

**МИНИСТЕРСТВО ОБРАЗОВАНИЯ РЕСПУБЛИКИ БЕЛАРУСЬ**  
**УЧРЕЖДЕНИЕ ОБРАЗОВАНИЯ**  
**«БРЕСТСКИЙ ГОСУДАРСТВЕННЫЙ ТЕХНИЧЕСКИЙ УНИВЕРСИТЕТ»**  
**КАФЕДРА ВЫСШЕЙ МАТЕМАТИКИ**

**PROBABILITY THEORY**  
**ELEMENTS OF MATHEMATICAL STATISTICS**

учебно-методическая разработка на английском языке

Брест 2014

УДК 519.2(076)=111

Настоящая методическая разработка предназначена для иностранных студентов технических специальностей. Данная разработка содержит необходимый материал по разделам «Теория вероятностей» и «Математическая статистика». Изложение теоретического материала по всем темам сопровождается рассмотрением большого количества примеров и задач, некоторые понятия и примеры проиллюстрированы.

Составители: Гладкий И.И., доцент  
Дворниченко А.В., старший преподаватель  
Каримова Т.И., к.ф.-м.н., доцент  
Лебедь С.Ф., к.ф.-м.н., доцент  
Шишко Т.В., преподаватель

Рецензент: Матысик О.В., заведующий кафедрой прикладной математики и технологий программирования учреждения образования «Брестский государственный университет им. А.С. Пушкина», к.ф.-м.н., доцент.

# 1. EVENT AND PROBABILITY

## 1.1 TRIAL AND EVENT

The probability theory is the science of rules of mass random phenomena. One can say that it's the learning dealing with regularities of mass random phenomena. Known sources of the probability theory are: a) investigations of demographic processes [of population laws]; b) games of chance [games of luck, hazards].

A trial and an event are the **main notions** of the probability theory.

**Def. 1.** A **trial** is a realization of some complex of conditions.

It's supposed that a trial can be arbitrary realized many times.

**Def. 2.** An **event** is every fact, which can occur [appear, happen] or not occur in a trial.

Ex. 1 (see the table).

Trial	Events
1. Coin flip [coin tossing]	"head" (occurrence of a head), "tail"
2. Dice toss(ing), fair dice rolling	"1", "2", "3", "4", "5", "6"
3. Drawing a ball from an urn containing $a$ white and $b$ black balls	"white ball", "black ball"

Events are usually denoted by capitals ( $A, B, C, \dots$ ). There are impossible, certain and random events.

**Def. 3.** An event is called **impossible** if it can't occur in any trial.

Ex. 2. Occurrence of a head **and** a tail in one coin tossing.

**Def. 4.** An event is called **certain** if it necessarily occurs in any trial.

Ex. 3. Occurrence of a head **or** a tail in one coin tossing. Occurrence of at least one of the digits 1, 2, 3, 4, 5, 6 in one dice rolling.

**Def. 5.** An event is called **random** if it can occur or not occur in a trial.

Ex. 4. All events fixed in Ex. 1.

There are joint or disjoint events.

**Def. 6.** Events  $A, B$  are called **joint** [compatible] if they can occur together [or simultaneously] in a trial.

Ex. 5. "head", "head"; "tail", "tail"; "head", "tail"; "tail", "head" if a trial implies double coin tossing.

**Def. 7.** Events  $A, B$  are called **disjoint** [incompatible, non-compatible] if they can't occur together [or simultaneously] in a trial.

Ex. 6. "head", "tail" in one coin toss.

Ex. 7. The events "1", "2", "3", "4", "5", "6" are pairwise disjoint in one dice rolling.

There are **dependent** and **independent** events. See the corresponding definitions below.

**Def. 8.** One says that events  $A, B, \dots, C$  form a total [complete] group (of events) [ $A, B, \dots, C$  are only possible events,  $A, B, \dots, C$  are exhaustive events] if at least one of them occurs in any trial.

Ex. 8. Events "head" and "tail" in one coin toss. All the events "1", "2", "3", "4", "5", "6" in one dice rolling.

**Def. 9.** Two events  $A$  and  $\bar{A}$  (non  $A$ ) are called **opposite** if they are disjoint and form a total group.

Ex. 9. If  $A$  is "head", then  $\bar{A}$  (non  $A$ ) is "tail" (in one coin toss). If  $A$  is "1", then  $\bar{A}$  (non  $A$ ) is the occurrence of at least one of events "2", "3", "4", "5", "6",  $A = \{\text{"2" or "3", or "4", or "5", or "6"}\}$  (in one fair dice rolling).

## 1.2. ELEMENTS OF COMBINATORICS

**Theorem 1 (fundamental principle of combinatorics).** Let an action  $A_1$  be able to be done by  $n_1$  ways, an action  $A_2$  by  $n_2$  ways, ..., an action  $A_k$  by  $n_k$  ways, then all these actions can be done together [or simultaneously] by  $n_1 \cdot n_2 \cdot n_3 \cdots n_k$  ways.

We'll illustrate the validity of this statement with the help of the next example.

Ex. 10. Let's suppose that one has  $a$  coins and  $b$  dice. Then he can take a coin and a die by  $a \cdot b$  ways.

■ Indeed, each coin generates  $1 \cdot b = b$  pairs "coin-die". Therefore,  $a$  coins generate  $a \cdot b$  pairs. ■

Ex. 11. One has 2 coins, 3 ties and 5 books. He can take one coin, one tie and one book by  $2 \cdot 3 \cdot 5 = 30$  ways.

### Main notions of combinatorics

Let there be given some set  $M$  containing  $n$  elements.

**Def. 10. Arrangement** of  $n$  elements (taken)  $k$  at a time [ $k$ -fold arrangement of  $n$  elements] is called any ordered  $k$ -fold subset of the  $n$ -fold set  $M$ .

Various arrangements differ by at least one element or by the order of their elements.

**Def. 11. Permutation** of  $n$  elements is called any arrangement of all  $n$  elements of the  $n$ -fold set  $M$ .

Distinct permutations differ by the order of (the same) elements.

One can say that permutation of  $n$  elements is the ordered set of all elements of the set  $M$ .

**Def. 12. Combination** of  $n$  elements (taken)  $k$  at a time [ $k$ -fold combination of  $n$  elements] is called any  $k$ -fold subset of the  $n$ -fold set  $M$ .

Every combination differs from another one by at least one element.

**Theorem 2.** Numbers of all  $k$ -fold arrangements, of all permutations, of all  $k$ -fold combinations of  $n$  elements are respectively equal

$$A_n^k = n \cdot (n-1) \cdot (n-2) \cdots (n-k+1) = \frac{n!}{(n-k)!} \quad (1)$$

$$P_n = n! \quad (2)$$

$$C_n^k = \frac{n \cdot (n-1) \cdot (n-2) \cdots (n-k+1)}{k!} = \frac{n!}{k!(n-k)!} \quad (3)$$

Ex. 12.  $A_8^3 = \frac{8!}{(8-3)!} = \frac{8!}{5!} = 6 \cdot 7 \cdot 8 = 336$ ,  $P_5 = 5! = 1 \cdot 2 \cdot 3 \cdot 4 \cdot 5 = 120$ ,

$$C_{10}^4 = \frac{10!}{4!(10-4)!} = \frac{10!}{4! \cdot 6!} = \frac{7 \cdot 8 \cdot 9 \cdot 10}{2 \cdot 3 \cdot 4} = 210.$$

Ex. 13. A group containing 25 students can elect the leader and their assistant by

$$A_{25}^2 = \frac{25!}{(25-2)!} = \frac{25!}{23!} = 24 \cdot 25 = 600$$

ways because these two students form 2-fold arrangement of 25 elements.

Ex. 14. One can invite any 4 students of the same group to do some work by

$$C_{25}^4 = \frac{25!}{4!(25-4)!} = \frac{25!}{4! \cdot 21!} = \frac{22 \cdot 23 \cdot 24 \cdot 25}{2 \cdot 3 \cdot 4} = 12650$$

ways because these 4 students form 4-fold combinations of 25 elements.

Ex. 15. 15 competitors of a chess tournament must play

$$C_{15}^2 = \frac{15!}{2!(15-2)!} = \frac{15!}{2! \cdot 13!} = \frac{14 \cdot 15}{2} = 105$$

games in one lap (every two chess players form 2-fold combination of 15 elements).

Ex. 16. 8 books can be placed in a bookshelf by

$$P_8 = 8! = 1 \cdot 2 \cdot 3 \cdot 4 \cdot 5 \cdot 6 \cdot 7 \cdot 8 = 40320$$

ways because they form a permutation of 8 elements.

### 1.3. CLASSIC DEFINITION OF PROBABILITY

There are events for which we can subtract a set of **elementary events (chances, possibilities)** that is a total group of pairwise disjoint and equally possible events. A chance is called **favourable** for an event  $A$  if  $A$  occurs when this chance occurs.

Let  $n$  be the number of all chances [of all elementary events, of all possibilities] and  $m$  be the number of those favourable for some event  $A$ . In this case the probability of this event is expressed by the next ratio:

$$P(A) = \frac{m}{n} \quad (4)$$

Ex. 17. Find the probability of occurrence of the head in one coin-tossing.

Solution. Let  $A$  be an event which means that a head occurs. We can subtract the next  $n = 2$  chances [elementary events, possibilities]: "head", "tail". There is  $m = 1$  favourable chance, namely "head". By the formula (4)

$$P(A) = \frac{m}{n} = \frac{1}{2} = 0.5$$

Ex. 18. Find the probability of occurrence of an even number in one dicerolling.

Solution. Let  $A$  be an event which consists in occurrence of even number in one dicerolling. The chances [elementary events, possibilities] connected with the event  $A$  are "1", "2", "3", "4", "5", "6",  $n = 6$ . The favourable chances are "2", "4", "6",  $m = 3$ . By the formula (4)

$$P(A) = \frac{m}{n} = \frac{3}{6} = 0.5$$

Ex. 19. There are 6 white and 14 black balls in some urn. One takes 10 balls at random. Find the probability of drawing of 4 white and 6 black balls.

Solution. Let  $A$  be an event consisting in drawing of 4 white and 6 black balls.

The chances [elementary events, possibilities] for the event  $A$  are various sets of 10 balls, that is 10-fold combinations of 20 elements. Therefore, the number of all chances is equal to

$$n = C_{20}^{10}$$

that is to number of all 10-fold combinations of 20 elements.

To determine the number  $m$  of favourable chances we must take into account that one can take 4 white balls (4-fold combination of 6 elements) by  $C_6^4$  ways and 6 black balls (6-fold combination of 14 elements) by  $C_{14}^6$  ways. Therefore, he can take 4 white and 6 black balls together by virtue of the fundamental principle of combinatorics by  $C_6^4 \cdot C_{14}^6$  ways. It means that

$$m = C_6^4 \cdot C_{14}^6$$

Hence,

$$P(A) = \frac{m}{n} = \frac{C_6^4 \cdot C_{14}^6}{C_{20}^{10}} = \frac{4! \cdot 2! \cdot 6! \cdot 12!}{20!} \approx 0.24.$$

Ex. 20. It's necessary to place 8 books on a bookshelf. Find the probability for two certain books  $A, B$  to stand side by side.

Solution. Let an event  $CC$  be the required position of our books. The chances [elementary events, possibilities] are their various locations which are permutations of 8 elements. Therefore, the number of all chances is that of all possible permutations of 8 elements,

$$n = P_8 = 8! = 1 \cdot 2 \cdot 3 \cdot 4 \cdot 5 \cdot 6 \cdot 7 \cdot 8 = 40320$$

To find the number of favourable chances (that is that of required positions of books) we'll introduce the next table

Actions to place $A, B$ side by side	Number of ways to do these actions
1. Finding places for $A, B$	7
2. Location of $A, B$ on these places (permutation of 2 elements)	$P_2 = 2!$
3. Disposition of the other 6 books (permutation of 6 elements)	$P_6 = 6!$
Getting of required disposition of all 8 books	$7 \cdot P_2 \cdot P_6$ (by virtue of the main combinatorial principle)

On the base of classical definition of probability

$$P(A) = \frac{m}{n} = \frac{7 \cdot P_2 \cdot P_6}{P_8} = \frac{1}{4} = 0.25.$$

#### 1.4. STATISTIC DEFINITION OF PROBABILITY

Let some event  $A$  be studied and there be fulfilled very large number  $N$  of independent trials on  $A$ . Let's denote as  $N(A)$  the number of occurrences of  $A$  in these trials. The ratio

$$p^* = P_N^*(A) = \frac{N(A)}{N} \quad (5)$$

is called a relative frequency (or sometimes frequency) of the event  $A$ .

Let's fulfil series of very large numbers  $N_1, N_2, \dots$  of independent trials on  $A$  and denote by

$$p_1^* = P_{N_1}^*(A), p_2^* = P_{N_2}^*(A), \dots$$

corresponding relative frequencies of  $A$ .

There are many events for which relative frequencies possess a property of **statistic stability**, that is they are approximately equal to some number  $p$

$$p_1^* \approx p, p_2^* \approx p, \dots$$

If our event  $A$  possesses such a stability property, we say that it has a probability (so-called **statistic probability**), and this probability equals

$$P(A) = p \quad (6)$$

Ex. 21. Many scientists have performed series of very large numbers of coin-tossings (see the table).

Scientist	N	N("head")	$p^* = P_N^*(A)$ ("head")
Buffon G.L.L. (1777)	4040	2048	0.507
de Morgan A. (at the beginning of the 19 <sup>th</sup> century)	4092	2048	0.5005
Pearson K. (at the beginning of the 20 <sup>th</sup> century)	12000	6019	0.5016
Pearson K.	24000	12012	0.50005

On the base of these results we conclude that the statistical probability of the event "head" (in one coin-tossing) equals  $p = P$  ("head") = 0.5, that is coincides with its "classic" probability.

## 2. MAIN RULES OF EVALUATING PROBABILITIES

### 2.1. SUM AND PRODUCT OF EVENTS

**Def. 1.** Sum  $A + B$  ( $A$  or  $B$ ) of two events  $A$  and  $B$  is called an event which consists in occurrence of at least one of them [which means that at least one of these events occurs] ( $A$  but not  $B$  or  $B$  but not  $A$  or  $A$  and  $B$  together).

Ex. 1. Sum of an event  $A$  and its opposite one  $\bar{A}$  is a certain event.

Ex. 2. If an event  $A$  is "1" in one dice rolling, then the opposite event  $\bar{A}$  is the sum  $\bar{A} = "2" + "3" + "4" + "5" + "6"$ .

**Def. 2.** Product  $AB$  ( $A$  and  $B$ ) of two events  $A$  and  $B$  is called an event consisting in occurrence of both these events [an event which means that both these events occur] together.

Ex. 3. Product of an event  $A$  and its opposite one  $\bar{A}$  is an impossible event.

Ex. 4 (*Euler circles*). Let  $M$  and  $N$  be two circles having non-empty intersection  $V = M \cap N$ , also let

$U = M \setminus N$ ,  $W = N \setminus M$ , and so  $M = U \cup V$ ,  $N = V \cup W$  (see fig. 1). If an event  $A$  means that a point  $P$  belongs to  $M$ , and  $B$  means that  $P$  belongs to  $N$ , then

$$A + B = \{P \in (U \cup V \cup W)\},$$

$$AB = \{P \in V\}.$$

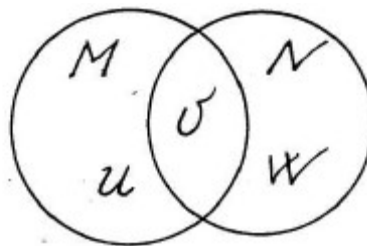


Fig. 1

Solving probabilistic problems it is sometimes useful to represent an event in question as a sum of product of other events (with pairwise disjoint summands).

Ex. 5.  $A + B = A \cdot \bar{B} + \bar{A} \cdot B + A \cdot B$ .

Ex. 6. Let events  $A, B, C$  mean that the first, second, third device (correspondingly) works. In this case the events

$$D = ABC, E = \bar{A}\bar{B}\bar{C}, F = \bar{A}BC + A\bar{B}\bar{C} + \bar{A}\bar{B}C, G = \bar{A}BC + A\bar{B}\bar{C} + \bar{A}\bar{B}C$$

$F + G + D = \bar{E}$  mean respectively that all the three devices work, none device works, only one device works (and the other two don't work), two devices work (and one doesn't work), at least one device works.

## 2.2. AXIOMS OF PROBABILITY THEORY. COROLLARIES

We state axioms of the probability theory on the base of statistic definition of probability ( $P(A) \approx P_N^*(A)$  for a large number  $N$  of trials).

1. If  $A$  is an **impossible** event, then its probability equals zero,  $P(A) = 0$  ( $A$  is impossible).
2. If  $A$  is a **certain** event, then its probability equals unity,  $P(A) = 1$  ( $A$  is certain).
3. If  $A$  is a **random** event, then its probability is contained between zero and unity,  $0 \leq P(A) \leq 1$  ( $A$  is random).
4. If  $A$  and  $B$  are two disjoint events, then the probability of their sum is equal to the sum of probabilities of these events,  $P(A + B) = P(A) + P(B)$  ( $A, B$  are disjoint).

To formulate the last axiom let's introduce the notion of a **conditional probability** of an event. Namely,  $P(B / A)$  is the probability of an event  $B$  by condition that an event  $A$  occurs. Analogous is the probability of  $A$  if  $B$  occurs.

5. Probability of a product of two events equals the product of the probability of one event and the condition probability of the other,  $P(AB) = P(A) \cdot P(B / A) = P(B) \cdot P(A / B)$ .

Ex. 7. An urn contains 3 white and 2 black balls. One takes two balls successively and at random. Find the probability that they are white.

Solution. Let's denote by  $A$  an event which means that two drawn balls are white. Also let's denote by  $B$  and  $C$  events which mean that the first and the second drawn balls are white respectively. It's evident that  $A = BC$ , hence, by virtue of the fifth axiom (and classic definition of probability)

$$P(BC) = P(B) \cdot P(C / B) = P(C) \cdot P(B / C) = \frac{3}{5} \cdot \frac{2}{4} = 0.3.$$

### Some corollaries

1. If some events  $A, B, C$  are **pairwise disjoint**, then  $P(A + B + C) = P(A) + P(B) + P(C)$ . If, moreover, they **form a total group**, then  $P(A + B + C) = P(A) + P(B) + P(C) = 1$ .

2. The sum of probabilities of two opposite events equals 1,

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

because the events  $A, \bar{A}$  are disjoint and form a total group.

3. Probabilities of a product of three, four etc events are equal to

$$P(ABC) = P(A) \cdot P(B / A)P(C / AB)$$

4. For two arbitrary events  $A$  and  $B$  the probability of their sum equals

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

**Def 3.** Two events  $A, B$  are called independent if probability of one of them doesn't depend on occurrence (or non-occurrence) of the other.

For 3, 4 ... events one introduces a notion of mutual independence (independence in the aggregate, collectionwise independence).



**Def 4.**  $n$  events (for  $n \geq 2$ ) are called mutually independent if the probability of one of them doesn't depend on occurrence or non-occurrence of any group of the other.

5. If  $A, B$  are independent events, then  $P(B / A) = P(B), P(A / B) = P(A)$  and so

$$P(AB) = P(A) \cdot P(B),$$

that is the probability of a product of two **independent** events is equal to the product of their probabilities.

6. If  $A, B, C$  are mutually independent events, then

$$P(ABC) = P(A) \cdot P(B) \cdot P(C)$$

Ex. 8. To pass an exam successfully a student has to know the proofs of 50 theorems but he knows only 40 of them. What is the probability for him to pass an exam if exam tasks contain 3 theorems?

Solution. Let  $A$  be an event which means that a student will pass an exam. Let's introduce the next three auxiliary events:  $B_1$  that is a student knows the proof of the first theorem,  $B_2$  of the second,  $B_3$  of the third. Then by the fourth corollary (and classic definition of probability)

$$A = B_1 B_2 B_3 \Rightarrow P(A) = P(B_1 B_2 B_3) = P(B_1) \cdot P(B_2 / B_1) \cdot P(B_3 / B_1 B_2)$$

$$P(A) = \frac{40}{50} \cdot \frac{39}{49} \cdot \frac{38}{48} \approx 0.5.$$

Ex. 9. A device consists of 2 independent modules. Probability for these modules to work are 0.95 and 0.9 respectively. Find the probability that the device doesn't work because of: a) only one module; b) at least one module.

Solution. Let an event  $A$  mean that a device doesn't work because of only one module, and an event  $B$  because of at least one module. Let events  $C_1, C_2$  mean that the first, the second module works. By condition

$$P(C_1) = 0.95, P(C_2) = 0.9$$

and so by the corollary 3

$$P(\overline{C_1}) = 1 - P(C_1) = 1 - 0.95 = 0.05, P(\overline{C_2}) = 1 - P(C_2) = 1 - 0.9 = 0.1$$

We represent the events  $A, B$  and  $\overline{B}$  as follows

$$A = C_1 \overline{C_2} + \overline{C_1} C_2, B = C_1 \overline{C_2} + \overline{C_1} C_2 + \overline{C_1} \overline{C_2} \Rightarrow \overline{B} = C_1 C_2.$$

All summands are pairwise disjoint, and all factors are independent in summands.

Therefore,

$$P(A) = P(C_1 \overline{C_2} + \overline{C_1} C_2) = P(C_1 \overline{C_2}) + P(\overline{C_1} C_2) = 0.05 \cdot 0.9 + 0.1 \cdot 0.95 = 0.14$$

$$P(\overline{B}) = P(C_1 C_2) = 0.95 \cdot 0.9 = 0.86 \Rightarrow P(B) = 1 - P(\overline{B}) = 1 - 0.86 = 0.14.$$

Ex. 10. Three independently working engines are installed in a workshop. Probabilities to work at a given time equal for them 0.6, 0.9, 0.7 respectively. Find probabilities of the next events: a) only one engine works; b) at least one engine works.

