,21 | 3,57 | 5,23 | 6,30 | 18,55 | 21,03 | 26,22 | 32,91 |
| 13 | 2,62 | 4,11 | 5,89 | 7,04 | 19,81 | 22,36 | 27,69 | 34,53 |
| 14 | 3,04 | 4,66 | 6,57 | 7,79 | 21,06 | 23,69 | 29,14 | 36,12 |
| 15 | 3,48 | 5,23 | 7,26 | 8,55 | 22,31 | 25,00 | 30,58 | 37,70 |
| 16 | 3,94 | 5,81 | 7,96 | 9,31 | 23,54 | 26,30 | 32,00 | 39,25 |
| 17 | 4,42 | 6,41 | 8,67 | 10,09 | 24,77 | 27,59 | 33,41 | 40,79 |
| 18 | 4,90 | 7,02 | 9,39 | 10,86 | 25,99 | 28,87 | 34,81 | 42,31 |
| 19 | 5,41 | 7,63 | 10,12 | 11,65 | 27,20 | 30,14 | 36,19 | 43,82 |
| 20 | 5,92 | 8,26 | 10,85 | 12,44 | 28,41 | 31,41 | 37,57 | 45,31 |

| | | | | | | | | |
|-----|-------|-------|-------|-------|--------|--------|--------|--------|
| 21  | 6,45  | 8,90  | 11,59 | 13,24 | 29,62  | 32,67  | 38,93  | 46,80  |
| 22  | 6,98  | 9,54  | 12,34 | 14,04 | 30,81  | 33,92  | 40,29  | 48,27  |
| 23  | 7,53  | 10,20 | 13,20 | 14,85 | 32,01  | 35,17  | 41,64  | 49,73  |
| 24  | 8,08  | 10,86 | 13,85 | 15,66 | 33,19  | 36,42  | 43,98  | 51,18  |
| 25  | 8,65  | 11,52 | 14,61 | 16,47 | 34,38  | 37,65  | 44,31  | 52,62  |
| 26  | 9,22  | 12,20 | 15,37 | 17,29 | 35,56  | 38,89  | 45,64  | 54,05  |
| 27  | 9,80  | 12,88 | 16,15 | 18,11 | 36,74  | 40,11  | 46,96  | 55,48  |
| 28  | 10,39 | 13,56 | 16,93 | 18,94 | 37,92  | 41,34  | 48,28  | 56,89  |
| 29  | 10,99 | 14,26 | 17,71 | 19,77 | 39,09  | 42,56  | 49,59  | 58,30  |
| 30  | 11,59 | 14,95 | 18,49 | 20,60 | 40,26  | 43,77  | 50,89  | 59,70  |
| 40  | 17,92 | 22,16 | 26,51 | 29,05 | 51,81  | 55,76  | 63,69  | 73,40  |
| 50  | 24,67 | 29,71 | 34,76 | 37,69 | 63,17  | 67,51  | 76,15  | 86,66  |
| 100 | 61,92 | 70,07 | 77,93 | 82,36 | 118,50 | 124,34 | 135,81 | 149,45 |

REFERENCES

1. P. Bocharov, A. Probability Theory. Mathematical statistics. - M .: Gardarica, 1998. - 328 p.
2. Gmurman V.E. Theory of Probability and mathematical statistics - M .: High society. HQ., 2004. - 479 p.
3.Eremenko V.O., Shynkaryk M.I. Theory of probability. - Economic thought, 2000. - 176 p.
4. Eremenko V.O., Shynkaryk M.I. Theory of probability. - Economic thought, 2002. - 247 p.