

Министерство образования и науки Российской Федерации



Федеральное государственное бюджетное образовательное учреждение  
высшего образования

Санкт-Петербургский горный университет

**Кафедра высшей математики**

## **Расчетно-графическая работа**

По дисциплине Математика  
(наименование учебной дисциплине согласно учебному плану)

Тема работы: Основы математической статистики

**Вариант 106**

Выполнил: ст. группы КРС-21 \_\_\_\_\_ Крюков Т.В.  
(должность) (подпись) (Ф.И.О)

Дата: \_\_\_\_\_

Проверил: доцент \_\_\_\_\_ Мансурова С.Е.  
(должность) (подпись) (Ф.И.О)

Санкт-Петербург  
2022

## Задача 1.

### Статистическая обработка массива данных.

### Выборочные оценки генеральной совокупности.

#### Условие задачи:

1) Составить сгруппированный статистический ряд:

а) определить нужное кол-во интервалов по формуле Стерджесса

$k \approx 1 + \log_2 N$ , где  $N$ -объем выборки;

б) найти наибольшее и наименьшее значение измеряемой величины;

в) определить ширину интервала. Используемая формула:

$$\Delta x = \left[ \frac{x_{\max} - x_{\min}}{k} \right] + 1, \text{ где } [\ ] - \text{округление вверх}$$

г) интервалы берутся открытыми слева и закрытыми справа:  $(\alpha_i; \beta_i]$ ,  $i=1..k$

д) концы интервалов вычисляются по формулам:

$$\alpha_i = x_{\min} - \frac{\Delta x}{2}; \alpha_{i+1} = \alpha_i + \Delta x, \beta_i = \alpha_i + \Delta x$$

е) определить середины интервалов:  $x_i = \frac{\alpha_i + \beta_i}{2}$

2) Для сгруппированного статистического ряда вычислить выборочное среднее, выборочную дисперсию, выборочное среднее квадратическое отклонение, исправленную дисперсию, исправленное среднее квадратическое отклонение.

3) Считая генеральную совокупность нормально распределенной, найти интервальные оценки математического ожидания и среднего квадратического отклонения с надежностью 0,95

4) Построить гистограмму, полигон частот и эмпирическую функцию распределения.

**Исходные данные:**

A	28	13	55	9	79	37
Вар.	29	15	56	29	80	26
106	30	25	57	16	81	24
	31	-6	58	25	82	31
	32	27	59	32	83	27
39	33	42	60	28	84	43
2	34	21	61	18	85	-5
17	35	27	62	22	86	41
29	36	22	63	12	87	31
17	37	18	64	-2	88	15
18	38	-5	65	3	89	6
2	39	24	66	17	90	3
6	40	16	67	39	91	20
17	41	5	68	16	92	4
27	42	22	69	-2	93	4
1	43	15	70	2	94	42
48	44	34	71	13	95	16
4	45	5	72	41	96	13
26	46	13	73	21	97	20
18	47	33	74	22	98	27
8	48	29	75	29	99	25
42	49	22	76	16	100	1
25	50	11	77	44	101	5
36	51	29	78	23	102	25
47	52	6	79	37	103	21
	53	19	80	26	104	
	54	37	81	24	105	
	55	9	82	21	106	

**Решение:**

- 1) Объем выборки N=100. Определим кол-во интервалов для группировки по формуле Стеджесса:

$k \approx \sqrt[n]{N}$ , где [...] – символ округления (вверх) до целого.

Найдем наибольшее и наименьшее значения измеряемой величины

$x_{min}$	$x_{max}$
-6	48

Определим ширину интервала.

$$\Delta x = \left[ \frac{x_{max} - x_{min}}{k} \right] + 1 = \left[ \frac{48 + 6}{8} \right] + 1 = [7] + 1 = 8$$

Составим таблицу со сгруппированным статистическим рядом:

Номер интервала	Интервал	Середина интервала	Частота	Относительная частота
i	$(\alpha_i; \beta_i]$	$x_i$	$m_i^*$	$p_i^*$
1	$(-10; -2]$	-6	5	0,05
2	$(-2; 6]$	2	17	0,17
3	$(6; 14]$	10	8	0,08
4	$(14; 22]$	18	27	0,27
5	$(22; 30]$	26	22	0,22
6	$(30; 38]$	34	9	0,09
7	$(38; 46]$	42	10	0,1
8	$(46; 54]$	50	2	0,02
Контрольные суммы			100	1