Solution. Let an event  $A$  mean that only one engine works and an event  $B$  mean that at least one engine works. Our problem is to find the probabilities of these events.

Let's introduce three auxiliary events, namely  $C_1$  which means that the first engine works,  $C_2$  the second engine works and  $C_3$  the third engine works. By conditions of the problem

$$P(C_1) = 0.6, P(\overline{C_1}) = 1 - P(C_1) = 1 - 0.6 = 0.4$$

$$P(C_2) = 0.9, P(\overline{C_2}) = 1 - P(C_2) = 1 - 0.9 = 0.1$$

$$P(C_3) = 0.7, P(\overline{C_3}) = 1 - P(C_3) = 1 - 0.7 = 0.3$$

a) The event  $A$  can be represented as the sum of products

$$A = C_1\overline{C_2}C_3 + \overline{C_1}C_2C_3 + \overline{C_1}C_2\overline{C_3}$$

with pairwise disjoint summands and independent factors in every summand. Hence the probability of the event  $A$  equals

$$P(A) = P(C_1\overline{C_2}C_3 + \overline{C_1}C_2C_3 + \overline{C_1}C_2\overline{C_3}) = P(C_1\overline{C_2}C_3) + P(\overline{C_1}C_2C_3) + P(\overline{C_1}C_2\overline{C_3})$$

$$P(A) = 0.6 \cdot 0.1 \cdot 0.3 + 0.4 \cdot 0.9 \cdot 0.3 + 0.4 \cdot 0.1 \cdot 0.7 = 0.154.$$

b) To find the probability of the event  $B$  we'll evaluate at first the probability of its opposite one  $\overline{B}$  (which means that all three engines don't work,  $\overline{B} = \overline{C_1}C_2\overline{C_3}$ ).

We'll obtain

$$P(\overline{B}) = P(\overline{C_1}C_2\overline{C_3}) = 0.4 \cdot 0.1 \cdot 0.3 = 0.012$$

whence it follows that

$$P(B) = 1 - P(\overline{B}) = 1 - 0.012 = 0.988.$$

### 2.3. FORMULAE OF TOTAL PROBABILITY AND BAYES

#### **The formula of total probability**

In practice we often deal with the next situation. An event  $A$  can occur only together with one of pairwise disjoint events  $H_1, H_2, \dots, H_n$ , which form a total group.

Let's call these events **hypotheses**. Their probabilities and corresponding conditional probabilities of the event  $A$  are known. In this case the probability of the event  $A$  can be found with the help of the next formula (the **formula of total probability**):

$$P(A) = P(H_1)P(A/H_1) + P(H_2)P(A/H_2) + \dots + P(H_n)P(A/H_n) \quad (4)$$

$$(P(H_1) + P(H_2) + \dots + P(H_n) = 1)$$

#### **Bayes formulae**

Let an event  $A$ , which can occur only together with one of given hypotheses  $H_1, H_2, \dots, H_n$ , occur. In this case the next probabilities  $P(H_k / A)$  of its occurrence together with each of these hypotheses can be evaluated with the help of the known Bayes formulas

$$P(H_k / A) = \frac{P(H_k)P(A/H_k)}{P(A)}, \quad k = \overline{1, n}. \quad (5)$$

Bayes formulae (5) state the probability that namely the  $k$ -th hypothesis has occurred together with the event  $A$  in question.

Ex. 11. There are 6 white and 2 black balls in the first urn and 8 white and 3 black balls in the second urn. One moves a ball from the first urn to the second one at random, and then he takes a ball from the second urn (also at random).

1. Find the probability for him to take a white ball from the second urn.

2. Let a white ball be taken from the second urn. A ball of which colour was most probably moved from the first urn?

1. Solution of the first problem. Let an event  $A$  mean that one will take a white ball from the second urn. We can introduce the next two hypotheses:  $H_1$  means that one has moved a white ball from the first urn;  $H_2$  that he has moved a black ball from there. Their probabilities equal

$$P(H_1) = \frac{6}{8} = \frac{3}{4}; P(H_2) = \frac{2}{8} = \frac{1}{4}$$

by condition, and corresponding conditional probabilities of the event  $A$  equal

$$P(A/H_1) = \frac{9}{12} = \frac{3}{4}; P(A/H_2) = \frac{8}{12} = \frac{2}{3}.$$

On the base of the formula (4) of total probability the probability of the event  $A$  equals

$$P(A) = P(H_1)P(A/H_1) + P(H_2)P(A/H_2) = \frac{3}{4} \cdot \frac{3}{4} + \frac{1}{4} \cdot \frac{2}{3} = 0.73.$$

2. To solve the second problem we must find and compare the next conditional probabilities  $P(H_1/A)$ ,  $P(H_2/A)$ . On the base of Bayes formulae (5)

$$P(H_1/A) = \frac{P(H_1)P(A/H_1)}{P(A)} = \frac{\frac{3}{4} \cdot \frac{3}{4}}{0,73} \approx 0.77$$

$$P(H_2/A) = \frac{P(H_2)P(A/H_2)}{P(A)} = \frac{\frac{1}{4} \cdot \frac{2}{3}}{0,73} \approx 0,23.$$

We see that  $0.77 > 0.23$ , therefore one has the most probably moved a white ball from the first urn to the second one.

Ex. 12. There are 10 and 15 products of the first and second factories respectively in the storage. The first factory makes 5% and the second 7% of defective products. One takes a product at random.

1) Find the probability of its defectiveness.

2) Suppose that this product is defective. Which factory has most probably made it?

Solution. Let an event  $A$  mean that a product taken at random is defective. Let's introduce the next two hypotheses:  $H_1$  this product was done by the first factory;  $H_2$  it was done by the second factory.

On the base of the classical definition of probability

$$P(H_1) = \frac{10}{25} = 0.4; P(H_2) = \frac{15}{25} = 0.6.$$

Corresponding conditional probabilities of the event  $A$  equal

$$P(A/H_1) = \frac{5}{100} = 0.05; P(A/H_2) = \frac{7}{100} = 0.07.$$

1) Using the formula of total probability we'll get

$$P(A) = P(H_1)P(A/H_1) + P(H_2)P(A/H_2) = 0.4 \cdot 0.05 + 0.6 \cdot 0.07 = 0.062.$$

2) Now with the help of Bayes formulas we find

$$P(H_1/A) = \frac{P(H_1)P(A/H_1)}{P(A)} = \frac{0.4 \cdot 0.05}{0,062} \approx 0.32$$

$$P(H_2/A) = \frac{P(H_2)P(A/H_2)}{P(A)} = \frac{0.6 \cdot 0.07}{0,062} \approx 0,68.$$

Thus,

$$P(H_2/A) > P(H_1/A)$$

Therefore, the taken product was most probably made by the second factory.

### Exercise Set 1, 2.

1. The first worker makes 40% of the second-class parts, and the second makes 30%. Two parts are taken from each worker at random. Find the probability that: a) all the four parts are second-class; b) at least three parts are second-class; c) less than three parts are second-class.

2. A worker services three machine tools. The probability that during his shift the machine tools will claim his attention is equal to 0.7 for the first one, 0.65 for the second one, 0.55 for the third one. Find the probability that during his shift his attention will be claimed by: a) two machine tools; b) not less than two machine tools; c) at least one machine tool.

3. The probability that the student will pass the examinations is equal to 0.8 for the first one, 0.7 for the second one, 0.65 for the third one. Find the probability that the student will pass: a) two examinations; b) not less than two examinations; c) at least one examination.

4. A shooter fires at the target three times. Probabilities of hitting the target are respectively equal to  $p_1$ ,  $p_2$ ,  $p_3$ . Find the probabilities of the following events: 1) hitting the target three times; 2) hitting the target not less than two times; 3) hitting the target at least once.

	$p_1$	$p_2$	$p_3$
A	0,8	0,85	0,9
B	0,7	0,8	0,9
C	0,9	0,75	0,8
D	0,6	0,8	0,7

5. Articles from three conveyors enter for assembling. The number of articles given for assembling is  $\alpha\%$  for the first one,  $\beta\%$  for the second one,  $\gamma\%$  for the third one. On the average, the number of defective articles is  $\delta_1\%$  from the first conveyor;  $\delta_2\%$  from the second conveyor;  $\delta_3\%$  from the third conveyor. Find the probability that a defective part has entered for assembling. What is the probability that the part from the  $i$ -th of conveyor is defective?

	$\alpha$	$\beta$	$\gamma$	$\delta_1$	$\delta_2$	$\delta_3$	$i$
A	30	15	55	2	2	3	3
B	40	40	20	1	3	2	1
C	50	30	20	2	4	3	2
D	20	45	35	3	5	2	3
E	35	35	30	2	3	5	2

### Homework Problems

**Exercise 1.** A fair die is rolled 5 times and the sequence of scores recorded.

(a) How many outcomes are there?

(b) Find the probability that first and last rolls are 6.

**Exercise 2.** If a 3-digit number (000 to 999) is chosen at random, find the probability that exactly one digit will be larger than 5.

**Exercise 3.** A license plate is made of 3 numbers followed by 3 letters.

(a) What is the total number of possible license plates ?

(b) What is the number of license plates that start with an A?

**Exercise 4.** A lottery is played as follows: the player picks six numbers out of  $f_1, 2, \dots, 54g$ . Then, six numbers are drawn at random out of the 54. You win the first prize if you have 6 correct numbers and the second prize if you get 5 of them.

- (a) What is the probability to win the first prize ?
- (b) What is the probability to win the second prize ?

**Exercise 5.** Another lottery is played as follows: the player picks five numbers out of  $f_1, 2, \dots, 50g$  and two other numbers from the list  $f_1, \dots, 9g$ . Then, five numbers are drawn at random from the first list and two from the random list.

- (a) You win the first prize if all numbers are correct. What is the probability to win the first prize ?
- (b) Which lottery would you choose to play between this one and the one from the previous problem ?

**Exercise 6.** An urn contains 3 red, 8 yellow and 13 green balls; another urn contains 5 red, 7 yellow and 6 green balls. We pick one ball from each urn at random. Find the probability that both balls are of the same color.

**Exercise 7.** Suppose that there are 5 duck hunters, each a perfect shot. A flock of 10 ducks fly over, and each hunter selects one duck at random and shoots. Find the probability that 5 ducks are killed.

**Exercise 8.** A conference room contains  $m$  men and  $w$  women. These people seat at random in  $m + w$  seats arranged in a row. Find the probability that all the women will be adjacent.

**Exercise 9.** If a box contains 75 good light bulbs and 25 defective bulbs and 15 bulbs are removed, find the probability that at least one will be defective.

**Exercise 10.** Find the probability that a five-card poker hand (i.e. 5 out of a 52-card deck) will be :

- (a) Four of a kind, that is four cards of the same value and one other card of a different value (xxxxy shape).
- (b) Three of a kind, that is three cards of the same value and two other cards of different values (xxxzy shape).
- (c) A straight flush, that is five cards in a row, of the same suit (ace may be high or low).
- (d) A flush, that is five cards of the same suit, but not a straight flush.
- (e) A straight, that is five cards in a row, but not a straight flush (ace may be high or low).

**Exercise 11.** An urn contains 10 balls numbered from 1 to 10. We draw five balls from the urn, without replacement. Find the probability that the second largest number drawn is 8.

**Exercise 12.** Eight cards are drawn without replacement from an ordinary deck. Find the probability of obtaining exactly three aces or exactly three kings (or both).

**Exercise 13.** How many possible ways are there to seat 8 people (A,B,C,D,E,F,G and H) in a row, if:

- (a) No restrictions are enforced;
- (b) A and B want to be seated together;
- (c) assuming there are four men and four women, men should be only seated between women and the other way around;
- (d) assuming there are five men, they must be seated together;
- (e) assuming these people are four married couples, each couple has to be seated together.

**Exercise 14.** John owns six discs: 3 of classical music, 2 of jazz and one of rock (all of them different). How many possible ways does John have if he wants to store these discs on a shelf, if:

- (a) No restrictions are enforced;
- (b) The classical discs and the jazz discs have to be stored together;
- (c) The classical discs have to be stored together, but the jazz discs have to be separated.

**Exercise 15.** How many (not necessarily meaningful) words can you form by shuffling the letters of the following words: (a) bike; (b) paper; (c) letter; (d) minimum.

**Exercise 16.** An urn contains 30 white and 15 black balls. If 10 balls are drawn with (respectively without) replacement, find the probability that the first two balls will be white, given that the sample contains exactly six white balls.

**Exercise 17.** In a certain village, 20% of the population has some disease. A test is administered which has the property that if a person is sick, the test will be positive 90% of the time and if the person is not sick, then the test will still be positive 30% of the time. All people tested positive are prescribed a drug which always cures the disease but produces a rash 25% of the time. Given that a random person has the rash, what is the probability that this person had the disease to start with?

**Exercise 18.** An insurance company considers that people can be split in two groups : those who are likely to have accidents and those who are not. Statistics show that a person who is likely to have an accident has probability 0.4 to have one over a year; this probability is only 0.2 for a person who is not likely to have an accident. We assume that 30% of the population is likely to have an accident.

(a) What is the probability that a new customer has an accident over the first year of his contract?

(b) A new customer has an accident during the first year of his contract. What is the probability that he belongs to the group likely to have an accident?

**Exercise 19.** A transmitting system transmits 0's and 1's. The probability of a correct transmission of a 0 is 0.8, and it is 0.9 for a 1. We know that 45% of the transmitted symbols are 0's.

(a) What is the probability that the receiver gets a 0?

(b) If the receiver gets a 0, what is the probability the transmitting system actually sent a 0?

**Exercise 20.** 46% of the electors of a town consider themselves as independent, whereas 30% consider themselves democrats and 24% republicans. In a recent election, 35% of the independents, 62% of the democrats and 58% of the republicans voted.

(a) What proportion of the total population actually voted?

(b) A random voter is picked. Given that he voted, what is the probability that he is independent? democrat? republican?

### 3. RANDOM VARIABLES

#### 3.1. A RANDOM VARIABLE

**Def. 1.** A random variable is a variable which takes on some value in any trial and this value isn't known beforehand [in advance].

We'll denote random variables by letters  $X, Y, Z, \dots$  and their possible values by  $x, y, z, \dots$ . We'll study discrete and continuous random variables.

**Def. 2.** A random variable is called **discrete** if it can take on only separate isolated possible values (with some probabilities).

Ex.1. A number  $X$  of occurrences of an event  $A$  in one trial:  $X = 1$  if  $A$  occurs, and  $X = 0$ , if  $A$  doesn't occur (that is an opposite event  $\bar{A}$  occurs);

$$P(X = 1) = P(A), P(X = 0) = P(\bar{A}) = 1 - P(A)$$

Ex. 2. The number of students on the lecture.

Ex. 3. The daily production of some factory (in items).

Definition of a **continuous** random variable will be given below. Now we'll only say that its possible values fill some interval completely.

Ex. 4. The human height and weight.

Ex. 5. The size of an item.

Ex.6. The error of measurement.

**Def. 3.** The distribution [the distribution law, the law of distribution, the law] of a random variable is a rule which sets a correspondence between its possible values and corresponding probabilities.

The distribution law of a random variable can be expressed:

1) analytically by a formula (for example  $P_n(X = m) = C_n^m p^m q^{n-m}$ ,  $q = 1 - p$ );

2) tabularly (by the **distribution table** for discrete random variables);

3) geometrically (by **distribution polygon** for discrete random variables, by graph of the distribution function or density (see below)).

**The distribution table** of a discrete random variable  $X$  (with finite number  $n$  of possible values) has the next form:

$X$	$x_1$	$x_2$	...	$x_n$
$P$	$p_1$	$p_2$	...	$p_n$

Its first row contains all possible values of the random variable, and the second row contains corresponding probabilities of these values. The notation  $X = x_i$  means that the random variable  $X$  takes on a value  $x_i$ . Events

$$(X = x_1), (X = x_2), \dots, (X = x_n)$$

are pairwise disjoint and form a total group. Therefore, the sum of their probabilities

$$p_1 = P(X = x_1), p_2 = P(X = x_2), \dots, p_n = P(X = x_n) \text{ equals } 1 (\sum p_i = 1).$$

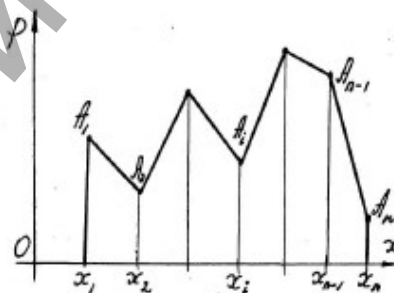


Fig. 2

**The distribution polygon** of a discrete random variable is a broken line [a polygonal line, an open polygon] which is generated by successive joining of the points

$$A_1(x_1, p_1), A_2(x_2, p_2), \dots, A_n(x_n, p_n) \text{ (fig. 2).}$$

Ex. 7. An urn contains 7 balls (namely 3 white and 4 black balls). One draws 3 balls at random. Find the distribution law of the number  $X$  of white balls which can be taken from the urn.

Solution. Possible values of the random variable  $X$  are 0, 1, 2, 3. We determine corresponding probabilities

$$p_1 = P(X = 0), p_2 = P(X = 1), p_3 = P(X = 2), p_4 = P(X = 3)$$

with the help of the classical definition of probability. Elementary events (chances) for every of these four cases are sets of 3 balls that is 3-fold combinations of 7 elements. Hence the general number of chances equals

$$n = C_7^3 = \frac{7!}{3! \cdot 4!} = 35.$$

Numbers of favourable chances are represented in the table

Event	Number of favorable chances	Explication
$X = 0$	$m_1 = C_4^3 = C_4^1 = 4$	One can take 0 white (and so 3 black) balls by the number of 3-fold combinations of 4 elements
$X = 1$	$m_2 = 3 \cdot C_4^2 = 18$	One can take 1 white ball by 3 ways and 2 black balls by $C_4^2$ ways
$X = 2$	$m_3 = 4 \cdot C_3^2 = 12$	One can take 2 white balls by $C_3^2$ ways and 1 black ball by 4 ways
$X = 3$	$m_4 = 1$	One can take 3 white balls in one unique way

Numbers  $m_2, m_3$  are calculated with the help of the main principle of combinatorics.

The distribution law of the random variable  $X$  is represented by the next distribution table:

$X$	0	1	2	3
$p_i$	$\frac{4}{35}$	$\frac{18}{35}$	$\frac{12}{35}$	$\frac{1}{35}$

The sum of obtained probabilities equals 1:

$$\sum_{i=1}^4 p_i = \frac{4 + 18 + 12 + 1}{35} = 1.$$

The most probable value of the random variable is  $X = 1$ .

The distribution polygon of the random variable  $X$  is shown on fig.3.

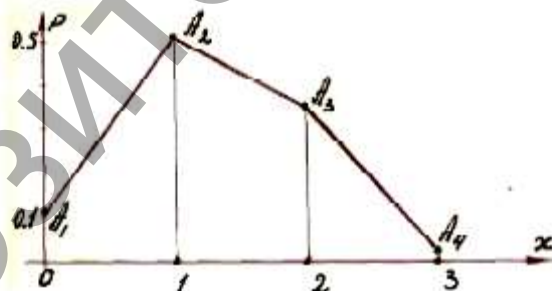


Fig.3

### 3.2. BERNOULLI [BINOMIAL] DISTRIBUTION

**Def.4.** Let a random variable  $X$  be the number of successes (the number of occurrences of some event  $A$ ) in  $n$  independent trials with constant probability of the success  $A$  in any trial

$$p = P(A), P(\bar{A}) = 1 - p = q$$

One says that  $X$  is distributed binomially (by Bernoulli [binomial] law) or simply:  $X$  is Bernoulli (binomial) distribution (briefly:  $X$  **distr. B**).

Let's find a probability  $P(X = m) = P_n(m)$ , that is the probability of  $m$  successes.

We'll get so-called Bernoulli formula:

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



**Note 1.** For  $m = 0, 1, 2, \dots$  the value of the probability (1) first increases and then decreases, and therefore there exists the so-called most probable number  $m_0$  of successes.

It is defined by the next double inequality

$$np - q \leq m_0 \leq np + p, P_n(m_0) = \max \quad (2)$$

**Note 2.** Probability of no less than  $k_1$  and no greater than  $k_2$  successes equals

$$P_n(k_1 \leq m \leq k_2) = P_n(k_1) + P_n(k_1 + 1) + P_n(k_1 + 2) + \dots + P_n(k_2) \quad (3)$$

Ex. 8. 6 independently working engines are installed in a shop. Probability for any engine to work at a given moment is 0.8. Find the distribution of a random variable  $X$  which is a number of working engines at this moment. Find the probabilities of at least one engine to work, of no less than 2 and no greater than 5 engines to work and the most probable value of  $X$ .

Solution. We can consider setting of an engine as a trial. So we have  $n = 6$  independent trials. Let a success  $A$  mean that an engine works.

$$p = P(A) = 0.8, P(\bar{A}) = 1 - p = q = 0.2$$

The random variable  $X$  has Bernoulli distribution (briefly " $X$  distributed  $B$ "), it can take on the values 0, 1, 2, 3, 4, 5, 6, which one calculates by Bernoulli formula (1). For example

$$P(X = 0) = P_6(0) = C_6^0 p^0 q^6 = 0.2^6 = 0.00006$$

$$P(X = 1) = P_6(1) = C_6^1 p^1 q^5 = 6 \cdot 0.8 \cdot 0.2^5 = 0.00154$$

By the same way the other probabilities are calculated, and the distribution table of the random variable is the next one

$X$	0	1	2	3	4	5	6
$p_i$	0.00006	0.00154	0.01536	0.08192	0.24576	0.39321	0.26214

The probability of at least one engine to work equals

$$P(X > 0) = P(X \geq 1) = 1 - P(X = 0) = 1 - q^6 = 0.99994$$

The most probable value of  $X$ ,  $m_0 = 5$ , we see in the table. Evaluation of  $m_0$  by the formula (2) gives  $np - q \leq m_0 \leq np + p \Rightarrow 6 \cdot 0.8 - 0.2 \leq m_0 \leq 6 \cdot 0.8 + 0.8$  whence it follows that

$$4.6 \leq m_0 \leq 5.6 \Rightarrow m_0 = 5.$$

Probability of no less than 2 and no greater than 5 engines to work on the base of the formula (3) equals

$$\begin{aligned} P_6(2 \leq m \leq 5) &= P_6(2) + P_6(3) + P_6(4) + P_6(5) = \\ &= 0.01536 + 0.08192 + 0.24576 + 0.39321 + 0.73625 \approx 0.74 \end{aligned}$$

Bernoulli formula (1) isn't convenient for large number  $n$  of trials. There are some approximate formulas.

### 3.3. POISSON FORMULA AND DISTRIBUTION

Let a random variable  $X$  be distributed  $B$ . Let's suppose that the number  $n$  of trials tends to infinity, the probability  $p$  of a success  $A$  goes to zero, but a product  $np$  retains constant,

$$n \rightarrow \infty, p \rightarrow 0, np = \text{const} = \lambda.$$

In this case the limit of the probability  $P(X = m) = P_n(m)$ , which is defined by Bernoulli formula (1), equals

$$P_n(m) \approx \frac{\lambda^m}{m!} e^{-\lambda} \quad (4)$$

**Def. 5.** One says that a discrete random variable  $X$  (with non-negative integer possible values) has Poisson distribution with a parameter  $\lambda$  if its distribution law is given by the next formula (Poisson formula):

$$P(X = m) = P_n(m) = \frac{\lambda^m}{m!} e^{-\lambda} \quad (5)$$

Let a random variable  $X$  have Bernoulli distribution (short:  $X$  is distributed  $B$ ), the number of trials  $n$  is large, the probability  $p$  of a success is small and  $np \leq 10$ . In this case we can consider  $X$  as having Poisson distribution with the parameter  $\lambda = np$  and use the formula (5) instead of (1). By this reason Poisson distribution is sometimes called the **law of rare events**.

Table of the values of the function  $P_k = \frac{a^k}{k!} e^{-a}$

$k \backslash a$	0,1	0,2	0,3	0,4	0,5	0,6	0,7	0,8	0,9
0	0,9048	0,8187	0,7408	0,6703	0,6065	0,5488	0,4966	0,4493	0,4066
1	0,0905	0,1638	0,2222	0,2681	0,3033	0,3293	0,3476	0,3595	0,3659
2	0,0045	0,0164	0,0333	0,0536	0,0758	0,0988	0,1217	0,1438	0,1647
3	0,0002	0,0019	0,0033	0,0072	0,0126	0,0198	0,0284	0,0383	0,0494
4		0,0001	0,0002	0,0007	0,0016	0,0030	0,0050	0,0077	0,0111
5				0,0001	0,0002	0,0004	0,0007	0,0012	0,0020
6						0,0001	0,0001	0,0002	0,0003

$k \backslash a$	1	2	3	4	5	6	7	8	9	10
0	0,3679	0,1353	0,0498	0,0183	0,0067	0,0025	0,0009	0,0003	0,0001	0,0000
1	0,3679	0,2707	0,1494	0,0733	0,0337	0,0149	0,0064	0,0027	0,0011	0,0005
2	0,1839	0,2707	0,2240	0,1465	0,0842	0,0446	0,0223	0,0107	0,0050	0,0023
3	0,0613	0,1804	0,2240	0,1954	0,1404	0,0892	0,0521	0,0286	0,0150	0,0076
4	0,0153	0,0902	0,1680	0,1954	0,1755	0,1339	0,0912	0,0572	0,0337	0,0189
5	0,0031	0,0361	0,1008	0,1563	0,1755	0,1606	0,1277	0,0916	0,0607	0,0378
6	0,0005	0,0120	0,0504	0,1042	0,1462	0,1606	0,1490	0,1221	0,0911	0,0631
7	0,0001	0,0037	0,0216	0,0595	0,1044	0,1377	0,1490	0,1396	0,1171	0,0901
8		0,0009	0,0081	0,0298	0,0653	0,1033	0,1304	0,1396	0,1318	0,1126
9		0,0002	0,0027	0,0132	0,0363	0,0688	0,1014	0,1241	0,1318	0,1251
10			0,0008	0,0053	0,0181	0,0413	0,0710	0,0993	0,1186	0,1251
11			0,0002	0,0019	0,0082	0,0225	0,0452	0,0722	0,0970	0,1137
12			0,0001	0,0006	0,0034	0,0126	0,0263	0,0481	0,0728	0,0948
13				0,0002	0,0013	0,0052	0,0142	0,0296	0,0504	0,0729
14				0,0001	0,0005	0,0022	0,0071	0,0169	0,0324	0,0521
15					0,0002	0,0009	0,0033	0,0090	0,0194	0,0347
16						0,0003	0,0014	0,0045	0,0109	0,0217
17						0,0001	0,0006	0,0021	0,0058	0,0128
18							0,0002	0,0009	0,0029	0,0071
19							0,0001	0,0004	0,0014	0,0037
20								0,0002	0,0006	0,0019
21								0,0001	0,0003	0,0009
22									0,0001	0,0004
23										0,0002
24										0,0001

Ex. 9. 500 items are sent. Probability of damage of an item in the trade is 0.002. Find the probabilities: a) 3, b) less than 3, c) more than 2 items are damaged; d) at least one item is damaged.

Solution. We have  $n = 500$  independent trials, a success  $A$  is a damage of an item,  $p = P(A) = 0.002 = \text{const}$ , a random variable  $X$  is a number of damaged items.  $X$  is distributed  $B$ , but  $n$  is large,  $p$  is small,  $np = 1$ . Therefore we can consider  $X$  as Poisson distribution with  $\lambda = 1$  and make use of Poisson formula (5).

a)  $P(X = 3) = 0.0613$

b)  $P(X < 3) = P(X \leq 2) = P(X = 0) + P(X = 1) + P(X = 2) = 0.3679 + 0.3679 + 0.1839 = 0.9197$

c)  $P(X > 2) = 1 - P(X \leq 2) = 1 - 0.9197 = 0.0803$

d)  $P(X \geq 1) = P(X > 0) = 1 - P(X = 0) = 1 - 0.3679 = 0.6321$ .

### 3.4 LAPLACE LOCAL AND INTEGRAL THEOREMS

**Laplace local theorem** gives an approximate value of the probability  $P(X = m) = P_n(m)$  of  $m$  successes in  $n$  independent trials (with constant probability  $p$  of the success in any trial). Namely, for large  $n$  we can substitute Bernoulli formula (1) by the next approximate one:

$$P(X = m) = P_n(m) \approx \frac{1}{\sqrt{npq}} \varphi(x) \quad (6)$$

where

$$x = \frac{m - np}{\sqrt{npq}}$$

and

$$\varphi(x) = \frac{1}{\sqrt{2\pi}} e^{-x^2/2} \quad (7)$$

is so-called small Laplace function. The function is even one, that is  $\varphi(-x) = \varphi(x)$  and therefore its graph is symmetric with respect to the Oy-axis. The function is tabulated one. We can suppose  $\varphi(x) = 0, |x| \geq 4$ .

**Integral Laplace theorem** gives an approximate value of the probability of hitting of values of Bernoulli distribution in a segment  $[m_1, m_2]$ . Namely for large  $n$  we can substitute the exact formula (6) by the next approximate formula:

$$P_n(m_1 \leq m \leq m_2) = \Phi(x_2) - \Phi(x_1) \quad (8)$$

where

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

and a function

$$\Phi(x) = \frac{1}{\sqrt{2\pi}} \int_0^x e^{-t^2/2} dt \quad (9)$$

is so-called **Laplace function** (or **normed Laplace function**). The function is an odd one, that is  $\Phi(-x) = -\Phi(x)$ . Laplace function is tabulated. We can suppose

$$\Phi(x) = 0.5, x \geq 5.$$

Table of the values of the function  $\varphi(x) = \frac{1}{\sqrt{2\pi}} e^{-x^2/2}$

x	0	1	2	3	4	5	6	7	8	9
0,0		3989	3989	3988	3986	3984	3982	3980	3977	3973
0,1	0,3989	3965	3961	3956	3951	3945	3939	3932	3925	3918
0,2	3970	3902	3894	3885	3876	3867	3857	3847	3836	3825
0,3	3910	3802	3790	3778	3765	3752	3739	3726	3712	3697
0,4	3814	3668	3653	3637	3621	3605	3589	3572	3555	3538
0,5	3683	3503	3485	3467	3448	3429	3410	3391	3372	3352
0,6	3521	3312	3292	3271	3251	3230	3209	3187	3166	3144
0,7	3332	3101	3079	3056	3034	3011	2989	2966	2943	2920
0,8	3123	2874	2850	2827	2803	2780	2756	2732	2709	2685
0,9	2897	2637	2613	2589	2565	2541	2516	2492	2468	2444
	2661									
1,0		2396	2371	2347	2323	2299	2275	2251	2227	2203
1,1	0,2420	2155	2331	2107	2083	2059	2036	2012	1989	1965
1,2	2179	1919	1895	1872	1849	1826	1804	1781	1758	1736
1,3	1942	1691	1669	1647	1626	1604	1582	1561	1539	1518
1,4	1714	1476	1456	1435	1415	1394	1374	1354	1334	1315
1,5	1497	1276	1257	1238	1219	1200	1182	1163	1145	1127
1,6	1295	1092	1074	1057	1040	1023	1006	0989	0973	0957
1,7	1109	0925	0909	0893	0878	0863	0848	0833	0818	0804
1,8	0940	0775	0761	0748	0734	0721	0707	0694	0681	0669
1,9	0790	0644	0632	0620	0608	0596	0584	0573	0562	0551
	0656									
2,0		0529	0519	0508	0498	0488	0478	0468	0459	0449
2,1	0,0540	0431	0422	0413	0404	0396	0387	0379	0371	0363
2,2	0440	0347	0339	0332	0325	0317	0310	0303	0297	0290
2,3	0355	0277	0270	0264	0258	0252	0246	0241	0235	0229
2,4	0283	0219	0213	0208	0203	0198	0194	0189	0184	0180
2,5	0224	0171	0167	0163	0158	0154	0151	0147	0143	0139
2,6	0175	0132	0129	0126	0122	0119	0116	0113	0110	0107
2,7	0136	0101	0099	0096	0093	0091	0088	0086	0084	0081
2,8	0104 0079	0077	0075	0073	0071	0069	0067	0065	0063	0061
2,9	0060	0058	0056	0055	0053	0051	0050	0048	0047	0046
3,0	0,0044	0043	0042	0040	0039	0038	0037	0036	0035	0034
3,1	0033	0032	0031	0030	0029	0028	0027	0026	0025	0025
3,2	0024	0023	0022	0022	0021	0020	0020	0019	0018	0018
3,3	0017	0017	0016	0016	0015	0015	0014	0014	0013	0013
3,4	0012	0012	0012	0011	0011	0010	0010	0010	0009	0009
3,5	0009	0008	0008	0008	0008	0007	0007	0007	0007	0006
3,6	0006	0006	0006	0005	0005	0005	0005	0005	0005	0004
3,7	0004	0004	0004	0004	0004	0004	0003	0003	0003	0003
3,8	0003	0003	0003	0003	0003	0002	0002	0002	0002	0002
3,9	0002	0002	0002	0002	0002	0002	0002	0001	0001	0001

Table of the values of the function  $\Phi(x) = \frac{1}{\sqrt{2\pi}} \int_0^x e^{-\frac{t^2}{2}} dt$

x	$\Phi(x)$	x	$\Phi(x)$	x	$\Phi(x)$	x	$\Phi(x)$	x	$\Phi(x)$	x	$\Phi(x)$
0,00	0,0000	0,45	0,1736	0,90	0,3159	1,35	0,4115	1,80	0,4641	2,50	0,4938
0,01	0,0040	0,46	0,1772	0,91	0,3186	1,36	0,4131	1,81	0,4649	2,52	0,4941
0,02	0,0080	0,47	0,1808	0,92	0,3212	1,37	0,4147	1,82	0,4656	2,54	0,4945
0,03	0,0120	0,48	0,1844	0,93	0,3238	1,38	0,4162	1,83	0,4664	2,56	0,4948
0,04	0,0160	0,49	0,1879	0,94	0,3264	1,39	0,4177	1,84	0,4671	2,58	0,4951
0,05	0,0199	0,50	0,1915	0,95	0,3289	1,40	0,4192	1,85	0,4678	2,60	0,4953
0,06	0,0239	0,51	0,1950	0,96	0,3315	1,41	0,4207	1,86	0,4686	2,62	0,4956
0,07	0,0279	0,52	0,1985	0,97	0,3340	1,42	0,4222	1,87	0,4693	2,64	0,4959
0,08	0,0319	0,53	0,2019	0,98	0,3365	1,43	0,4236	1,88	0,4699	2,66	0,4961
0,09	0,0359	0,54	0,2054	0,99	0,3389	1,44	0,4251	1,89	0,4706	2,68	0,4963
0,10	0,0398	0,55	0,2088	1,00	0,3413	1,45	0,4265	1,90	0,4713	2,70	0,4965
0,11	0,0438	0,56	0,2123	1,01	0,3438	1,46	0,4279	1,91	0,4719	2,72	0,4967
0,12	0,0478	0,57	0,2157	1,02	0,3461	1,47	0,4292	1,92	0,4726	2,74	0,4969
0,13	0,0517	0,58	0,2190	1,03	0,3485	1,48	0,4306	1,93	0,4732	2,76	0,4971
0,14	0,0557	0,59	0,2224	1,04	0,3508	1,49	0,4319	1,94	0,4738	2,78	0,4973
0,15	0,0596	0,60	0,2257	1,05	0,3531	1,50	0,4332	1,95	0,4744	2,80	0,4974
0,16	0,0636	0,61	0,2291	1,06	0,3554	1,51	0,4345	1,96	0,4750	2,82	0,4976
0,17	0,0675	0,62	0,2324	1,07	0,3577	1,52	0,4357	1,97	0,4756	2,84	0,4977
0,18	0,0714	0,63	0,2357	1,08	0,3599	1,53	0,4370	1,98	0,4761	2,86	0,4979
0,19	0,0753	0,64	0,2389	1,09	0,3621	1,54	0,4382	1,99	0,4767	2,88	0,4980
0,20	0,0793	0,65	0,2422	1,10	0,3643	1,55	0,4394	2,00	0,4772	2,90	0,4981
0,21	0,0832	0,66	0,2454	1,11	0,3665	1,56	0,4406	2,02	0,4783	2,92	0,4982
0,22	0,0871	0,67	0,2486	1,12	0,3686	1,57	0,4418	2,04	0,4793	2,94	0,4984
0,23	0,0910	0,68	0,2517	1,13	0,3708	1,58	0,4429	2,06	0,4803	2,96	0,4985
0,24	0,0948	0,69	0,2549	1,14	0,3729	1,59	0,4441	2,08	0,4812	2,98	0,4986
0,25	0,0987	0,70	0,2580	1,15	0,3749	1,60	0,4452	2,10	0,4821	3,00	0,4987
0,26	0,1026	0,71	0,2611	1,16	0,3770	1,61	0,4463	2,12	0,4830	3,20	0,4993
0,27	0,1064	0,72	0,2642	1,17	0,3790	1,62	0,4474	2,14	0,4838	3,40	0,4997
0,28	0,1103	0,73	0,2673	1,18	0,3810	1,63	0,4484	2,16	0,4846	3,60	0,4998
0,29	0,1141	0,74	0,2703	1,19	0,3830	1,64	0,4495	2,18	0,4854	3,80	0,4999
0,30	0,1179	0,75	0,2734	1,20	0,3849	1,65	0,4515	2,20	0,4861	4,00	0,4999
0,31	0,1217	0,76	0,2764	1,21	0,3869	1,66	0,4505	2,22	0,4868	4,50	0,5000
0,32	0,1255	0,77	0,2794	1,22	0,3883	1,67	0,4525	2,24	0,4875	5,00	0,5000
0,33	0,1293	0,78	0,2823	1,23	0,3907	1,68	0,4535	2,26	0,4881		
0,34	0,1331	0,79	0,2852	1,24	0,3925	1,69	0,4545	2,28	0,4887	↓	↓
0,35	0,1368	0,80	0,2881	1,25	0,3944	1,70	0,4554	2,30	0,4893	+∞	0,5
0,36	0,1406	0,81	0,2910	1,26	0,3962	1,71	0,4564	2,32	0,4898		
0,37	0,1443	0,82	0,2939	1,27	0,3980	1,72	0,4573	2,34	0,4904		
0,38	0,1480	0,83	0,2967	1,28	0,3997	1,73	0,4582	2,36	0,4909		
0,39	0,1517	0,84	0,2995	1,29	0,4015	1,74	0,4591	2,38	0,4913		
0,40	0,1554	0,85	0,3023	1,30	0,4032	1,75	0,4599	2,40	0,4918		
0,41	0,1591	0,86	0,3051	1,31	0,4049	1,76	0,4608	2,42	0,4922		
0,42	0,1628	0,87	0,3078	1,32	0,4066	1,77	0,4616	2,44	0,4927		
0,43	0,1654	0,88	0,3106	1,33	0,4082	1,78	0,4625	2,46	0,4931		
0,44	0,1700	0,89	0,3133	1,34	0,4099	1,79	0,4633	2,48	0,4934		

Ex. 10. The probability of the appearance of an event in each of 245 independent trials is constant and equal to 0,25. Find the probability that an event will begin exactly 50 times.

Solution. Using Laplace local and integral theorems we have

$$n = 245, m = 50, p = 0,25, q = 1 - p = 0,75$$

$$x = \frac{50 - 245 \cdot 0,25}{\sqrt{245 \cdot 0,75 \cdot 0,25}} = \frac{-11,25}{\sqrt{45,9375}} = -\frac{11,25}{6,778} \approx -1,66$$

The function is an even one  $\varphi(-1,66) = \varphi(1,66)$ .

$$P_{245}(50) \approx \frac{1}{\sqrt{npq}} \varphi(1,66) = \frac{0,1006}{6,778} \approx 0,0148$$

Ex. 11. The probability of the appearance of an event in each of 245 independent trials is constant and equal to 0,25. Find the probability that an event will begin not less than 45 times and not more than 60 times.

Solution. Using integral Laplace theorem we have

$$n = 245, m_1 = 45, m_2 = 60, p = 0,25, q = 1 - p = 0,75$$

$$x_1 = \frac{45 - 245 \cdot 0,25}{\sqrt{245 \cdot 0,75 \cdot 0,25}} = \frac{-16,25}{\sqrt{45,9375}} \approx -2,40$$

$$x_2 = \frac{60 - 245 \cdot 0,25}{\sqrt{245 \cdot 0,75 \cdot 0,25}} = \frac{-1,25}{\sqrt{45,9375}} \approx -0,18$$

The function is an odd one. Using the table of Laplace function, we get

$$\Phi(-2,40) = -\Phi(2,40) \approx -0,4918$$

$$\Phi(-0,18) = -\Phi(0,18) \approx -0,0714$$

$$P_{245}(45 \leq m \leq 60) \approx \Phi(-0,18) - \Phi(-2,40) = -0,0714 + 0,4918 = 0,4204$$

**The probability of the deviation of the relative frequency of an event  $A$  from its probability  $p = P(A)$  (in  $n$  independent trials with constant probability  $p = P(A) = \text{const}$  of the event) can be find by the next formula**

$$P\left(\left|\frac{m}{n} - p\right| \leq \varepsilon\right) \approx 2\Phi\left(\varepsilon \sqrt{\frac{n}{pq}}\right) \quad (10)$$

### Exercise Set 3.

A. The probability of the appearance of an event in each of  $n$  of independent trials is constant and equal to  $p$ . Find the probability that an event will begin exactly  $m$  times.

Exercise	$n$	$m$	$p$
1	144	120	0,8
2	110	18	0,15
3	220	140	0,6
4	112	13	0,1
5	99	17	0,2
6	117	85	0,7
7	240	80	0,3
8	115	100	0,9
9	62	5	0,1
10	154	90	0,6

B. The probability of the appearance of an event in each of  $n$  of independent trials is constant and equal to  $p$ . Find the probability that an event will begin not less than  $m_1$  times and not more than  $m_2$  times.

Exercise	$n$	$m_1$	$m_2$	$p$
11	144	115	125	0,8
12	110	15	20	0,15
13	220	130	145	0,6
14	112	10	14	0,1
15	99	15	20	0,2
16	117	80	100	0,7
17	240	70	90	0,3
18	115	100	110	0,9
19	62	5	10	0,1
20	154	80	100	0,6

C. The probability of the production of a defective part is equal to  $p$ . Find the probability that from the tested  $n$  parts  $m$  are defective.

Exercise	$n$	$m$	$p$
21	1000	6	0,008
22	2500	2	0,001
23	1500	10	0,006
24	3500	5	0,002
25	10000	4	0,0005
26	8000	6	0,0008
27	4500	5	0,0008
28	2000	1	0,0001
29	5000	3	0,0008
30	7000	4	0,0006

## 4. THE DISTRIBUTION FUNCTION AND DENSITY. NUMBER CHARACTERISTICS OF RANDOM VARIABLES

### 4.1. THE DISTRIBUTION FUNCTION OF A RANDOM VARIABLE

The distribution function is the most general form of the distribution law of a random variable  $X$ .