Рис.1. Таблица со сгруппированным статистическим рядом

Начало первого интервала вычисляется по формуле:  $\alpha_i = x_{min} - \frac{\Delta x}{2} = -6 - \frac{8}{2} = -10$ ;  
 для остальных :  $\alpha_{i+1} = \alpha_i + \Delta x$  ,  $\beta_i = \alpha_i + \Delta x$

## 2) Числовые характеристики

(точечные оценки параметров генеральной совокупности)

Выборочное среднее:

$$\bar{x}^* = \frac{1}{N} \sum_{i=1}^4 x_i m_i^* \quad \text{или} \quad \bar{x}^* = \sum_{i=1}^4 x_i p_i^* .$$

Выборочная дисперсия:

$$D^*(x) = \frac{1}{N} \sum_{i=1}^4 x_i^2 m_i^* - (\bar{x}^*)^2 \quad \text{или} \quad D^*(x) = \sum_{i=1}^4 x_i^2 p_i^* - (\bar{x}^*)^2 .$$

Выборочное среднее квадратическое отклонение:

$$\sigma_x^* = \sqrt{D^*(x)} .$$

Несмещенная дисперсия:

$$s_x^2 = \frac{n}{n-1} D^*(x)$$

Исправленное среднее квадратическое отклонение

$$s_x = \sqrt{s_x^2}$$

Полученные результаты из Excel:

Выборочное среднее	Выборочная дисперсия	Выборочное среднее квадратичное отклонение	Несмещенная дисперсия	Исправление среднее квадратичного отклонения
$\bar{x}^*$	$D^*(x)$	$\delta_x^*$	$S_x^2$	$S_x$
19,68	189,8	13,78	191,7	13,85

### 3) Интервальные оценки параметров генеральной совокупности

Теория: 1. Оценка мат. ожидания при неизвестной дисперсии

*Интервальной* называют оценку, которая определяется двумя числами — концами интервала, покрывающего оцениваемый параметр.

*Доверительным* называют интервал, который с заданной надежностью  $\gamma$  покрывает заданный параметр.

1. *Интервальной оценкой (с надежностью  $\gamma$ ) математического ожидания  $a$  нормально распределенного количественного признака  $X$  по выборочной средней  $\bar{x}_B$  при известном среднем квадратическом отклонении  $\sigma$  генеральной совокупности служит доверительный интервал*

$$\bar{x}_B - t(\sigma/\sqrt{n}) < a < \bar{x}_B + t(\sigma/\sqrt{n}),$$

где  $t(\sigma/\sqrt{n}) = \delta$  — точность оценки,  $n$  — объем выборки,  $t$  — значение аргумента функции Лапласа  $\Phi(t)$  (см. приложение 2), при котором  $\Phi(t) = \gamma/2$ ; при неизвестном  $\sigma$  (и объеме выборки  $n < 30$ )

$$\bar{x}_B - t_\gamma(s/\sqrt{n}) < a < \bar{x}_B + t_\gamma(s/\sqrt{n}),$$

где  $s$  — «исправленное» выборочное среднее квадратическое отклонение,  $t_\gamma$  находят по таблице приложения 3 по заданным  $n$  и  $\gamma$ .

Для решаемой задачи: объем выборки  $N=100$ ; надежность  $\gamma=0,95$ ; исправленное квадратическое отклонение  $s=13,85$ . По таблице значений для  $t_\gamma$  находим:  $t_\gamma=1,984$  (таблица значений ниже)

Тогда точность оценки (радиус доверительного интервала):

$$\varepsilon_a = \frac{t_\gamma \cdot s_x}{\sqrt{N}} = \frac{1,984 \cdot 13,85}{\sqrt{100}} = 2,75$$

Таким образом, с вероятностью 0,95 математическое ожидание генеральной совокупности принадлежит интервалу  $(19,68-2,75; 19,68+2,75)$ , т.е.  $16,866 < M(\epsilon) < 22,494$

Таблица значений  $t_{\gamma} = t(\gamma, n)$

$\gamma \backslash n$	0,95	0,99	0,999	$\gamma \backslash n$	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				

Теория: 2. Оценка ср.кв. отклонения.