**Def. 1.** The distribution function of a random variable  $X$  is called a function:

$$F(x) = P(X < x) = P(-\infty < X < x) \quad (1)$$

The distribution function  $F(x)$  is the probability for the random variable  $X$  to take on values which are less than  $x$  or the probability of hitting of the random variable in the infinite interval (the half-axis)  $(-\infty, x)$ .

#### **Properties of the distribution function**

1. The distribution function, being a probability, lies between 0 and 1 [ranges from 0 to 1]:

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

2.  $\lim_{x \rightarrow -\infty} F(x) = 0$ ,  $\lim_{x \rightarrow +\infty} F(x) = 1$

$$3. F(x_1) \leq F(x_2) \text{ if } x_1 < x_2$$

$$4. P(\alpha < X < \beta) = F(\beta) - F(\alpha) \quad (2)$$

Ex.1. From the urn, which contains 3 white and 5 black spheres, 3 spheres are extracted. Let random variable  $X$  be the number of taken out black spheres. Find the distribution law. Plot the graph of function of distribution.

Solution. Possible values of the random variable  $X$  are 0, 1, 2, 3. We determine corresponding probabilities

$$p_1 = P(X = 0), p_2 = P(X = 1), p_3 = P(X = 2), p_4 = P(X = 3)$$

with the help of the classical definition of probability.

$$p_1 = P(X = 0) = \frac{C_3^3}{C_8^3} = \frac{1 \cdot 2 \cdot 3}{8 \cdot 7 \cdot 6} = \frac{1}{56}.$$

$$p_2 = P(X = 1) = \frac{C_5^1 \cdot C_3^2}{C_8^3} = \frac{5 \cdot 3}{56} = \frac{15}{56}.$$

$$p_3 = P(X = 2) = \frac{C_5^2 \cdot C_3^1}{C_8^3} = \frac{5 \cdot 4 \cdot 3}{2 \cdot 56} = \frac{15 \cdot 2}{56} = \frac{30}{56}.$$

$$p_4 = P(X = 3) = \frac{C_5^3}{C_8^3} = \frac{5 \cdot 4 \cdot 3}{2 \cdot 3} \cdot \frac{1}{56} = \frac{10}{56}.$$

The distribution law of the random variable  $X$  is represented by the next distribution table:

X	0	1	2	3
P	$\frac{1}{56}$	$\frac{15}{56}$	$\frac{30}{56}$	$\frac{10}{56}$

The distribution function of the random variable  $X$ , namely of the number of shots which can be done in reality and its graph are given as follows:

$$\text{If } x \in (0; 1], \text{ then } F(x) = P(X = 0) = \frac{1}{56}.$$

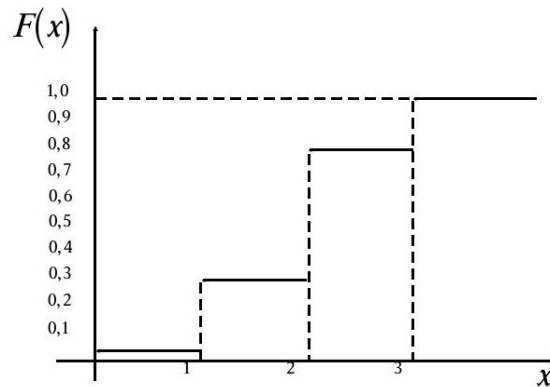
$$\text{If } x \in (1; 2], \text{ then } F(x) = P(X = 0) + P(X = 1) = \frac{1}{56} + \frac{15}{56} = \frac{16}{56}.$$

$$\text{If } x \in (2; 3], \text{ then } F(x) = \frac{16}{56} + P(X = 2) = \frac{16}{56} + \frac{30}{56} = \frac{46}{56}.$$

$$\text{If } x \in (3; +\infty], \text{ then } F(x) = \frac{46}{56} + P(X = 3) = \frac{46}{56} + \frac{10}{56} = 1.$$

$$F(x) = \begin{cases} 0, & \text{if } x \leq 0, \\ \frac{1}{56} = 0,018, & \text{if } 0 < x \leq 1, \\ \frac{16}{56} = 0,286, & \text{if } 1 < x \leq 2, \\ \frac{46}{56} = 0,821, & \text{if } 2 < x \leq 3, \\ 1, & \text{if } 3 < x < +\infty. \end{cases}$$





## 4.2. THE DISTRIBUTION DENSITY OF A RANDOM VARIABLE

Let  $X$  be a continuous random variable and  $F(x)$  its distribution function. The probability of hitting of the random variable in an infinitely small interval  $(x, x + \Delta x)$  (fig. 4)



Fig. 4

by the formula (2) equals

$$P(x < X < x + \Delta x) = F(x + \Delta x) - F(x) = \Delta F(x)$$

The average density of this probability on the interval  $(x, x + \Delta x)$  equals

$$f_{av}(x) = \frac{\Delta F(x)}{\Delta x}$$

**Def. 2.** The distribution density (the density of probability) of a continuous random variable  $X$  is called the derivative of its distribution function,

$$f(x) = F'(x) \quad (3)$$

It follows from the definition 2 that:

1. The distribution function of a continuous random variable is a primitive (an antiderivative) of its distribution density.

2. Accurate [with an accuracy] to infinitely small of higher order

$$P(x < X < x + \Delta x) = f(x)\Delta x \quad (4)$$

the expression in the right side of the formula (4), namely  $f(x)dx$ , (5) is called a **probability element**. It's the differential of the distribution function of the random variable which we consider,

$$f(x)dx = F'(x)dx \quad (5)$$

### Properties of the distribution density

1. The distribution density is a non-negative function,  $f(x) \geq 0$

$$2. P(a < X < b) = \int_a^b f(x)dx \quad (6)$$

$$3. F(x) = \int_{-\infty}^x f(t)dt \quad (7)$$

$$4. \int_{-\infty}^{\infty} f(x)dx = 1 \quad (8)$$

Ex. 2. Let there be given a function  $f(x) = a e^{-|x|}$ . Find the value of the parameter  $a$  so that the function can be the distribution density of some continuous random variable.

Solution. The function in question is a non-negative one. To be the distribution density it must satisfy the condition (8). We find the value of  $a$

$$\int_{-\infty}^{\infty} a e^{-|x|} dx = 1, \quad a \int_{-\infty}^0 e^x dx + a \int_0^{\infty} e^{-x} dx = 1$$

$$a e^x \Big|_{-\infty}^0 - a e^{-x} \Big|_0^{\infty} = a(1-0) - a(0-1) = 2a = 1, \quad a = \frac{1}{2}.$$

Ex. 3. Let there be given a function

$$f(x) = \begin{cases} \frac{4}{x^5}, & \text{if } x \geq 1, \\ 0, & \text{if } x < 1. \end{cases}$$

Find its distribution function.

Solution. By virtue of the formula (7) the distribution function of the random variable is

$$F(x) = \int_1^x f(t) dt = \int_1^x \frac{4}{t^5} dt = 4 \cdot \frac{t^{-4}}{-4} \Big|_1^x = \left( -\frac{1}{t^4} \right) \Big|_1^x = 1 - \frac{1}{x^4}$$

$$F(x) = \begin{cases} 0, & \text{if } x \leq 1, \\ 1 - \frac{1}{x^4}, & \text{if } x > 1. \end{cases}$$

### 4.3. THE MATHEMATICAL EXPECTATION OF A RANDOM VARIABLE

#### Definition of the mathematical expectation

Let's suppose that  $n$  independent trials on a random variable  $X$  are fulfilled and obtained results are represented by the next table:

$X$	$x_1$	$x_2$	...	$x_n$
$P$	$p_1$	$p_2$	...	$p_n$

According to the statistic definition of probability we introduce the next definition.

**Def. 3.** The mathematical expectation of a discrete random variable  $X$  is defined by the next expression

$$M(X) = \sum_{i=1}^n x_i p_i \quad (9)$$

which is the sum of products of its possible values and corresponding probabilities of these values.

Let  $X$  be a continuous random variable with the distribution density  $f(x)$ . We get the mathematical expectation of a continuous random variable in the form of the improper integral

$$M(X) = \int_{-\infty}^{\infty} x f(x) dx, \quad \text{if } X \in (-\infty, \infty) \quad (10)$$

$$M(X) = \int_a^b x f(x) dx, \quad \text{if } X \in (a; b). \quad (11)$$

**Probability [probabilistic] sense of the mathematical expectation of a random variable:** it is its mean [average] value.

**Properties of the mathematical expectation**

1.  $M(C) = C$ .
2.  $M(CX) = C M(X)$ ,  $C = const$ .
3.  $M(X + Y) = M(X) + M(Y)$ ,  $X$  and  $Y$  – random variable ;
4.  $M(X \cdot Y) = M(X) \cdot M(Y)$ ,  $X$  and  $Y$  – are independent random variables.

Ex. 4. The mathematical expectation of the next discrete random variable

X	0	1	2	3
P	$\frac{1}{56}$	$\frac{15}{56}$	$\frac{30}{56}$	$\frac{10}{56}$

on the base of the same formula (9) equals

$$M(X) = 0 \cdot \frac{1}{56} + 1 \cdot \frac{15}{56} + 2 \cdot \frac{30}{56} + 3 \cdot \frac{10}{56} = \frac{15 + 60 + 30}{56} = \frac{105}{56} = 1,875$$

Ex. 5. The distribution density of a continuous random variable  $X$  equals (see Ex. 3)

$$f(x) = \begin{cases} \frac{4}{x^5}, & \text{if } x \geq 1, \\ 0, & \text{if } x < 1. \end{cases}$$

In this case the mathematical expectation must be calculated by the formula (10),

$$M(X) = \int_1^{\infty} xf(x) dx = \int_1^{\infty} \frac{4}{x^4} dx = 4 \cdot \frac{x^{-3}}{-3} \Big|_1^{\infty} = -\frac{4}{3x^3} \Big|_1^{\infty} = \frac{4}{3}$$

**4.4 THE DISPERSION AND ROOT-MEAN-SQUARE DEVIATION**

**Def. 4.** The deviation of a random variable  $X$  from its mathematical expectation is the next random variable

$$X - M(X) \tag{12}$$

**Def. 5.** The dispersion<sup>1</sup> of a random variable  $X$  is the mathematical expectation of the square of its deviation (12) from its mathematical expectation

$$D(X) = M((X - M(X))^2) \tag{13}$$

**Probability [probabilistic] sense of the dispersion:** the dispersion characterizes [describes] a dissipation of a random variable about its mathematical expectation, that is about its mean [or average] value.

**Def. 6.** The root-mean-square deviation<sup>1</sup> of a random variable  $X$  is the square root of its dispersion, that is

$$\sqrt{D(X)} = \sigma(X) \tag{14}$$

**Properties of the dispersion**

1. The dispersion of a constant quantity  $C$  equals zero,  $D(C) = 0$ .
2.  $D(CX) = C^2 D(X)$ ,  $C = const$ .
3.  $D(X + Y) = D(X) + D(Y)$ , where  $X$  and  $Y$  are independent random variables.
4.  $D(X) = M(X^2) - M^2(X)$ .

Let  $X$  be a continuous random variable with the distribution density  $f(x)$ . We get the dispersion of a continuous random variable in the form of the improper integral

$$D(X) = M((X - M(X))^2) = \int_{-\infty}^{\infty} (x - M(X))^2 f(x) dx = \int_{-\infty}^{\infty} x^2 f(x) dx - M^2(X) \quad (15)$$

The dispersion of a discrete random variable  $X$  is the next expression

$$D(X) = M(X^2) - M^2(X) = \sum_{i=1}^n x_i^2 p_i - M^2(X) \quad (16)$$

Ex. 6. Calculate the dispersion and root-mean-square deviation of the random variable  $X$  of Ex. 4.

Solution. The distribution tables of  $X$  and its square (see Ex. 4) are

X	0	1	2	3
P	$\frac{1}{56}$	$\frac{15}{56}$	$\frac{30}{56}$	$\frac{10}{56}$

By the formulas (16) and (14) we obtain

$$M(X^2) = 1 \cdot \frac{15}{56} + 4 \cdot \frac{30}{56} + 9 \cdot \frac{10}{56} = \frac{15 + 120 + 90}{56} = \frac{225}{56} = 4,018.$$

$$D(X) = 4,018 - 1,875^2 = 0,5024.$$

$$\sigma(X) = \sqrt{0,5024} = 0,71.$$

Ex. 7. Calculate the dispersion and root-mean-square deviation of the random variable  $X$  of Ex. 5.

Solution. The distribution density of the random variable is

$$f(x) = \begin{cases} \frac{4}{x^5}, & \text{if } x \geq 1, \\ 0, & \text{if } x < 1. \end{cases}$$

The integral of the formula (15) equals

$$D(X) = M(X^2) - \left(\frac{4}{3}\right)^2 = \int_1^{\infty} x^2 f(x) dx - \frac{16}{9} = \int_1^{\infty} \frac{4}{x^3} dx - \frac{16}{9} = 4 \cdot \frac{x^{-2}}{-2} \Big|_1^{\infty} - \frac{16}{9} = -\frac{4}{2x^2} \Big|_1^{\infty} - \frac{16}{9} = 2 - \frac{16}{9} = \frac{2}{9}$$

Therefore, the root-mean-square deviation of the random variable  $X$  equals

$$\sigma(X) = \sqrt{D(X)} = \sqrt{\frac{2}{9}} \approx 0.471.$$

#### 4.5 MOMENTS OF A RANDOM VARIABLE

**Def. 7.** The  $n$ th order **initial moment** of a random variable  $X$  is the mathematical expectation of its  $n$ th power,

$$\alpha_n = M(X^n) \quad (17)$$

**Def. 8.** The  $n$ th order central moment of a random variable  $X$  is the mathematical expectation of the  $n$ th power of its centered random variable,

$$\mu_n = M((X - M(X))^n) \quad (18)$$

**Theorem.** If the distribution of a random variable is symmetric about its mathematical expectation, then all its odd-order central moments are equal to zero.

Central moments can be expressed in terms of initial moments, for example those of the second, third and fourth orders are equal

$$\begin{aligned}\mu_2 &= \alpha_2 - 2\alpha_1\alpha_1 + \alpha_1^2 = \alpha_2 - \alpha_1^2, \\ \mu_3 &= \alpha_3 - 3\alpha_2\alpha_1 + 3\alpha_1\alpha_1^2 - \alpha_1^3 = \alpha_3 - 3\alpha_2\alpha_1 + 2\alpha_1^3, \\ \mu_4 &= \alpha_4 - 4\alpha_3\alpha_1 + 6\alpha_2\alpha_1^2 - 4\alpha_1\alpha_1^3 + \alpha_1^4 = \alpha_4 - 4\alpha_3\alpha_1 + 6\alpha_2\alpha_1^2 - 3\alpha_1^4.\end{aligned}$$

There are two quantities which one introduces side by side with the third and fourth central moments of a random variable, namely its asymmetry and excess. The asymmetry of a random variable  $X$  is defined by a quotient

$$A = \frac{\mu_3}{\sigma^3} \quad (19)$$

and describes the symmetry or non-symmetry of its distribution law.

The excess of  $X$  is given by a quotient

$$E = \frac{\mu_4}{\sigma^4} - 3 \quad (20)$$

and describes so-called disnormality of  $X$ , that is the deviation of the distribution law of  $X$  from the normal distribution (see the next lecture).

#### Exercise Set 4.

Find the distribution law. Plot the graph of function of distribution. Find number characteristics of a random variable.

1. In the lot of 6 components 4 are standard. 2 components are selected at random. Compose the law of distribution of the random variable  $X$ , which is the number of standard parts among those selected.

2. Probabilities of hitting the target of the first, second and third shooters are respectively equal to 0,4; 0,3 and 0,6. The random variable  $X$  is the number of shots on the target. Find the distribution law of the random variable  $X$ .

3. The probability to hitting the target with one shot is equal to 0,6. The random variable  $X$  is the number of hitting the target with 5 shots. Find the distribution law of the random variable  $X$ ; find  $M(X)$ ,  $D(X)$ ,  $\sigma(X)$ .

Let's suppose that  $n$  independent trials on a random variable  $X$  are fulfilled, and the obtained results are represented by the next tables:

4.	X	-2	0	1	5
	P	0,5	0,2	0,1	0,2
5.	X	-3	-1	2	4
	P	0,3	0,2	0,4	0,1
6.	X	-1	2	3	5
	P	0,1	0,3	0,4	0,2

Build the plotted function of distribution. Find number characteristics of a random variable. Let there be given functions.

$$7. \quad f(x) = \begin{cases} 0, & x \leq 0, \\ A(x+2), & 0 < x \leq 2, \\ 0, & x > 2. \end{cases}$$

$$8. \quad f(x) = \begin{cases} 0, & x \leq 2, \\ A(x-2), & 2 < x \leq 4, \\ 0, & x > 4. \end{cases}$$

$$9. \quad f(x) = \begin{cases} 0, & x \leq 0, \\ Ax^2, & 0 < x \leq 2, \\ 0, & x > 2. \end{cases}$$

$$10. \quad f(x) = \begin{cases} 0, & x \leq 0, \\ Ax, & 0 < x \leq 1, \\ 0, & x > 1. \end{cases}$$

$$11. \quad f(x) = \begin{cases} 0, & x \leq -1, \\ A(x+1), & -1 < x \leq 2, \\ 0, & x > 2. \end{cases}$$

$$12. \quad f(x) = \begin{cases} 0, & x \leq 1, \\ A(x-1), & 1 < x \leq 2, \\ 0, & x > 2. \end{cases}$$

Find the value of the parameter  $A$  so that the function can be the distribution density of some continuous random variable. Find its distribution function. Find number characteristics of the random variable. Plot the graph of function of distribution.

## 5. SOME REMARKABLE DISTRIBUTIONS

### 5.1. THE UNIFORM DISTRIBUTION

**Def. 1.** One says that a random variable  $X$  has a uniform distribution over an interval  $(a, b)$  ( $X$  is uniformly distributed or simply  $X$  is the uniform distribution over an interval  $(a, b)$ ) if its distribution density is constant inside and equals zero outside this interval,

$$f(x) = \begin{cases} c, & \text{if } x \in [a; b], \\ 0, & \text{if } x \notin [a; b]. \end{cases} \quad (1)$$

We have to find the value of the constant  $C$ , the distribution function and number characteristics of the uniform distribution.

A. Finding the value of  $C$ .

On the base of property 4 of the distribution density we must have

$$\int_{-\infty}^{\infty} f(x) dx = 1 \Rightarrow \int_{-\infty}^{\infty} c dx = \int_a^b c dx = c(b-a) = 1 \Rightarrow c = \frac{1}{b-a}$$

and so the distribution density of the uniform distribution is the next one:

$$f(x) = \begin{cases} \frac{1}{b-a}, & \text{if } x \in [a; b], \\ 0, & \text{if } x \notin [a; b]. \end{cases} \quad (2)$$

B. Finding the distribution function of the uniform distribution.

Using property 3 of the distribution density, we must study three cases.

$$F(x) = \begin{cases} 0, & \text{if } x \leq a, \\ \frac{x-a}{b-a}, & \text{if } a < x \leq b, \\ 1, & \text{if } b < x < +\infty. \end{cases} \quad (3)$$

C. Finding the number characteristics of the uniform distribution.

We'll limit ourselves to the mathematical expectation, dispersion and root-mean-square deviation. For this purpose we'll make use of the formulas (11), (14), (15).

$$M(X) = \frac{a+b}{2}, \quad D(X) = \frac{(b-a)^2}{12}, \quad \sigma(X) = \frac{b-a}{\sqrt{12}}, \quad P(c < X < d) = \frac{d-c}{b-a} \quad (4)$$

Ex. 1. The time interval of the trolleybus service equals 5 minutes. Find the probability that one will wait a trolleybus no longer than 2 minutes.

Solution. The waiting time  $T$  is a random variable uniformly distributed over the interval  $(0,5)$ , and we have to find the probability  $P(0 < X < 2)$

$$P(0,2) = \frac{2-0}{5-0} = 0.4$$

## 5.2. THE NORMAL DISTRIBUTION

**Def. 2.** One says that a random variable  $X$  has a normal distribution with parameters  $a$ ,  $\sigma$  ( $\sigma > 0$ ) (or that  $X$  is distributed  $N(a, \sigma)$ ) if its distribution density is the next function:

$$p(x; a, \sigma) = f(x) = \frac{1}{\sigma\sqrt{2\pi}} e^{-\frac{(x-a)^2}{2\sigma^2}} \quad (5)$$

The graph of the function  $f(x)$ , i.e. that of the distribution density of the normal distribution, is represented on figure 5.

**Def. 3.** The graph of the distribution density of the normal distribution is called a normal curve. The normal curve has another and a very fine name, namely the bell-like [or the bell-shaped] curve.

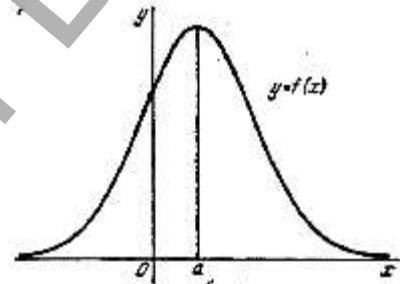


Fig. 5

The normal distribution is often called Gauss distribution, and the corresponding normal curve is called Gauss curve.

Let's consider two important facts connected with the distribution density  $f(x)$  of the normal distribution.

Let the parameter  $\sigma$  of the normal distribution tend to 0. For  $\sigma \rightarrow 0$  the normal curve stretches along the straight line  $x = a$  and simultaneously presses to the  $Ox$  axis. (See figure 6)

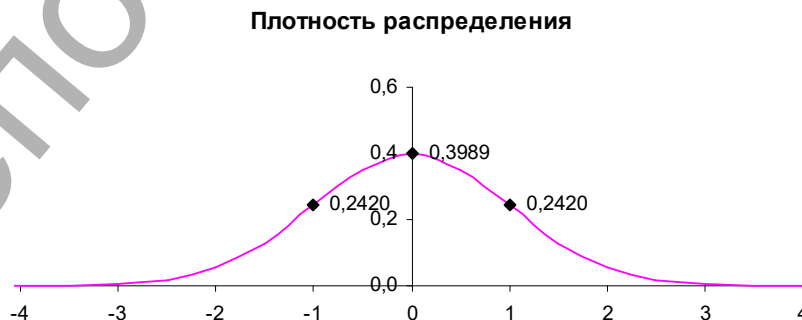


Fig. 6

Let a random variable  $X$  be normally distributed with parameters  $a$ ,  $\sigma$ ,  $\sigma > 0$  ( $X$  distributed  $N(a, \sigma)$ ). We assert that its number characteristics, namely the mathematical expectation, dispersion, root-mean-square deviation are represented by the next formulas:

$$a = M(X), \quad \sigma^2 = D(X), \quad \sigma(X) = \sigma \quad (6)$$

Let a random variable  $X$  be distributed  $N(a, \sigma)$ . The probability of its hitting on an interval  $(\alpha, \beta)$  can be calculated by the next formula:

$$P(\alpha < X < \beta) = \Phi\left(\frac{\beta - a}{\sigma}\right) - \Phi\left(\frac{\alpha - a}{\sigma}\right) \quad (7)$$

where  $\Phi(x) = \frac{1}{\sqrt{2\pi}} \int_0^x e^{-t^2/2} dt$  is known as Laplace function.

The probability of the deviation of a random variable  $X$ , which is distributed  $N(a, \sigma)$ , from its mathematical expectation  $a$  is given by the next formula:

$$P(|X - M(X)| < \delta) = P(|X - a| < \delta) = 2\Phi\left(\frac{\delta}{\sigma}\right) \quad (8)$$

For example, let  $\delta = 3\sigma$ . The formula (8) gives  $P(|X - a| < 3\sigma) = 2\Phi(3) = 0,9973$ .

We've got so-called **3 $\sigma$  - rule**: with a very large probability 0.9973 all values of the normal distribution are concentrated in the interval  $(a - 3\sigma; a + 3\sigma)$ .

Ex.2. A plant makes balls for the bearings. The nominal diameter of the balls is equal to 6 (mm). As a result of an inaccuracy in the production of the balls its actual diameter is a random variable, distributed according to the normal law with an average value of 6 (mm) and mean-square deviation of 0,04 (mm). The balls, whose diameter varies from the nominal by more than 0,1 (mm), are inspected out. Find: 1) what percentage of balls will be rejected on average; 2) probability that the actual diameter of balls will be contained in the range from 5,97 to 6,05 (mm).

Solution. Let the random variable  $X$  be an actual diameter of a ball. It is distributed according to the normal law, i.e.  $X \in N(a; \sigma)$ . If  $a = d_0 = 6$  and  $\sigma = 0,04$ , then  $X \in N(6; 0,04)$ .

Since according to the condition of the task the balls whose diameter differs from the nominal by more than 0,1 (mm) are inspected out, then let us examine the event  $|X - 6| > 0,1$ . To find the probability of this event we will use the opposite event  $|X - 6| \leq 0,1$ . Since random variable  $X$  is continuous, then

$$P(|X - 6| \leq 0,1) = P(|X - 6| < 0,1)$$

$$P(|X - a| < \varepsilon) = 2\Phi\left(\frac{\varepsilon}{\sigma}\right) \Rightarrow P(|X - 6| < 0,1) = 2\Phi\left(\frac{0,1}{0,04}\right) = 2\Phi(2,5)$$

Using the table of the values of the function of Laplace, let us find that  $\Phi(2,5) \approx 0,4938$

$$P(|X - 6| < 0,1) \approx 2 \cdot 0,4938 = 0,9876$$

$$\text{If } P(|X - 6| < 0,1) + P(|X - 6| > 0,1) = 1,$$

$$\text{then } P(|X - 6| > 0,1) = 1 - P(|X - 6| < 0,1) = 1 - 0,9876 = 0,0124.$$

Consequently, 1,24% of the balls will be rejected on average.

The probability of its hitting on an interval  $(\alpha, \beta)$  can be calculated by the next formula:

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

If  $X \in N(6; 0,04)$  and  $\alpha = 5,97$ ,  $\beta = 6,05$ ,



$$\text{then } P(5,97 < X < 6,05) = \Phi\left(\frac{6,05 - 6}{0,04}\right) - \Phi\left(\frac{5,97 - 6}{0,04}\right) = \Phi(1,25) + \Phi(0,75)$$

Using the table of the values of the function of Laplace, let us find that  $\Phi(1,25) \approx 0,3944$  and  $\Phi(0,75) \approx 0,2734$ .

$$P(5,97 < X < 6,05) \approx 0,3944 + 0,2734 = 0,6678$$

### 5.3. THE EXPONENTIAL DISTRIBUTION

We often deal with a call flow [a flow of calls] in a queuing system. Let's denote by  $\lambda$  an intensity of the flow, that is the number of calls which take place (on average) per unit of time. Let  $X$  be a number of calls during time  $t$ . There are many flows (so-called poissonian flows) for which  $X$  has Poisson distribution with the parameter  $a = \lambda t$ . In particular, the probability that  $X$  will take on a value  $m$  equals

$$P(X = m) = P_n(m) = \frac{(\lambda t)^m}{m!} e^{-\lambda t}$$

**Def. 4.** Let a random variable  $X$  be the time interval between two successive calls of some poissonian call flow. One says that  $X$  has the exponential distribution.

Our task is to find the distribution function and density and number characteristics of the exponential distribution.

For positive values of  $t$  the events  $(X < x)$  and  $X \geq 1$  coincide. They mean that during time  $x$  at least one call will occur. Hence, the distribution function of the random variable  $X$  for  $x > 0$  equals

$$F(x) = 1 - e^{-\lambda x} \quad (9)$$

The expression  $1 - e^{-\lambda x}$  tends to zero with  $x$ , and therefore we can define the distribution function in question as follows

$$F(x) = \begin{cases} 1 - e^{-\lambda x}, & \text{if } x \geq 0, \\ 0, & \text{if } x < 0. \end{cases} \quad (10)$$

It is continuous for all values of  $x$ , and therefore a random variable  $X$  which has the exponential distribution is a continuous one.

Differentiating the distribution function (10) we'll obtain the distribution density of the exponential distribution,

$$f(x) = \begin{cases} \lambda e^{-\lambda x}, & \text{if } x \geq 0, \\ 0, & \text{if } x < 0. \end{cases} \quad (11)$$

The graphs of both the functions  $F(x)$ ,  $f(x)$  are represented on fig. 7, 8.

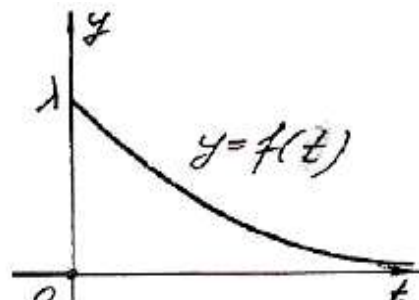
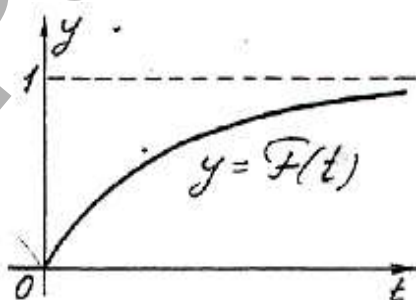


Fig. 7, 8

We assert that its number characteristics, namely the mathematical expectation, dispersion, root-mean-square deviation are represented by the next formulas:

$$\lambda > 0, \quad M(X) = \frac{1}{\lambda}, \quad D(X) = \frac{1}{\lambda^2}, \quad \sigma(X) = \frac{1}{\lambda} \quad (12)$$

The probability of its hitting on an interval  $(a, b)$  can be calculated by the next formula:

$$P(a < X < b) = e^{-\lambda a} - e^{-\lambda b} \quad (13)$$

Ex.3. Time  $t$  of the reliable work of radio-technical system is distributed according to the exponential law. Failure rate of the system is  $\lambda = 0,02$ . Find the mean time of failure-free operation and the probability of failure-free operation in 80 hours.

Solution. The density of probability distribution of this distribution takes the form

$$f(t) = \begin{cases} 0,02 e^{-0,02t}, & \text{if } t \geq 0, \\ 0, & \text{if } t < 0. \end{cases}$$

Mathematical expectation is this mean time of the reliable work of system.

$$M(T) = \frac{1}{\lambda} = \frac{1}{0,02} = \frac{100}{2} = 50 \text{ (hours)}.$$

Let us determine the probability of failure-free operation for 80 hours with the aid of the function of the reliability  $R(t) = e^{-\lambda t} = P(0 < X < t)$

$$R(80) = e^{-0,02 \cdot 80} = e^{-1,6} = 0,2019.$$

#### 5.4. BERNOULLI [BINOMIAL] DISTRIBUTION

**Def.5.** Let a random variable  $X$  be the number of successes (the number of occurrences of some event  $A$ ) in  $n$  independent trials with constant probability of the success  $A$  in any trial

$$p = P(A), \quad P(\bar{A}) = 1 - p = q$$

One says that  $X$  is distributed binomially (by Bernoulli [binomial] law) or simply:  $X$  is Bernoulli (binomial) distribution (briefly:  $X$  **distr. B**).

Let's find a probability  $P(X = m) = P_n(m)$ , that is the probability of  $m$  successes.

We'll get so-called Bernoulli formula:

$$P_n(m) = C_n^m p^m q^{n-m}, \quad q = 1 - p \quad (14)$$

The mathematical expectation, dispersion, root-mean-square deviation are represented by the next formulas:

$$M(X) = n p, \quad D(X) = n p q, \quad \sigma(X) = \sqrt{n p q} \quad (15)$$

Ex. 4. 6 independently working engines are installed in a shop. Probability for any engine to work at a given moment is 0.8. Find number characteristics of a random variable  $X$ , if the random variable  $X$  is a number of working engines at this moment.

Solution. We can consider setting of an engine as a trial. So we have  $n = 6$  independent trials. Let a success  $A$  mean that an engine works.

$$p = P(A) = 0.8, \quad P(\bar{A}) = 1 - p = q = 0.2$$

The random variable  $X$  has Bernoulli distribution (briefly " $X$  distributed  $B$ "), it can take on the values 0, 1, 2, 3, 4, 5, 6, which one calculates by Bernoulli formula (14). For example

$$P(X = 0) = P_6(0) = C_6^0 p^0 q^6 = 0.2^6 = 0.00006$$

$$P(X = 1) = P_6(1) = C_6^1 p^1 q^5 = 6 \cdot 0.8 \cdot 0.2^5 = 0.00154$$

Using formulas (15), we get

$$M(X) = np = 6 \cdot 0.8 = 4.8, \quad D(X) = npq = 6 \cdot 0.8 \cdot 0.2 = 0.96,$$

$$\sigma(X) = \sqrt{npq} = \sqrt{0.96} = 0.979$$

### 5.5. POISSON FORMULA AND DISTRIBUTION

Let a random variable  $X$  be distributed  $B$ . Let's suppose that the number  $n$  of trials tends to infinity, the probability  $p$  of a success  $A$  goes to zero, but a product  $np$  retains constant,

$$n \rightarrow \infty, p \rightarrow 0, np = \text{const} = \lambda.$$

In this case the limit of the probability  $P(X = m) = P_n(m)$ , which is defined by Bernoulli formula (1), equals

$$P_n(m) \approx \frac{\lambda^m}{m!} e^{-\lambda} \quad (16)$$

**Def. 6.** One says that a discrete random variable  $X$  (with non-negative integer possible values) has Poisson distribution with a parameter  $\lambda$  if its distribution law is given by the next formula (Poisson formula):

$$P(X = m) = P_n(m) = \frac{\lambda^m}{m!} e^{-\lambda} \quad (17)$$

The mathematical expectation, dispersion, root-mean-square deviation are represented by the next formulas:

$$M(X) = D(X) = \lambda, \sigma(x) = \sqrt{\lambda} \quad (18)$$

#### Exercise Set 5.

Random variable  $X$  is normally distributed with the mathematical expectation  $M(X)$  and the dispersion  $D(X)$ . Find the density of probability distribution  $f(x)$  and build the schematic graph of this function. Write down the interval of practically probable values of a random variable. Which is more probable  $X \in (\alpha; \beta)$  or  $X \in (\gamma; \delta)$ ?

Exercise	$M(X)$	$D(X)$	$\alpha$	$\beta$	$\gamma$	$\delta$
1	2	4	-1	4	5	6
2	-3	9	-5	-4	-2	-1
3	4	1	0	3	2	5
4	-5	16	-10	-6	0	4
5	3	4	0	5	6	7
6	-2	9	-6	-5	-4	0
7	5	1	0	3	2	6
8	-1	25	-3	1	0	4
9	1	9	-2	2	-1	5
10	-3	4	-4	0	1	6

The mean time of operation of each of the three elements, entering the technical device, is equal to  $T$  hours. For the reliable work of the device the failure-free operation of at least one of these three elements is necessary. Find the probability that the device will work from  $t_1$  to  $t_2$  hours, if the time of operation of each of the three elements independently is distributed according to the exponential law.

Exercise	T	$t_1$	$t_2$
11	800	650	700
12	1000	800	900
13	850	750	820
14	1200	900	1000
15	900	700	900
16	950	720	850

17. The probability of hitting the target with one shot is equal to 0,6. The random variable  $X$  is a number of striking the target with 5 shots. Find number characteristics of a random variable.

18. Someone expects a telephone call between 19.00 and 20.00. The waiting time of the ring is the random variable  $X$ , which has uniform distribution in section  $[19; 20]$ . Find probability that the telephone will ring in the period from 19 hours 22 minutes to 19 hours 46 minutes.

19. Radio equipment for 1000 hours of work goes out of order on the average one time. Find the probability of failure of the radio equipment for 200 hours of work, if the period of failure-free operation is a random variable, distributed according to the exponential law.

20. A basketball player makes three penalty throws. The probability of hit with each throw is equal to 0,7. Find number characteristics of a random variable  $X$ , if the random variable  $X$  is a number of shots at basket.

## 6. ELEMENTS OF MATHEMATICAL STATISTICS

### 6.1. GENERAL REMARKS. SAMPLING METHOD. VARIATION SERIES

Let's suppose that we study some random variable  $X$ .

We'll dwell upon three typical problems of the mathematical statistics.

1. Exact or approximate determination of the distribution law of a random variable (for example it can be stated or hypothesized that a random variable  $X$  is distributed normally).

2. Estimation (approximate calculation) of parameters of the distribution law of a random variable (for example estimation of  $M(X)$ ;  $\sigma(X)$ ;  $D(X)$ ;  $As(X)$  of a random variable  $X$ ).

3. Testing statistical hypotheses (for example testing a hypothesis that a given random variable  $X$  is distributed normally).

There are various methods of solving such the problems. One of the most widespread is the sampling method. Suppose that we have some population consisting a great number  $N$  of things (the population of the size  $N$ ) which must be studied with respect to some random variable  $X$ . We take at random things  $n \leq N$  from the population, fulfil their allround testing with respect to  $X$  and extend obtained results on the whole population.

On the language of the mathematical statistics we do a sampling of the size  $n$  ( $n \ll N$ ) getting the sample (of the size  $n$ ) which is subjected to thorough investigation with respect to a random variable in question. A sample must be representative, that is it must certainly represent the population. To be representative the sample must be random one.

A sample of the size  $n$ , which we obtain by a random sampling from the population, we study with the help of so-called variation (or statistical) series. There are variation series of two types namely those discrete and interval.

Table 1. A discrete variation series

$X, x_i$	$x_1$	$x_2$	...	$x_k$
$m_i$	$m_1$	$m_2$	...	$m_k$
$p_i^* = \frac{m_i}{n}$	$p_1^* = \frac{m_1}{n}$	$p_2^* = \frac{m_2}{n}$	...	$p_k^* = \frac{m_k}{n}$

1. A discrete variation series contains the row of observed values  $x_i$  of a random variable  $X$  to be the investigated (as the rule in increasing order), the row of numbers  $m_i$  of occurrences of these values and the row of their relative frequencies  $p_i^* = \frac{m_i}{n}$  (table 1).

It must be  $\sum_{i=1}^k m_i = n$  and  $\sum_{i=1}^k p_i^* = 1$ .

Such a discrete variation series can be represented geometrically with the help of a polygon of frequencies or a polygon of relative frequencies. The polygon of relative frequencies is a broken line which joins successively the next points:

$$A_1(x_1, p_1^*), A_2(x_2, p_2^*), \dots, A_n(x_n, p_n^*) \text{ (see fig. 9).}$$

The polygon of frequencies is a broken line joining successively the other points namely those with ordinates (frequencies)  $m_1, m_2, \dots, m_k$

$$B_1(x_1, p_1), B_2(x_2, p_2), \dots, B_n(x_n, p_n).$$

If there are a lot of distinct observed values of a random variable  $X$ , they can be united in intervals which generate so-called interval variation series.

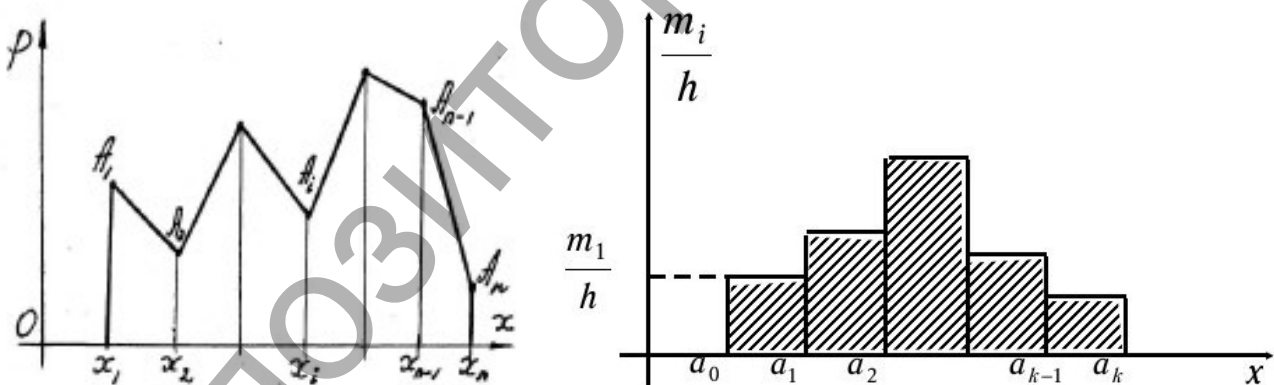


Fig. 9, 10

2. An interval variation series contains the row of intervals with all observed values of a random variable  $X$  in question, the row of numbers (frequencies)  $m_i$  of hitting of these values in the  $i$ -th interval and the row of corresponding relative frequencies  $p_i^* = \frac{m_i}{n}$  (table 2).

Table 2. An interval variation series

Intervals	$a_0 - a_1$	$a_1 - a_2$	...	$a_{k-1} - a_k$
Frequencies $m_i$	$m_1$	$m_2$	...	$m_k$
$p_i^* = \frac{m_i}{n}$	$p_1^* = \frac{m_1}{n}$	$p_2^* = \frac{m_2}{n}$		$p_k^* = \frac{m_k}{n}$

Just as in the case of a discrete variation series it must be

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

The interval variation series can be represented geometrically by a histogram of frequencies or relative frequencies. The histogram of relative frequencies is the set of rectangles with the bases  $(a_{k-1}; a_k]$ , of the lengths  $\Delta x_i = a_{i+1} - a_i$  and the areas  $S_i = p_i^*$  (fig. 10). The altitude of the  $i$ -th rectangle of such a histogram equals  $m_k/h$ .

The histogram of frequencies contains rectangles of the areas  $S_i = m_i$  with the same bases.

Often we form intervals of the same length  $\Delta x$  (in particular on fig. 10).

If one wants to generate  $k$  intervals, he can put

$$h = \frac{x_{\max} - x_{\min}}{k} \quad (1)$$

where  $x_{\max}$  is the greatest and  $x_{\min}$  is the least of observed values of a random variable  $X$ . In practice it's useful sometimes to take instead  $x_{\max}$  some greater number and instead  $x_{\min}$  some less number in the formula (1).

One can preset the length of intervals instead their number. For this purpose he can use the next known approximate formula:

$$h \cong \frac{x_{\max} - x_{\min}}{1 + 3,322 \cdot \lg n} \quad (2)$$

As the left point of the first interval he can take

$$a_0 = x_{\min} - \frac{h}{2} \quad (3)$$

and get then the other points as follows

$$a_i = x_{i-1} + h \quad (4)$$

The number  $k$  of intervals is determined by the condition

$$a_k \geq x_{\max} \quad (5)$$

which means that the last  $k$ -th interval must contain the value  $x_{\max}$  of the random variable  $X$  in question.

Having an interval variation series we often must compile a corresponding discrete variation series by taking some inner point  $x_i$  in the  $i$ th interval. For example we can take

$$x_i = \frac{a_i + a_{i+1}}{2} \quad (6)$$

and obtain the discrete variation series which is given by the table 3. We suppose conditionally that the point  $x_i$  represents the  $i$ -th interval and therefore that the value  $X = x_i$  were observed  $m_i$  times (with the relative frequency  $p_i^* = \frac{m_i}{n}$ ). It means that the tables 2 and 3 contain the same second and third rows.