2. Интервальной оценкой (с надежностью  $\gamma$ ) среднего квадратического отклонения  $\sigma$  нормально распределенного количественного признака  $X$  по «исправленному» выборочному среднему квадратическому отклонению  $s$  служит доверительный интервал

$$s(1-q) < \sigma < s(1+q) \text{ (при } q < 1),$$

$$0 < \sigma < s(1+q) \text{ (при } q > 1),$$

где  $q$  находят по таблице приложения 4 по заданным  $n$  и  $\gamma$ .

Таблица значений  $q = q(\gamma, n)$

$\gamma \backslash n$	0,95	0,99	0,999	$\gamma \backslash n$	0,95	0,99	0,999
5	1,37	2,67	5,64	20	0,37	0,58	0,88
6	1,09	2,01	3,88	25	0,32	0,49	0,73
7	0,92	1,62	2,98	30	0,28	0,43	0,63
8	0,80	1,38	2,42	35	0,26	0,38	0,56
9	0,71	1,20	2,06	40	0,24	0,35	0,50
10	0,65	1,08	1,80	45	0,22	0,32	0,46
11	0,59	0,98	1,60	50	0,21	0,30	0,43
12	0,55	0,90	1,45	60	0,188	0,269	0,38
13	0,52	0,83	1,33	70	0,174	0,245	0,34
14	0,48	0,78	1,23	80	0,161	0,226	0,31
15	0,46	0,73	1,15	90	0,151	0,211	0,29
16	0,44	0,70	1,07	100	0,143	0,198	0,27
17	0,42	0,66	1,01	150	0,115	0,160	0,211
18	0,40	0,63	0,96	200	0,099	0,136	0,185
19	0,39	0,60	0,92	250	0,089	0,120	0,162

По таблице значений для  $q(\gamma, N)$  находим:  $q(0,95; 100) = 0,143$ . Тогда концы доверительного интервала:

$$\alpha_{\sigma} = s(1 - q) = 13,85(1 - 0,143) = 11,867;$$

$$\beta_{\sigma} = s(1 + q) = 13,85(1 + 0,143) = 15,827.$$

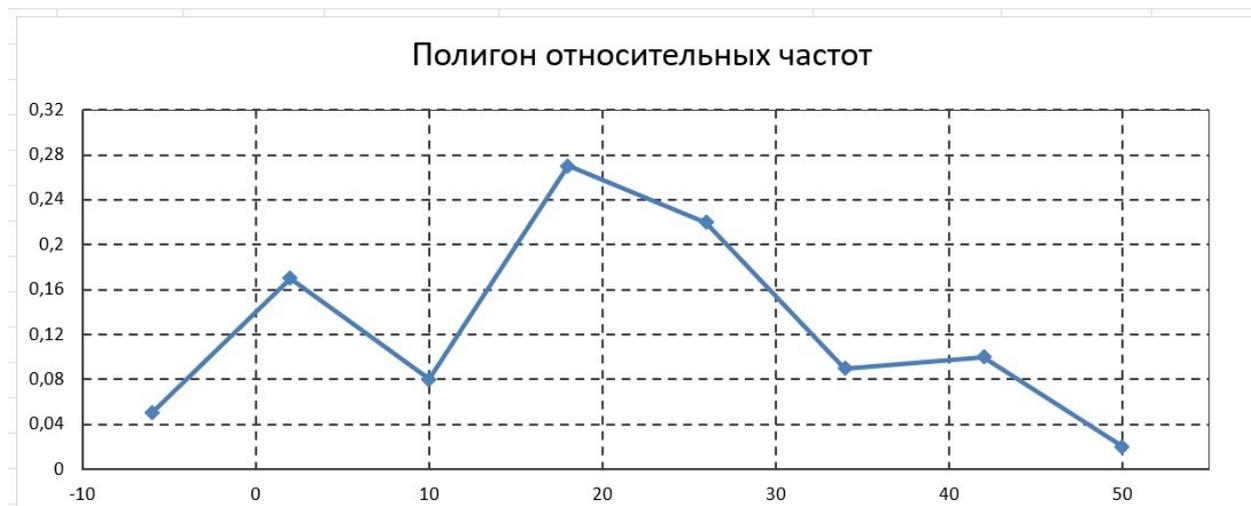
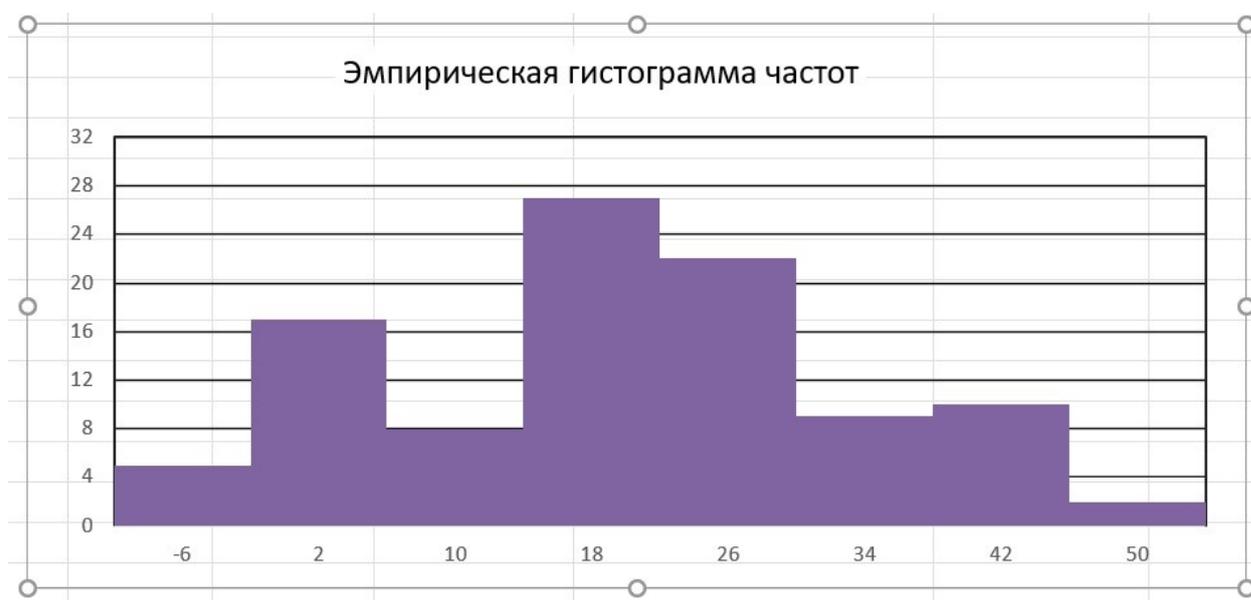
Таким образом, с вероятностью 0,95 среднее квадратическое отклонение генеральной совокупности принадлежит интервалу  $(11,867; 15,827)$ , т. е.

$$11,867 < \sigma(\varepsilon) < 15,827$$

Результаты из Excel:

$\varepsilon_{\alpha}$	$\alpha_{\sigma}$	$\beta_{\sigma}$
2,75	11,867	15,827

#### 4) Гистограмма, полигон частот и эмпирическая функция распределения.



Для построения эмпирической функции распределения, с учетом группировки, лучше использовать не кусочно-постоянную, а непрерывную

кусочно-линейную функцию, где до начала первого интервала и на левом конце первого интервала  $F^* = 0$ , на правом конце первого интервала  $F^i = p_i$  на отрезке  $[\alpha_1; \beta_1]$  строим функцию как отрезок, соединяющий конечные точки. На втором отрезке будет  $F^i(\alpha_2) = F^i(\beta_1) = p_1$ ;

$F^i(\beta_2) = p_1 + p_2$  и т.д. Таким образом, для последнего интервала получим:  $F^i(\beta_k) = 1$  и это же значение  $F^i$  будет иметь для всех значений  $x$ , больших, чем 54. При подобном рассмотрении эмпирической функции распределения вероятностей мы предполагаем, что изменения ее значений происходят не скачком в средних точках интервалов, а накапливаются "плавно" на каждом интервале.