Table 3. The discrete variation series which corresponds to that interval

$x_i$	$x_1$	$x_2$	...	$x_k$
$m_i$	$m_1$	$m_2$	...	$m_k$
$p_i^* = \frac{m_i}{n}$	$p_1^* = \frac{m_1}{n}$	$p_2^* = \frac{m_2}{n}$	...	$p_k^* = \frac{m_k}{n}$

Ex. 1 (the **Basic example**). To study a random variable  $X$  a sampling is fulfilled and the sample of the size  $n = 100$  is obtained (the table 4). Study the random variable.

Table 4. The sample with respect to the random variable  $X$  of the **Basic example**

24,8	26,2	25,6	24,0	26,4	25,2	26,7	25,4	25,3	26,1
24,3	25,3	25,6	26,7	24,5	26,0	25,7	25,0	26,4	25,9
24,4	25,4	26,1	23,4	26,5	25,9	23,9	25,7	27,1	24,9
23,8	25,6	25,2	26,4	24,2	26,5	25,7	24,7	26,0	25,8
24,3	25,5	26,7	24,9	26,2	26,7	24,6	26,0	25,4	25,0
25,4	25,3	24,1	26,6	24,8	25,6	23,7	26,8	25,2	26,1
24,5	25,4	25,1	26,2	24,2	26,4	25,7	23,9	27,2	25,0
23,9	25,6	24,9	24,5	26,2	26,7	24,3	26,1	27,7	25,8
25,6	25,2	24,2	26,0	24,7	26,5	23,5	25,4	27,1	24,0
26,2	24,2	25,5	26,0	25,7	26,4	24,6	27,0	25,2	26,9

$$x_{\max} - x_{\min} = 27,7 - 23,4 = 4,3$$

We have  $x_{\min} = 23,4$  and  $x_{\max} = 27,7$  here. Let's compile intervals of the length

$$h = \frac{x_{\max} - x_{\min}}{1 + 3,322 \cdot \lg n} = \frac{4,3}{1 + 3,322 \cdot \lg 100} = \frac{4,3}{7,644} \approx 0,56$$

$$a_1 = x_{\min} - \frac{h}{2} = 23,4 - 0,28 = 23,12$$

$$a_1 = 23,12; \quad a_2 = 23,68; \quad a_3 = 24,24; \quad a_4 = 24,8; \quad a_5 = 25,36;$$

$$a_6 = 25,92; \quad a_7 = 26,48; \quad a_8 = 27,04; \quad a_9 = 27,6; \quad a_{10} = 28,16.$$

Table 5 represents the interval variation series, inner points  $x_i$  for corresponding discrete variation series and altitudes  $h_i$  of rectangles for plotting the histogram of relative frequencies.

Table 5. The variation series for the **Basic example**

Intervals	$x_k$	Frequencies $m_k$	$p_i^* = \frac{m_i}{n}$	Altitudes $h_i = \frac{p_i^*}{\Delta x}$
23,12-23,68	23,40	2	0,02	0,04
23,68-24,24	23,96	11	0,11	0,20
24,24-24,80	24,52	14	0,14	0,25
24,80-25,36	25,08	14	0,14	0,25
25,36-25,92	25,64	23	0,23	0,41
25,92-26,48	26,20	20	0,20	0,36
26,48-27,04	26,76	12	0,12	0,21
27,04-27,60	27,32	3	0,03	0,05
27,60-28,16	27,88	1	0,01	0,02
$\Sigma$		100	1,00	

The histogram of relative frequencies for the interval and corresponding discrete variation series are represented on fig. 11. For the sake of geometric visualization we draw to different scales along the axes.

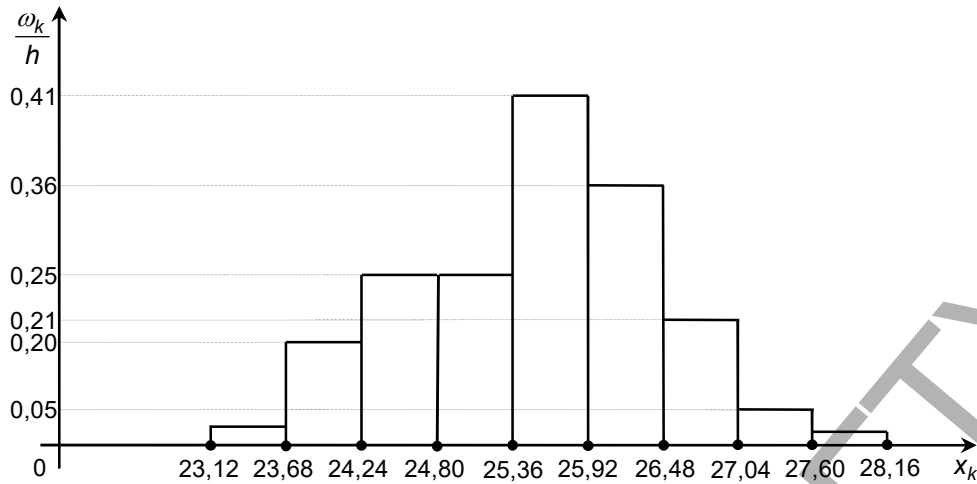


Fig. 11

## 6.2 APPROXIMATE DETERMINATION OF THE DISTRIBUTION LAW OF A RANDOM VARIABLE. ESTIMATION OF PARAMETERS OF THE DISTRIBUTION LAW OF A RANDOM VARIABLE

### The statistic distribution function

Def. 1. The statistic distribution function of a random variable  $X$  is called the relative frequency of hitting of its observed values in the interval  $(-\infty, x)$ ,

$$F^*(x) = W(X < x) = \frac{n_x}{n} \quad (7)$$

It's obvious that the statistic distribution function  $F^*(x)$  equals the sum of relative frequencies of those observed values of  $X$  which are less than  $x$ ,

$$F^*(x) = \sum_{x_j < x} p_j^* \quad (8)$$

To find the statistic distribution function we must distinguish the cases of discrete and interval variation series.

Ex. 2. Find the statistic distribution function of the random variable  $X$  which is studied in the **Basic example**.

Solution. Proceeding from the discrete variation series (see inner points  $x_j$  and corresponding relative frequencies  $p_j^* = \frac{m_j}{n}$  in the table 5) we obtain

$$F^*(x) = \begin{cases} 0, & x \leq 23,12; \\ 0,02, & 23,12 < x \leq 23,68; \\ 0,13, & 23,68 < x \leq 24,24; \\ 0,27, & 24,24 < x \leq 24,80; \\ 0,41, & 24,80 < x \leq 25,36; \\ 0,64, & 25,36 < x \leq 25,92; \\ 0,84, & 25,92 < x \leq 26,48; \\ 0,96, & 26,48 < x \leq 27,04; \\ 0,99, & 27,04 < x \leq 27,60; \\ 1, & 27,60 < x \leq 28,16; \\ 1, & x > 28,16. \end{cases}$$



Using the interval variation series we get approximate values of the statistic distribution function represent corresponding points on the  $xOy$ -plane and plot the approximate graph of the statistic distribution function in the form of a continuous line (fig. 12).

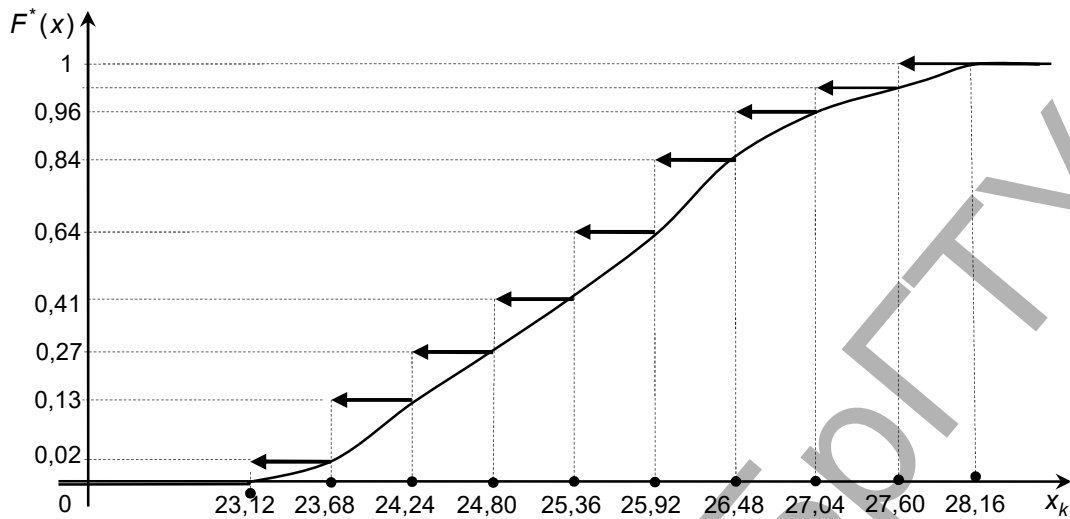


Fig. 12

Let  $\theta$  be a parameter to be estimated and  $\tilde{\theta}$  be its estimate [or estimator].

This latter is a function of results of  $n$  trials on  $X$ .

In theory we consider results of trials as random variables  $X_1, X_2, \dots, X_n$  with the same distribution law as  $X$ , in particular with the same mathematical expectation and dispersion

$$M(X_i) = M(X), D(X_i) = D(X), i = \overline{1, n}.$$

Respectively we suppose the estimate  $\tilde{\theta}$  to be a function of these random variables:

$$\tilde{\theta} = f(X_1, X_2, \dots, X_n).$$

In practice we express the estimate  $\tilde{\theta}$  with the help of results of trials on the random variable  $X$  represented by variation series.

For discrete variation series (with  $k$  observed various values  $x_i$  of the random variable  $X$ ) we can write

$$\tilde{\theta} = f(x_1, x_2, \dots, x_n)$$

For interval variation series (with  $k$  intervals) we use the inner points  $x_i^*$  of the intervals and so

$$\tilde{\theta} = f(x_1^*, x_2^*, \dots, x_n^*)$$

There exist three necessary requirements which must be laid to every estimate: consistency, unbiasedness and efficiency [effectiveness].

1. An estimate  $\tilde{\theta}$  is called a consistent estimator of the parameter  $\theta$  if it converges to  $\theta$  in probability ( $n \rightarrow \infty \quad P(|\tilde{\theta} - \theta| < \varepsilon) \rightarrow 1$ )

2. This estimate  $\tilde{\theta}$  is called an unbiased estimator of the parameter  $\theta$  if its mathematical expectation equals  $\theta$  ( $M(\tilde{\theta}) = \theta$ ).

3. The estimate  $\tilde{\theta}$  is called an efficient [effective] estimator of  $\theta$  if it has minimal dispersion in comparison with all other estimates.

There are pointwise and interval estimates of the parameters of distribution laws of random variables.

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

$$\sqrt{D(X)} = \sigma(X), A = \frac{\mu_3}{\sigma^3}, E = \frac{\mu_4}{\sigma^4} - 3$$

In mathematical statistics we try at first to estimate corresponding parameters by analogous formulas.

1. Estimate [estimator] of the mathematical expectation [estimation of expectation].

$$M(\bar{x}_s) = \bar{x}_g \quad (9)$$

$$\bar{x}_s = \bar{x} = \frac{1}{n} \sum_{k=1}^m n_k x_k \quad (10)$$

2. Estimate [estimator] of the dispersion [variance estimate].

$$D_s = \frac{1}{n} \sum_{k=1}^m n_k (x_k - \bar{x})^2 = \frac{1}{n} \sum_{k=1}^m n_k x_k^2 - (\bar{x})^2 = \overline{x^2} - (\bar{x})^2 \quad (11)$$

Therefore the sample dispersion  $D_s$  is consistent but biased [shifted] estimate of the dispersion. To have an unbiased consistent estimator one introduces so-called **corrected dispersion**

$$S^2 = \frac{n}{n-1} \cdot D_s, \quad M(S^2) = D_g \quad (12)$$

3. The root-mean-square deviation of the random variable  $X$  is estimated by the **sample root-mean-square deviation**

$$\sigma_s = \sqrt{D_s} \quad (13)$$

and the **corrected root-mean-square deviation** (or the **standard**)

$$s = \sqrt{S^2} \quad (14)$$

Ex. 3. Estimate the mathematical expectation, dispersion, root-mean-square deviation of the random variable  $X$  of the **Basic example**.

Corresponding calculations are represented in table 6.

Table 6.

$(a_k; a_{k+1})$	$x_k$	$m_k$	$x_k m_k$	$x_k^2 m_k$
23,12-23,68	23,40	2	46,8	1095,12
23,68-24,24	23,96	11	263,56	6314,898
24,24-24,80	24,52	14	343,28	8417,226
24,80-25,36	25,08	14	351,12	8806,09
25,36-25,92	25,64	23	589,72	15120,42
25,92-26,48	26,20	20	524	13728,8
26,48-27,04	26,76	12	321,12	8593,171
27,04-27,60	27,32	3	81,96	2239,147
27,60-28,16	27,88	1	27,88	777,2944
$\Sigma$		100	2549,44	65092,17

Thus the sample means of the random variable  $X$  and its square are equal to

$$\bar{x}_s = \frac{2549,44}{100} = 25,4944 \approx 25,49; \quad \overline{x^2} = 650,9217$$

and therefore by (11) and (13) the sample dispersion and root-mean-square deviation are equal to

$$D_s = 650,9217 - 25,4944^2 = 650,9217 - 649,9644 \approx 0,9573;$$

$$\sigma_s = \sqrt{0,9573} \approx 0,9784 \approx 0,98$$

The corrected dispersion and root-mean-square deviation equal by virtue of the formulas (12), (14)

$$S^2 = \frac{100}{99} \cdot 0,9573 \approx 0,967 \approx 0,97;$$

$$s = \sqrt{0,97} \approx 0,98.$$

### **Interval estimates [estimation by confidence interval]**

Let  $\theta$  be a parameter to be estimated,  $\tilde{\theta}$  its estimate,  $\delta$  some small positive number (so-called **accuracy**) and  $\gamma$  is some large probability (the **reliability**).

The next relation

$$P(|\theta - \tilde{\theta}| < \delta) = \gamma \quad (15)$$

means that with the reliability  $\gamma$

$$\theta \in (\tilde{\theta} - \delta, \tilde{\theta} + \delta)$$

therefore it is equivalent to the next one:

$$P(\tilde{\theta} - \delta < \theta < \tilde{\theta} + \delta) = P(|\theta - \tilde{\theta}| < \delta) = \gamma \quad (16).$$

This last equality shows that with the reliability  $\gamma$  the estimating parameter  $\theta$  is covered by a random interval  $(\tilde{\theta} - \delta, \tilde{\theta} + \delta)$ .

**Def.** An interval  $(\tilde{\theta} - \delta, \tilde{\theta} + \delta)$  which covers an estimated parameter  $\theta$  with the reliability  $\gamma$  is called **confidence** one.

Let a random variable  $X$  be normally distributed, and it's necessary to find the confidence interval for its mathematical expectation  $M(X) = a$ .

We can suppose that the sample mean  $\bar{x}_s$  and corrected root-mean-square deviation  $\sigma_s$  of  $X$  are found.

Let's study an auxiliary random variable

$$t = \sqrt{n} \frac{\bar{x}_s - a}{\sigma_s} \quad (17)$$

It can be proved that  $t$  has Student distribution (or  $t$ -distribution) with  $k = n - 1$  degrees of freedom. Therefore for the given reliability  $\gamma$  we can find a number  $t_\gamma$  such that

$$P(|t| < t_\gamma) < \gamma \quad (18)$$

There's a corresponding table for finding  $t_\gamma$  by known  $\gamma$  and  $n$  (see table 7).

The formulas (17) and (18) determine the confidence interval in question.

1. If  $\sigma$  is known

$$\bar{x}_s - \frac{t\sigma}{\sqrt{n}} < a < \bar{x}_s + \frac{t\sigma}{\sqrt{n}} \quad (19),$$

where  $2\Phi(t) = \gamma$  or  $P\left(\left|a - \bar{x}_B\right| < \frac{t\sigma}{\sqrt{n}}\right) = \gamma$ .

Table 7.

n \ \gamma		\gamma			n \ \gamma		\gamma		
		0,95	0,99	0,999			0,95	0,99	0,999
5		2,78	4,60	8,61	20	2,093	2,861	3,883	
6		2,57	4,03	6,86	25	2,064	2,797	3,745	
7		2,45	3,71	5,96	30	2,045	2,756	3,659	
8		2,37	3,50	5,41	35	2,032	2,720	3,600	
9		2,31	3,36	5,04	40	2,023	2,708	3,558	
10		2,26	3,25	4,78	45	2,016	2,692	3,527	
11		2,23	3,17	4,59	50	2,009	2,679	3,502	
12		2,20	3,11	4,44	60	2,001	2,662	3,464	
13		2,18	3,06	4,32	70	1,996	2,649	3,439	
14		2,16	3,01	4,22	80	1,991	2,640	3,418	
15		2,15	2,98	4,14	90	1,987	2,633	3,403	
16		2,13	2,95	4,07	100	1,984	2,627	3,392	
17		2,12	2,92	4,02	120	1,980	2,617	3,374	
18		2,11	2,90	3,97	$\infty$	1,960	2,576	3,291	
19		2,10	2,88	3,92					

2. If  $\sigma$  is unknown

$$\bar{x}_s - \frac{t_\gamma \cdot S}{\sqrt{n}} < a < \bar{x}_s + \frac{t_\gamma \cdot S}{\sqrt{n}} \quad (20),$$

where  $S^2 = \frac{n}{n-1} \cdot D_s$ ,  $s = \sqrt{S^2}$  and  $t_\gamma = t(\gamma, n)$ .

3. Confidence interval for root-mean-square deviation  $\sigma$  is

$$S(1-q) < \sigma < S(1+q), \text{ if } q < 1$$

$$0 < \sigma < S(1+q), \text{ if } q > 1$$

where we find the number  $q = q(\gamma, n)$  according to table 8.

Table 8.

n \ \gamma		\gamma		
		0,95	0,99	0,999
5		1,37	2,67	5,64
6		1,09	2,01	3,88
7		0,92	1,62	2,98
8		0,80	1,38	2,42
9		0,71	1,20	2,06
10		0,65	1,08	1,80
11		0,59	0,98	1,60
12		0,55	0,90	1,45
13		0,52	0,83	1,33
14		0,48	0,78	1,23
15		0,46	0,73	1,15
16		0,44	0,70	1,07
17		0,42	0,66	1,01
18		0,40	0,63	0,96
19		0,39	0,60	0,92

n \ \gamma		\gamma		
		0,95	0,99	0,999
20		0,37	0,58	0,88
25		0,32	0,49	0,73
30		0,28	0,43	0,63
35		0,26	0,38	0,56
40		0,24	0,35	0,50
45		0,22	0,32	0,46
50		0,21	0,30	0,43
60		0,188	0,269	0,38
70		0,174	0,245	0,34
80		0,161	0,226	0,31
90		0,151	0,211	0,29
100		0,143	0,198	0,27
150		0,115	0,160	0,211
200		0,099	0,136	0,185
250		0,089	0,120	0,162

Ex. 4. Supposing that the random variable  $X$  of the **Basic example** has a normal distribution find the confidence intervals for its mathematical expectation  $a$  with the reliabilities  $\gamma = 0,95$ .

For the reliability  $\gamma = 0,95$  and  $n = 100$  we have  $t_\gamma = t(\gamma, n) = t(0,95; 100) = 1,984$  table 7, hence the corresponding accuracy equals

$$\delta = \frac{s}{\sqrt{n}} t_\gamma, s = 0,98, \delta = \frac{0,98}{10} \cdot 1,984 \approx 0,1944 \approx 0,19$$

and the confidence interval is

$$(\bar{x} - \delta; \bar{x} + \delta) = (25,49 - 0,19; 25,49 + 0,19) = (25,30; 25,68).$$

The confidence interval for its root-mean-square deviation is calculated by the formula

$$s(1 - q) < \sigma < s(1 + q),$$

$$s - \delta < \sigma < s + \delta.$$

$$q = q(\gamma, n) = q(0,95; 100) = 0,143, \delta = s \cdot q = 0,98 \cdot 0,143 = 0,13299 \approx 0,13.$$

The confidence interval for its root-mean-square deviation is

$$(s - \delta; s + \delta) = (0,98 - 0,13; 0,98 + 0,13) = (0,85; 1,11).$$

### 6.3 TESTING STATISTIC HYPOTHESES

We'll limit ourselves to hypotheses about distribution law of a random variable which is investigated. For instance our **Basic example** has generated the hypothesis that  $X$  has a normal distribution.

Let us test a hypothesis that a random variable  $X$  has certain distribution law. For this purpose we introduce some non-negative random variable (**goodness-of-fit test**)  $K$  which is the measure of deviation of the theoretical assumption (based on the hypothesis) and the results of trials on the random variable. It's supposed that we know exact or approximate distribution law of  $K$ . On the base of the result of trials on  $X$  (for example on the base of the fulfilled sample) we find so-called calculated value  $K_{calc}$  of  $K$  and compare it with some well defined critical value  $K_{crit}$  of the same  $K$ .

Let for example the critical value  $K_{crit}$  of the goodness-of-fit test  $K$  be defined by the relation

$$P(K > K_{crit}) > \alpha$$

where  $\alpha$  is some small probability. This probability is often called the **significance level** in such a sense that the event  $K > K_{crit}$  can be considered as highly improbable or even practically impossible. Correspondingly we consider the occurrence of the result  $K_{calc} > K_{crit}$  as low-probability [unlikely] outcome.

Comparison of the calculated value  $K_{calc}$  of the test  $K$  with its critical value  $K_{crit}$  can give rise to two cases:  $K_{calc} > K_{crit}$  or  $K_{calc} \leq K_{crit}$ .

In the first case we say that the results of trials **contradict the hypothesis** because of occurrence of the low-probability event. Therefore we can reject the hypothesis in question.

In the second case we say that the results of trials **don't contradict the hypothesis**, and we can accept it.

Remark 1. It's necessary to understand that we ascertain contradiction or noncontradiction of the results of trials to the hypothesis but **we don't state its validity or invalidity**.

Henceforth we'll study some frequently used goodness-of-fit tests.

### **Pearson $\chi^2$ -goodness-of-fit test**

Let's introduce the next random variable (so-called Pearson distribution; it's often called Pearson  $\chi^2$ -goodness-of-fit test or simply  $\chi^2$ -goodness-of-fit test)

$$\chi^2_{\text{calc}} = \sum \frac{(m_i - m'_i)^2}{m'_i}, m'_i = n \cdot P_i, i = \overline{1, k} \quad (21)$$

where  $n$  is the number of trials on the random variable  $X$  (for example the size of a sample if the trials consist in sampling); the sense of the other quantities, namely the sense of  $m'_i, P_i, k$  depends on the form of a variation series which represents the results of trials.

a) In the case of a **discrete variation series**  $k$  is the number of different observed values of the random variable  $X$ ,  $m_i$  is the number of occurrences of the value  $x_i$  of  $X$ , and  $P_i$  is the probability of occurrence of this value  $P_i = P(X = x_i)$ . This latter is calculated on the base of the advanced hypothesis. If for example we've hypothesized that a random variable  $X$  has Poisson distribution, then

$$P(X = m) = P_n(m) = \frac{\lambda_s^m}{m!} e^{-\lambda_s} \quad (22)$$

where  $\lambda_s$  is the point estimate of the parameter  $\lambda$  (sample  $\lambda$ ).

b) In the case of an **interval variation series**  $k$  is the number of intervals,  $m_i$  the number of values of  $X$  which have hit in the  $i$ -th interval (that is the frequency of hitting of observed values of  $X$  in this interval), and  $P_i$  is the probability of hitting of  $X$  in this interval,  $P_i = P(a_{i-1} < X < a_i)$ , calculated by virtue of the hypothesis.

For example in the hypothesis of the normal distribution of a random variable  $X$  we can write

$$P_i = P(a_{i-1} < X < a_i) = \Phi\left(\frac{a_i - \bar{x}_s}{\sigma_s}\right) - \Phi\left(\frac{a_{i-1} - \bar{x}_s}{\sigma_s}\right) \quad (23)$$

The law of Pearson distribution (21) is known: it approximately coincides with the  $\chi^2$ -distribution. The number  $\nu$  of degrees of freedom of this distribution is proved to be the next:

$$\nu = k - r - 1 \quad (24)$$

where  $r$  is the number of independent parameters which we estimate as to our random variable  $X$ . Values of the parameter  $r$  for some known distributions are represented in table 8. For known number of degrees of freedom  $\nu$  and a small probability (a significance level)  $\alpha$  there exists the **critical value**  $\chi^2_{\text{crit}}$  of  $\chi^2$  such that

$$P(\chi^2 > \chi^2_{\text{crit}}) = \alpha \quad (25)$$

It can be found with the help of a corresponding table if we preset ourselves the significance level (see in table 9 near).

	Hypothesis	Parameters to be estimated	The corresponding value of the parameter
1	$X$ has a normal distribution	The mathematical expectation and dispersion $a, \sigma$	2
2	$X$ has Bernoulli distribution	The mathematical expectation	1

3	$X$ has Poisson distribution with the parameter $a$	The mathematical expectation $M(X) = a$	1
4	$X$ has an exponential distribution with the parameter $\lambda$	The mathematical expectation $M(X) = \sigma(X) = \frac{1}{\lambda}$	1
5	$X$ has a uniform distribution on an interval with endpoints $a$ and $b$	There aren't such the parameters	0

Table 9. The critical values of the  $\chi^2$ - distribution

$\alpha$ $\nu$	0,20	0,10	0,05	0,02	0,01	0,001
1	1,642	2,706	3,841	5,412	6,635	10,827
2	3,219	4,605	5,991	7,824	9,210	13,815
3	4,642	6,251	7,815	9,837	11,345	16,266
4	5,989	7,779	9,488	11,668	13,237	18,467
5	7,289	9,236	11,070	13,388	15,086	20,515
6	8,558	10,645	12,592	15,033	16,812	22,457
7	9,803	12,017	14,067	16,622	18,475	24,322
8	11,030	13,362	15,507	18,168	20,090	26,125
9	12,242	14,684	16,919	19,679	21,666	27,877
10	13,442	15,987	18,307	21,161	23,209	29,588
11	14,631	17,275	19,675	22,618	24,795	31,264
12	15,812	18,549	21,026	24,054	26,217	32,909
13	16,985	19,812	22,362	25,472	27,688	34,528
14	18,151	21,064	23,685	26,783	29,141	36,123
15	19,311	22,307	24,996	28,259	30,578	37,697
16	20,465	23,542	26,296	29,633	32,000	39,252
17	21,615	24,769	27,587	30,995	32,409	40,790
18	22,760	25,989	28,869	32,346	34,805	42,312
19	23,900	27,204	30,144	33,678	36,191	43,820
20	25,038	28,412	31,410	35,020	37,566	45,315
21	26,171	29,615	32,671	36,343	38,932	46,797
22	27,301	30,813	33,924	37,659	40,289	48,268
23	28,429	32,007	35,172	38,968	41,638	49,728
24	29,553	33,196	36,415	40,270	42,980	51,179
25	30,675	34,382	37,652	41,566	42,314	52,620
26	31,795	35,563	38,885	42,856	45,642	54,052
27	32,912	36,741	40,113	44,140	46,963	55,476
28	34,027	37,916	41,337	45,419	48,278	56,893
29	35,139	39,087	42,557	46,693	49,588	58,302
30	36,250	40,256	43,773	47,962	50,892	59,703

Application of the goodness-of-fit test  $\chi^2$ -is fulfilled by the next plan.

1. At first on the base of advanced hypothesis we find approximate values of all the probabilities  $P_i = P(a_{i-1} < X < a_i)$  and define the calculated value  $\chi_{calc}^2$ .
2. Knowing the number of degrees of freedom  $\nu$  and specifying some significance level (a small probability)  $\alpha$  we find the critical value  $\chi_{crit}^2$  from the table.
3. Now we compare  $\chi_{calc}^2$  with  $\chi_{crit}^2$ .
  - a) If  $\chi_{calc}^2 \leq \chi_{crit}^2$ , we say that results of trials don't contradict the hypothesis.

b) If  $\chi_{calc}^2 > \chi_{crit}^2$ , we say that results of trials contradict the hypothesis (because of the event  $\chi^2 > \chi_{crit}^2$ , which we regarded as practically impossible, has occurred).

In the case a) we can accept the hypothesis, and in the case b) we can reject it.

Ex. 5. Test the hypothesis that the random variable  $X$  of the **Basic example** has a normal distribution.

1. By virtue of the hypothesis we find approximate values of the probabilities of hitting of the random variable in all the intervals of the interval variation series using the formulas (23), (24). Corresponding evaluations up to finding the calculated value  $\chi_{calc}^2 \approx 4,1218$  are represented in tables 10, 11.

Table 10.

$a_k$	$\frac{a_k - \bar{x}_s}{\sigma_s}$	$\Phi\left(\frac{a_k - \bar{x}_s}{\sigma_s}\right)$	$\Phi\left(\frac{a_{k+1} - \bar{x}_s}{\sigma_s}\right) - \Phi\left(\frac{a_k - \bar{x}_s}{\sigma_s}\right)$
23,12	-2,41	-0,4921	$P_1 = 0,0243$
23,68	-1,85	-0,4678	$P_2 = 0,0681$
24,24	-1,28	-0,3997	$P_3 = 0,1417$
24,80	-0,70	-0,2580	$P_4 = 0,2063$
25,36	-0,13	-0,0517	$P_5 = 0,2217$
25,92	0,44	0,1700	$P_6 = 0,1738$
26,48	1,01	0,3438	$P_7 = 0,0991$
27,04	1,58	0,4429	$P_8 = 0,041$
27,60	2,15	0,4839	$P_9 = 0,0128$
28,16	2,72	0,4967	$\Sigma = 0,9888$

Table 11

№	$x_k$	$P_k$	$m_k$	$m'_k = 100P_k$	$(m_k - m'_k)^2$	$\frac{(m_k - m'_k)^2}{m'_k}$	
1	23,40	0,0243	2 } 13	2,43	14,1376	1,5300	
2	23,96	0,0681		11			6,81
3	24,52	0,1417	14	14,17	0,0289	0,0020	
4	25,08	0,2063	14	20,63	43,9569	2,1307	
5	25,64	0,2217	23	22,17	0,6889	0,0311	
6	26,20	0,1738	20	17,38	6,8644	0,3950	
7	26,76	0,0991	12 } 16	9,91	0,5041	0,0330	
8	27,32	0,041		3			4,1
9	27,88	0,0128		1			1,28
$\Sigma$		0,9888	100	98,88		$\chi_{calc}^2 \approx 4,1218$	



2. Let's choose the significance level (a small probability)  $\alpha = 0,05$ . Finding by the formula (25) the number of degrees of freedom

$$\nu = k - r - 1, k = 6, r = 2, \nu = 3$$

we find the critical value  $\chi_{crit}^2(\nu, \alpha) = \chi_{crit}^2(3, 0,05) = 7,815$  from table 9.

3. We've obtained  $\chi_{calc}^2 = 4,1218 < \chi_{crit}^2$ . It means that the results of trials on the random variable  $X$  in question don't contradict the hypothesis that its distribution law is a normal one.

### **Kolmogorov goodness-of-fit test**

Let  $X$  be a continuous random variable, and we advance the next hypothesis: the given function  $F(X)$  is the distribution function of this random variable.

Let a number  $D$  be the greatest value of the absolute value of a difference of the function  $F(X)$  and statistical distribution function  $F^*(X)$  that is

$$D_{calc.} = \sqrt{n} \cdot \max_{1 \leq i \leq k} |F^*(x_i) - F(x_i)| \quad (26)$$

Table 12 gives some values of the  $D_{crit}$ .

Table 12. The values of the  $D_{crit}$ .

$\alpha$	0,20	0,10	0,05	0,02	0,01	0,001
$D_{crit}$	1,073	1,224	1,358	1,520	1,627	1,950

a) If  $D_{calc} \leq D_{crit}$ , we say that results of trials don't contradict the hypothesis.

b) If  $D_{calc} > D_{crit}$ , we say that results of trials contradict the hypothesis (because of the event  $D_{calc} > D_{crit}$ , which we regarded as practically impossible, has occurred).

In the case a) we can accept the hypothesis, and in the case b) we can reject it.

### **Tasks for individual work on mathematical statistics**

**PROBLEM.** 100 independent trials are fulfilled on a random variable  $X$ , and the results of trials are represented by corresponding sample of the size  $n = 100$ .

1. Compile the interval variation series for the random variable  $X$ , plot the **histogram of relative frequencies** and approximate graph of the **distribution density**. Find the **statistical distribution function** for the interval variation series and construct its approximate graph.

2. On the base of the interval variation series form the **discreet variation series** by taking inner points in each interval. Construct the **polygon of relative frequencies**. Form the **statistical distribution function** for the discreet variation series and plot its graph.

3. Calculate the sample mean [the sample average], dispersion, root-mean-square deviation of the random variable. Draw a conclusion as to symmetry of its distribution law and deviation of this law from the normal distribution.

4. Find the corrected dispersion and root-mean-square deviation of the random variable. Compare their values with corresponding sample estimators.

5. Search out the confidence intervals for the mathematical expectation of the random variable  $X$  with reliabilities 0.95 and 0.99, basing on the hypothesis of its normal distribution.

6. Find approximate values of the probabilities of hitting of the random variable  $X$  on all the intervals of its variation series proceeding from the hypothesis of the normal distribution of  $X$ .

7. Test the hypothesis on the normal distribution of the random variable  $X$  with the significance levels  $\alpha$  0.01, 0.025, 0.05 making use of Pearson  $\chi^2$ -goodness-of-fit tests.

1.

26,65	26,55	26,25	26,20	26,15	26,00	26,00	25,95	25,85	26,25
25,80	25,75	25,70	25,60	25,50	25,50	25,35	25,10	25,10	25,65
25,35	25,50	25,55	25,65	25,70	25,70	25,75	25,75	25,85	25,60
25,85	25,95	26,00	26,15	26,20	26,25	26,45	26,55	26,65	26,00
26,65	26,55	26,45	26,25	26,15	26,10	26,00	26,00	25,85	26,25
25,75	25,70	25,65	25,55	25,50	25,40	25,10	26,85	25,70	25,60
25,30	25,45	25,50	25,65	26,70	25,70	25,75	25,85	25,85	25,55
26,00	26,00	26,10	26,25	25,55	26,55	26,65	25,65	26,25	26,20
26,65	26,55	26,25	26,00	25,85	25,75	25,65	25,00	25,30	26,10
26,65	26,25	26,20	25,90	25,80	25,70	25,65	25,00	25,50	26,00

2.

2,54	0,69	2,59	2,60	0,89	1,85	0,95	2,89	1,07	2,92
2,23	1,13	2,27	2,30	1,22	1,75	1,25	2,44	1,29	2,22
2,08	1,38	2,10	2,11	1,43	1,64	1,46	2,20	1,49	1,48
1,91	1,58	1,93	1,95	1,62	1,44	1,66	2,04	1,69	2,07
1,80	2,58	1,82	1,83	1,75	1,25	1,76	1,87	1,78	1,90
1,71	2,24	1,73	1,19	1,84	0,90	2,84	1,77	1,88	1,09
1,54	2,09	1,60	1,41	1,97	1,99	2,34	1,68	2,05	1,31
1,34	1,91	1,40	1,61	2,13	2,15	2,18	1,28	2,90	1,51
0,67	1,81	1,16	1,74	2,31	2,33	2,03	1,28	2,46	1,31
1,12	1,72	0,73	0,84	2,69	2,75	1,86	1,02	2,21	1,09

3.

25,45	25,50	26,00	25,85	25,65	25,70	26,00	26,00	25,65	25,85
25,65	25,50	25,70	26,00	25,85	25,75	26,10	26,35	25,90	25,95
25,50	25,85	25,50	25,40	25,75	26,10	25,50	25,75	25,25	25,70
25,80	26,25	25,75	25,75	26,15	26,10	26,00	27,00	25,85	26,00
25,00	25,65	26,25	26,10	25,85	26,65	26,35	25,50	25,65	26,30
26,15	25,85	25,50	26,10	25,25	25,75	25,55	26,75	25,50	25,50
25,50	25,75	25,90	26,10	25,75	25,85	26,10	25,55	26,35	26,50
25,65	26,25	25,70	25,05	25,75	26,35	25,75	25,75	25,85	25,65
25,75	25,70	26,25	25,90	26,10	25,50	26,15	25,50	25,50	25,40
25,85	25,75	25,50	26,10	25,75	26,00	25,50	26,05	26,05	25,60

4.

2,74	3,01	2,99	3,11	3,00	3,00	3,02	3,02	2,93	3,18
3,13	2,93	3,11	3,01	3,00	3,05	2,99	3,07	2,91	3,04
3,15	3,06	3,03	3,07	3,21	3,13	3,09	3,01	3,07	3,04
2,99	2,93	3,02	3,05	3,03	2,97	2,93	3,09	3,11	3,00
3,09	3,09	3,01	3,04	3,07	3,10	3,19	2,92	3,03	3,05
2,78	3,15	3,09	3,06	3,03	3,08	2,96	3,16	2,95	3,00
3,00	3,32	3,12	3,05	3,01	3,02	2,89	3,02	2,99	3,14
3,01	3,12	2,98	3,03	2,95	3,03	3,12	3,11	3,10	3,01
3,02	3,16	3,08	3,08	2,97	3,08	2,95	2,98	3,02	2,97
3,01	3,10	3,12	3,02	3,11	2,92	2,99	3,02	3,04	3,05

5.

66,3	66,0	55,5	57,9	59,3	61,7	47,3	49,7	63,5	65,9
54,0	56,4	15,3	17,7	26,3	38,7	26,8	29,2	51,3	53,7
39,8	42,2	80,3	82,7	68,3	70,7	60,3	62,7	61,8	54,1
62,5	64,9	55,2	57,6	54,1	56,5	74,3	76,7	35,3	37,7
61,9	64,3	29,1	31,5	55,2	57,7	28,3	30,7	53,1	55,5
52,7	55,1	58,5	59,0	62,9	66,3	37,7	50,1	68,7	62,1
42,9	45,3	38,8	41,2	55,3	57,7	45,0	47,4	62,8	65,2
51,6	76,8	54,8	57,2	32,0	34,4	73,9	76,3	36,4	37,0
35,3	70,7	82,0	54,0	42,4	88,4	44,6	76,9	35,3	47,0
49,3	44,3	49,1	37,7	70,1	72,5	44,1	46,5	47,6	50,0

6.

122	172	208	187	208	194	115	201	172	93
129	115	194	43	108	86	122	122	108	201
136	180	100	158	151	129	144	187	172	144
115	79	158	129	100	136	122	108	158	172
122	165	100	158	151	108	93	172	158	72
151	64	129	136	108	151	144	100	57	129
165	151	144	165	151	136	165	100	129	108
100	136	158	129	50	122	151	151	86	129
172	108	100	129	115	100	180	136	93	122
172	136	115	122	151	129	172	144	86	144

7.

66,3	66,0	55,5	57,9	59,3	61,7	46,3	48,7	47,3	49,7
54,0	56,4	15,3	17,7	26,3	28,7	50,3	52,7	26,8	29,2
39,8	42,2	80,3	82,7	68,3	70,7	20,6	24,0	60,3	62,7
62,5	64,9	55,2	57,6	54,1	56,5	55,0	57,4	74,3	76,7
61,9	64,3	29,1	31,5	55,2	57,7	76,5	78,9	28,3	30,7
52,7	55,1	58,5	59,0	62,9	66,3	33,7	46,1	47,7	50,1
42,9	45,3	38,8	41,2	55,3	57,7	32,3	34,7	45,0	47,4
51,6	76,8	54,8	57,2	32,0	34,4	59,3	61,7	73,9	76,3
35,3	70,7	82,0	54,0	42,4	88,4	74,5	44,8	44,6	76,9
49,3	44,3	49,1	37,7	70,1	72,5	27,3	29,7	44,1	46,5

8.

59,3	61,7	46,3	48,7	90,1	92,5	47,3	49,7	81,8	84,2
26,3	28,7	50,3	52,7	47,8	50,2	60,3	62,7	54,3	56,7
68,3	70,7	20,6	24,0	68,1	70,5	83,8	81,2	72,3	74,7
54,1	56,5	55,0	57,4	27,7	30,1	57,3	59,7	75,3	78,1
55,2	57,7	76,5	78,9	56,1	58,5	83,6	86,0	32,3	34,7
62,9	66,3	33,7	46,1	54,1	56,5	54,1	56,5	54,7	65,2
55,3	57,7	32,3	34,7	23,7	26,1	40,9	40,3	33,7	36,1
32,0	34,4	59,3	61,7	36,5	38,9	43,1	51,3	74,4	51,7
42,4	88,4	74,5	44,8	51,4	50,7	93,5	53,8	68,3	97,7
70,1	72,5	27,3	29,7	46,3	48,7	23,1	25,5	35,3	37,7

9.

55,5	57,9	59,3	61,7	46,3	48,7	47,3	49,7	63,5	65,9
15,3	17,7	26,3	28,7	50,3	52,7	26,8	29,2	51,3	53,7
80,3	82,7	68,3	70,7	20,6	24,0	60,3	62,7	61,8	54,1
55,2	57,6	54,1	56,5	55,0	57,4	74,3	76,7	35,3	37,7
29,1	31,5	55,2	57,7	76,5	78,9	28,3	30,7	53,1	55,5
58,5	59,0	62,9	66,3	33,7	46,1	47,7	50,1	68,7	62,1
38,8	41,2	55,3	57,7	32,3	34,7	45,0	47,4	62,8	65,2
54,8	57,2	32,0	34,4	59,3	61,7	73,9	76,3	36,4	37,0
82,0	54,0	42,4	88,4	74,5	44,8	44,6	76,9	35,3	47,0
49,1	37,7	70,1	72,5	27,3	29,7	44,1	46,5	47,6	50,0

10.

81,8	54,3	72,3	75,3	32,3	62,8	33,7	74,4	68,3	35,3
84,2	56,7	74,7	78,1	34,7	65,2	36,1	51,7	97,7	37,7
64,6	49,3	59,3	37,3	52,7	50,0	50,1	72,5	68,0	48,3
70,1	67,0	51,7	60,7	39,7	55,1	59,4	54,4	70,4	37,7
63,5	51,3	61,8	35,3	53,1	68,7	62,8	36,4	35,3	47,6
65,9	53,7	54,1	37,7	55,5	62,1	65,2	37,0	47,0	50,0
46,3	50,3	20,6	55,0	76,5	33,7	32,3	59,3	74,5	27,3
48,7	52,7	24,0	57,4	78,9	46,1	34,7	61,7	44,8	29,7
47,3	60,3	83,8	57,3	83,6	54,7	40,9	43,1	93,5	23,1
25,5	53,8	51,3	40,3	57,1	86,0	59,7	81,2	62,7	49,7

11.

4,03	4,05	4,13	4,26	4,25	4,13	4,12	4,25	4,14	4,05
4,07	4,14	4,13	4,22	4,14	4,15	4,19	4,29	4,30	4,16
4,16	4,27	4,05	4,25	4,22	4,14	4,28	4,08	4,17	4,25
4,25	4,30	4,03	4,03	4,15	4,19	4,30	4,17	4,13	4,30
4,17	4,15	4,14	4,14	4,15	4,16	4,26	4,15	4,16	4,07
4,06	4,08	4,13	4,05	4,13	4,13	4,13	4,12	4,14	4,15
4,19	4,15	4,04	4,11	4,11	4,19	4,17	4,16	4,26	4,21
4,20	4,25	4,12	4,26	4,25	4,28	4,28	4,14	4,09	4,03
4,28	4,25	4,23	4,24	4,18	4,22	4,18	4,20	5,14	4,14
4,10	4,11	4,30	4,17	4,20	4,10	4,22	4,17	4,11	4,15

12.

5,50	5,53	5,55	5,60	5,53	5,57	5,55	5,49	5,60	5,57
5,50	5,50	5,45	5,53	5,55	5,62	5,57	5,65	5,62	5,55
5,49	5,60	5,50	5,57	5,60	5,40	5,53	5,55	5,65	5,33
5,60	5,45	5,57	5,50	5,62	5,68	5,50	5,55	5,60	5,57
5,62	5,45	5,50	5,53	5,60	5,40	5,55	5,57	5,53	5,49
5,68	5,55	5,60	5,68	5,57	5,70	5,55	5,45	5,57	5,55
5,60	5,55	5,60	5,49	5,50	5,53	5,57	5,55	5,62	5,53
5,60	5,70	5,53	5,55	5,60	5,49	5,50	5,62	5,53	5,55
5,60	5,47	5,57	5,55	5,55	5,49	5,53	5,57	5,60	5,68
5,57	5,60	5,49	5,53	5,60	5,62	5,49	5,50	5,67	5,50

13.

226	113	307	356	259	307	291	291	324	221
437	421	405	421	405	421	437	405	437	405
388	307	372	388	324	374	372	307	291	307
372	359	307	226	324	243	197	194	259	340
324	324	340	259	372	275	388	324	307	340
372	324	324	340	356	291	327	307	243	324
243	243	356	275	453	486	291	226	340	291
307	372	356	356	307	178	145	129	259	324
372	388	324	259	259	275	210	162	162	178
275	243	388	243	340	226	275	243	240	291

14.

43,9	43,7	68,7	43,6	17,5	45,0	27,2	43,2	40,0	23,7
47,7	14,7	56,7	42,5	43,7	52,3	43,7	20,4	70,4	37,5
34,7	38,7	39,0	43,4	64,9	22,1	20,7	47,7	30,8	58,5
35,7	15,2	48,7	15,7	69,9	62,3	33,4	36,1	16,7	32,7
33,0	32,5	23,0	51,2	48,1	23,7	41,5	40,2	51,9	39,7
37,7	53,0	46,7	25,7	41,1	45,4	41,5	56,4	23,7	36,0
58,4	36,7	24,9	12,1	42,5	44,5	16,1	56,6	36,2	78,5
34,7	39,8	37,5	43,1	72,0	45,6	72,1	48,7	35,7	29,3
70,2	42,7	40,7	64,1	20,7	51,2	22,1	37,7	83,7	11,5
23,6	56,7	62,8	31,3	41,1	50,3	50,9	28,2	52,0	42,4

15.

66,3	66,0	29,7	27,3	63,5	65,9	90,1	48,7	47,3	49,7
54,0	56,4	44,8	74,5	51,3	53,7	47,8	50,7	26,8	29,2
39,8	42,2	61,7	59,3	61,8	54,1	68,1	38,9	60,3	62,7
62,5	64,9	34,7	32,3	35,3	37,7	27,7	26,1	74,3	76,7
61,9	64,3	46,1	33,7	53,1	55,5	56,1	56,5	28,3	30,7
52,7	55,1	78,9	76,5	68,7	62,1	54,1	58,5	47,7	50,1
42,9	45,3	57,4	55,0	62,8	65,2	23,7	30,1	45,0	47,4
51,6	76,8	24,0	20,6	36,4	37,0	36,5	70,5	73,9	76,3
35,3	70,7	52,7	50,3	35,3	47,0	51,4	50,2	44,6	76,9
49,3	44,3	48,7	46,3	47,6	50,0	46,3	92,5	44,1	46,5

## 7. ELEMENTS OF CORRELATION THEORY

There is deterministic [functional, stiff] dependence between random variables  $X, Y$  (for example linear dependence  $Y = aX + b$ ). And there is undetermined [non-functional, non-stiff, statistic, correlation] dependence between  $X, Y$ , for example dependence between labour productivity and living standard, between a state of health and a productivity of a worker, between height and weight of a man.

Correlation dependence between random variables  $X$  and  $Y$  is that between values of one variable and corresponding mean value [average value, distribution centre, mathematical expectation] of the other. Such dependence is defined by introducing conditional distributions of random variables and conditional mathematical expectations.

**Def.1.** The conditional mathematical expectation of the random variable  $Y$ , that is  $f(x)$ , is called the **regression function** of  $Y$  on  $X$ , its graph is called the **regression curve** [regression line] of  $Y$  on  $X$  and the equation

$$\bar{y}_x = f(x) \quad (1)$$

is called the **regression equation** of  $Y$  on  $X$ .

By analogy the regression function  $\phi(y)$ , the regression curve [regression line] and the regression equation are defined

$$\bar{x}_y = \phi(y) \quad (2)$$

**Def. 2.** A correlation dependence between random variables  $X, Y$  is called a functional dependence between possible values of one random variable and corresponding regression function (average [mean] value) of the other.

### **Main problems of correlation theory (for the case of two random variables)**

1. Determine the **form of correlation dependence** between random variables.

If regression functions of random variables  $X, Y$  are linear, then one says about a linear correlation (or a linear correlation dependence) between these random variables. Otherwise one says about non-linear (or curvilinear) correlation.

2. Determine the **closeness of relation between** random variables  $X, Y$ . Closeness problem is resolved with the help of the correlation coefficient  $r_s$ , which is the measure of a linear dependence between  $X, Y$ , and the correlation ratios  $\rho_{xy}, \rho_{yx}$ , which are the measure of a functional (not necessary linear) dependence between  $X, Y$ .

The **correlation coefficient** of random variables  $X, Y$  is defined by the formula

$$r_s = \frac{\sum xym_{xy} - n\bar{x}\bar{y}}{n\sigma_x\sigma_y} \quad (3)$$

where

$$\sum xym_{xy} - n\bar{x}\bar{y} \quad (4)$$

is the correlation moment of  $X, Y$ .

If random variables are independent, then  $r_s = 0$ . The converse isn't true in general: there are dependent random variables with  $r_s = 0$ .

**Def. 3.** Random variables  $X, Y$  are called those **correlated** if their correlation coefficient doesn't equal zero ( $r_s \neq 0$ ), and **non-correlated** otherwise ( $r_s = 0$ ).

It's known that

$$|r_s| \leq 1, \text{ and } |r_s| = 1$$

if  $X, Y$  are connected by linear (functional) dependence  $Y = aX + b$ .

### Linear correlation

Let the regression functions of random variables  $X, Y$  be linear, that is there is a linear correlation between  $X, Y$ . It can be proved that the regression functions are given by the next formulas

$$\bar{y}_x - \bar{y} = r_s \cdot \frac{\sigma_y}{\sigma_x} (x - \bar{x}), \bar{x}_y - \bar{x} = r_s \cdot \frac{\sigma_x}{\sigma_y} (y - \bar{y}) \quad (5)$$

Corresponding regression straight lines  $\bar{y}_x = f(x)$  and  $\bar{x}_y = \varphi(y)$ :

a) intersect at the point  $(\bar{x}, \bar{y})$ ;

b) have the slopes  $k_1 = r_s \cdot \frac{\sigma_y}{\sigma_x}, k_2 = r_s \cdot \frac{\sigma_x}{\sigma_y}$ ;

c) coincide if  $|r_s| = 1$  that is if there is a linear functional dependence between the random variables  $X$  and  $Y$ ;

d) are perpendicular respectively to the  $Ox$ -axis and  $Oy$ -axis if  $r_s = 0$  that is if  $X$  and  $Y$  aren't correlated.

**The simplest case. Every pair of random variables was observed only one time.**

Let each pair  $(X = x_i, Y = y_i)$  of the random variables  $X, Y$  occur only one time. In this case we'll use next formulas

$$\bar{x}_s = \frac{\sum x_i}{n}, \bar{x}^2 = \frac{\sum x_i^2}{n}, \sigma_s^2(x) = \bar{x}^2 - (\bar{x}_s)^2, \bar{y}_s = \frac{\sum y_i}{n}, \bar{y}^2 = \frac{\sum y_i^2}{n},$$

$$\sigma_s^2(y) = \sqrt{\bar{y}^2 - \bar{y}_s^2}, \overline{xy} = \frac{\sum x_i y_i}{n}, r_s = \frac{\sum x y m_{xy} - n \bar{x} \bar{y}}{n \sigma_x \sigma_y}.$$

### General case

Now we'll consider the **general case** when as a rule several values of a random variable  $Y$  correspond to an arbitrary value of a random variable  $X$  or when an arbitrary pair of values  $(X = x_i, Y = y_j)$  of random variables  $X, Y$  appears  $n_{ij}$  times. We study correlative dependence between random variables  $X, Y$  with the help of the next table, which is called the **correlation table (tab. 13)**.

$X \backslash Y$	$y_1$	$y_2$	...	$y_p$	$m_x$
$x_1$	$n_{11}$	$n_{12}$	...	$n_{1p}$	$m_{x1}$
$x_2$	$n_{21}$	$n_{22}$	...	$n_{2p}$	$m_{x2}$
...	...	...	...	...	...
$x_k$	$n_{k1}$	$n_{k2}$	...	$n_{kp}$	$m_{xk}$
$m_y$	$m_{y1}$	$m_{y2}$	...	$m_{yp}$	$n$

In this case we'll use next formulas

$$\bar{x}_s = \frac{\sum x_i m_x}{n}, \bar{x}^2 = \frac{\sum x_i^2 m_x}{n}, \sigma_s^2(x) = \bar{x}^2 - (\bar{x}_s)^2, \sigma_s^2(y) = \bar{y}^2 - (\bar{y}_s)^2$$

$$\bar{y}_s = \frac{\sum y_j m_y}{n}, \bar{y}^2 = \frac{\sum y_j^2 m_y}{n}, \overline{xy} = \frac{\sum \sum x_i y_j m_{xy}}{n}, r_s = \frac{\sum x y m_{xy} - n \bar{x} \bar{y}}{n \sigma_x \sigma_y}.$$

Ex. 4. Results of  $n = 100$  trials on random variables  $X$  and  $Y$  are represented by the table. Estimate number characteristics of  $X$ ,  $Y$ , regression functions of  $X$  on  $Y$  and  $Y$  on  $X$ , correlation coefficient.

X \ Y	10	20	30	40	50	60	$m_x$
2	2	4					6
8		3	7				10
14		1	48	10	2		61
20			2	7	5		14
26				1	2	2	5
32					2	2	4
$m_y$	2	8	57	18	11	4	$n=100$

Solution.

We insert preliminary calculations into an "extended" table.

X \ Y	10	20	30	40	50	60	$m_x$	$xm_x$	$x^2m_x$
2	2	4					6	12	24
8		3	7				10	80	640
14		1	48	10	2		61	854	11956
20			2	7	5		14	280	5600
26				1	2	2	5	130	3380
32					2	2	4	128	4096
$m_y$	2	8	57	18	11	4	$n = 100$	1484	25696
$ym_y$	20	160	1710	720	550	240	3400		
$y^2m_y$	200	3200	51300	28800	27500	14400	125400		

Estimates of number characteristics of the random variables  $X$  and  $Y$ .

$$\bar{x} = \frac{\sum xm_x}{n} = \frac{1484}{100} = 14,84;$$

$$D_x = \frac{\sum x^2m_x}{n} - \bar{x}^2 = \frac{25696}{100} - 14,84^2 \approx 36,73; \quad \sigma_x = \sqrt{36,73} \approx 6,06;$$

$$\bar{y} = \frac{\sum ym_y}{n} = \frac{3400}{100} = 34;$$

$$D_y = \frac{\sum y^2m_y}{n} - \bar{y}^2 = \frac{125400}{100} - 34^2 = 98, \quad \sigma_y = \sqrt{98} \approx 9,90.$$

$$\sum xym_{xy} = 2(10 \cdot 2 + 20 \cdot 4) + 8(20 \cdot 3 + 30 \cdot 7) + 14(20 \cdot 1 + 30 \cdot 48 + 40 \cdot 10 + 50 \cdot 2) + 20(30 \cdot 2 + 40 \cdot 7 + 50 \cdot 5) + 26(40 \cdot 1 + 50 \cdot 2 + 60 \cdot 2) + 32(50 \cdot 2 + 60 \cdot 3) = 55400.$$

Let's pass to the estimation of the correlation coefficient of random variables  $X$  and  $Y$ .

$$r_s = \frac{\sum xym_{xy} - n \cdot \bar{x} \cdot \bar{y}}{n \cdot \sigma_x \sigma_y} = \frac{55400 - 100 \cdot 14,84 \cdot 34}{100 \cdot 6,06 \cdot 9,9} \approx 0,824.$$

This latter is sufficiently large, so we can say that  $X$  and  $Y$  are connected by essential linear dependence.

By the formula (5) we find the estimates of values of the regression function  $X$  on  $Y$

$$\bar{y}_x - \bar{y} = r_s \cdot \frac{\sigma_y}{\sigma_x} (x - \bar{x}),$$

$$\bar{y}_x - 34 = 0,824 \cdot \frac{9,9}{6,06} (x - 14,84),$$

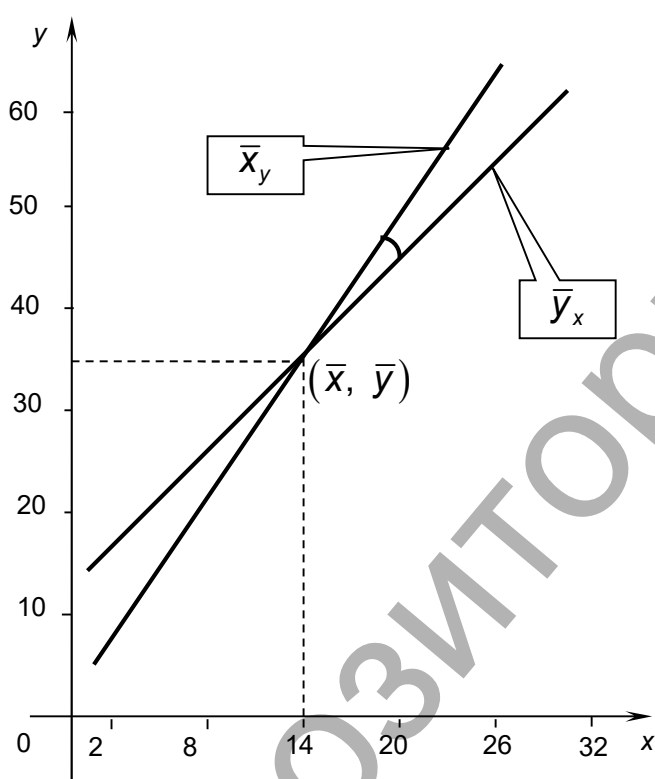
$$\bar{y}_x = 1,35x + 13,97.$$

and  $Y$  on  $X$ :

$$\bar{x}_y - \bar{x} = r_s \cdot \frac{\sigma_x}{\sigma_y} (y - \bar{y})$$

$$\bar{x}_y - 14,84 = 0,824 \cdot \frac{6,06}{9,9} (y - 34)$$

$$\bar{x}_y = 0,5y - 2,16.$$



Let us build the graphs of the obtained straight lines on one drawing.

The nearer to zero the acute angle between them (marked by the arc), the closer the connection between the features. But if this angle is close to  $90^\circ$ , this indicates that the connection is weak or is absent at all.

### Tasks for individual work on elements of correlation theory

Results of  $n = 100$  trials on random variables  $X$  and  $Y$  are represented by the table. Estimate number characteristics of  $X$ ,  $Y$ , regression functions of  $X$  on  $Y$  and  $Y$  on  $X$ , correlation coefficient.

7.01

$X \backslash Y$	2	5	8	11	14	17	$m_x$
15	5	3	6				14
25		7	8	11			26
35			9	10	12		31
45				9	9		18
55					7	4	11
$m_y$	5	10	23	30	28	4	$n = 100$



## 7.02

X \ Y	3	5	7	9	11	13	$m_x$
12	4	3	5				12
14	6	7	8				21
16		10	12	11			33
18			8	8	5		21
20				4	5	4	13
$m_y$	10	20	33	23	10	4	$n = 100$

## 7.03

X \ Y	4	8	12	16	20	24	$m_x$
1	3	2	9				14
4		7	10	9			26
7			12	10	5		27
10				9	8	5	22
13					6	5	13
$m_y$	3	9	31	28	19	10	$n = 100$

## 7.04

X \ Y	11	14	17	20	23	26	$m_x$
10	4	6	3				13
15		7	9	10			26
20			13	9	7		29
25				12	6	3	21
30					6	5	11
$m_y$	4	13	25	31	19	8	$n = 100$

## 7.05

X \ Y	4	12	20	28	36	44	$m_x$
1				5	4	6	15
5				11	6		17
9			5	14	8		27
13		9	8	7			24
17	4	6	7				17
$m_y$	4	15	20	37	18	6	$n = 100$

## 7.06

X \ Y	6	8	10	12	14	16	$m_x$
11				6	5	3	14
16			10	8	6		24
21		7	12	9			28
26	6	8	10				24
31	3	7					10
$m_y$	9	22	32	23	11	3	$n = 100$

## 7.07

X \ Y	2	8	14	20	26	32	$m_x$
10					7	3	10
20			6	9	7	2	24
30		8	12	10			30
40	5	7	12				24
50	4	6	2				12
$m_y$	9	21	32	19	14	5	$n = 100$

## 7.08

X \ Y	10	15	20	25	30	35	$m_x$
5				5	5	3	13
12			5	9	7		21
19		9	13	6			28
26	9	8	10				27
33	4	7					11
$m_y$	13	24	28	20	12	3	$n = 100$

## 7.09

X \ Y	2	4	6	8	10	12	$m_x$
13				5	4	3	12
17			5	10	6		21
21			9	14	5		28
25		6	12	6			24
29	5	4	6				15
$m_y$	5	10	32	35	15	3	$n = 100$

## 7.10

X \ Y	14	21	28	35	42	49	$m_x$
3				3	6	5	14
4				9	5	3	17
5			16	8	4		28
6		7	10	5			22
7	6	13					19
$m_y$	6	20	26	25	15	8	$n = 100$

## Literature

1. Davar Khoshnevisan, Firas Rassoul-Agha Math 5010. Introduction to Probability. – 2012.
2. J.F. Kosolapov Probability theory and mathematical statistics. – 2008.

## Contents

<b>1</b>	<b>EVENT AND PROBABILITY</b> .....	3
1.1	Trial and event.....	3
1.2	Elements of combinatorics.....	4
1.3	Classic definition of probability.....	5
1.4	Statistic definition of probability.....	6
<b>2</b>	<b>MAIN RULES OF EVALUATING PROBABILITIES</b> .....	7
2.1	Sum and product of events.....	7
2.2	Axioms of probability theory. corollaries.....	8
2.3	Formulae of total probability and bayes.....	10
	Exercise Set 1, 2.....	12
	Homework Problems.....	12
<b>3</b>	<b>RANDOM VARIABLES</b> .....	14
3.1	A random variable.....	14
3.2	Bernoulli [binomial] distribution.....	16
3.3	Poisson formula and distribution.....	17
3.4	Laplace local and integral theorems.....	19
	Exercise Set 3.....	22
<b>4</b>	<b>THE DISTRIBUTION FUNCTION AND DENSITY. NUMBER CHARACTERISTICS OF RANDOM VARIABLES</b> .....	23
4.1	The distribution function of a random variable.....	23
4.2	The distribution density of a random variable.....	25
4.3	The mathematical expectation of a random variable.....	26
4.4	The dispersion and root-mean-square deviation.....	27
4.5	Moments of a random variable.....	28
	Exercise Set 4.....	29
<b>5</b>	<b>SOME REMARKABLE DISTRIBUTIONS</b> .....	30
5.1	The uniform distribution.....	30
5.2	The normal distribution.....	31
5.3	The exponential distribution.....	33
5.4	Bernoulli [binomial] distribution.....	34
5.5	Poisson formula and distribution.....	35
	Exercise Set 5.....	35
<b>6</b>	<b>ELEMENTS OF MATHEMATICAL STATISTICS</b> .....	36
6.1	General remarks. sampling method. Variation series.....	36
6.2	Approximate determination of the distribution law of a random variable. Estimation of parameters of the distribution law of a random variable.....	40
6.3	Testing statistic hypotheses.....	45
	Tasks for individual work on mathematical statistics.....	49
<b>7</b>	<b>ELEMENTS OF CORRELATION THEORY</b> .....	53
	Tasks for individual work on elements of correlation theory.....	56
	Literature.....	58

**УЧЕБНОЕ ИЗДАНИЕ**

Составители:

*Гладкий Иван Иванович  
Дворниченко Александр Валерьевич  
Каримова Татьяна Ивановна  
Лебедь Светлана Федоровна  
Шишко Татьяна Витальевна*

**PROBABILITY THEORY  
ELEMENTS OF MATHEMATICAL STATISTICS**

учебно-методическая разработка на английском языке

Ответственный за выпуск: Дворниченко А.В.

Редактор: Боровикова Е.А.

Компьютерная верстка: Кармаш Е.Л.

Корректор: Шишко Т.В.

---

Подписано к печати 24.09.2014 г. Формат 60x84 <sup>1</sup>/<sub>16</sub>. Бумага «Снегурочка».

Усл. п. л. 3,49. Уч.-изд. л. 3,75. Заказ № 790. Тираж 40 экз.

Отпечатано на ризографе Учреждения образования  
«Брестский государственный технический университет».

224017, г. Брест, ул. Московская, 267.