Составим таблицу для построения описанной функции распределения:

	"Запас"	$\alpha_1$	$\alpha_2 = \beta_1$	$\alpha_3 = \beta_2$	$\alpha_4 = \beta_3$	$\alpha_5 = \beta_4$	$\alpha_6 = \beta_5$	$\alpha_7 = \beta_6$	$\alpha_8 = \beta_7$	$\beta_8$	"Запас"
x	-20	-10	-2	6	14	22	30	38	46	54	70
$F^*(x)$	0	0	0,05	0,22	0,3	0,57	0,79	0,88	0,98	1	1



## Задача 2. Проверка статистических гипотез.

### Условие задачи:

По данным из задачи 1:

- 1) Оценить с помощью критерия Пирсона хи-квадрат согласие данных с нормальным распределением при уровне значимости  $\alpha = 0,05$
- 2) Построить сравнительный график эмпирических (ломаная линия с маркерами) и теоретических (сглаженная линия без маркеров) частот. График рисуется на белом фоне линиями разных цветов, добавляются линии сетки координат и легенда.

### Решение:

Оценим с помощью критерия Пирсона хи-квадрат согласие эмпирических данных с нормальным распределением при уровне значимости  $\alpha=0,05$ .

Основная гипотеза  $H_0$ : рассматриваемая генеральная совокупность  $\xi$  распределена по нормальному закону с параметрами  $a=\bar{x}^i=19,68$  и  $\delta=\delta^*=13,78$ ; альтернативная гипотеза  $H_1$ : генеральная совокупность распределена по какому-то другому закону.

Для этого дополним таблицу со сгруппированным статистическим рядом значениями теоретических вероятностей теоретических частот.

Теоретические вероятности попадания значения генеральной совокупности в интервал вычисляются как  $p_i = \varphi(x_i) \cdot \Delta x$ , где  $\varphi(x) = \frac{1}{\sqrt{2\pi}} e^{-\frac{x^2}{2}}$  плотность нормального распределения,  $i=1..k$ . Теоретические частоты вычисляются по формуле  $m_i = [p_i \cdot N]$ , где [...] — округление до целого.

При этом контрольные суммы для теоретических вероятностей и частот должны **приблизительно** равняться 1 и объему выборки.

Вычислив необходимые величины, получим:

Номер интервала	Середина интервала	Эмпирические значения:		Теоритические значения		
		Относительная частота	Частота	Плотность нормального	Вероятность попадания в	Частоты
i	x <sub>i</sub>	p <sub>i</sub> *	m <sub>i</sub> *	φ(x <sub>i</sub> )	p <sub>i</sub>	m <sub>i</sub>
1	-6	0,05	5	0,00510	0,0408	4
2	2	0,17	17	0,01271	0,1017	10
3	10	0,08	8	0,02262	0,1810	18
4	18	0,27	27	0,02874	0,2299	23
5	26	0,22	22	0,02606	0,2085	21
6	34	0,09	9	0,01687	0,1350	13
7	42	0,1	10	0,00780	0,0624	6
8	50	0,02	2	0,00257	0,0206	2
Σ=	100	1	100		0,9798	97

Т.к. критерий Пирсона не "работает" при малых частотах, объединим 1-й интервал со 2-м и 7-й с 8-м. Таким образом, количество интервалов уменьшится до  $k=6$ .

Добавим в таблицу столбец с наблюдаемыми значениями распределения  $\chi^2$ , где

$$\chi_i^2 = \frac{(m_i^2 - m_i)^2}{m_i}, \chi_{набл}^2 = \sum_{i=1}^k \chi_i^2$$

Получим:

Номер интервала	Интервал	Эмпирическая частота	Теоритическая частота	Наблюдаемые значения $\chi^2$	$\chi^2$
i	$(\alpha_i; \beta_i]$	$m_i^*$	$m_i$	$\chi_i^2$	7,81473
1	(-10; 6]	22	14	4,571428571	
2	(6; 14]	8	18	5,555555556	
3	(14; 22]	27	23	0,695652174	
4	(22; 30]	22	21	0,0476	
5	(30; 38]	9	13	1,230769231	
6	(38; 54]	12	8	2	
$\Sigma=$		100	97	14,101	

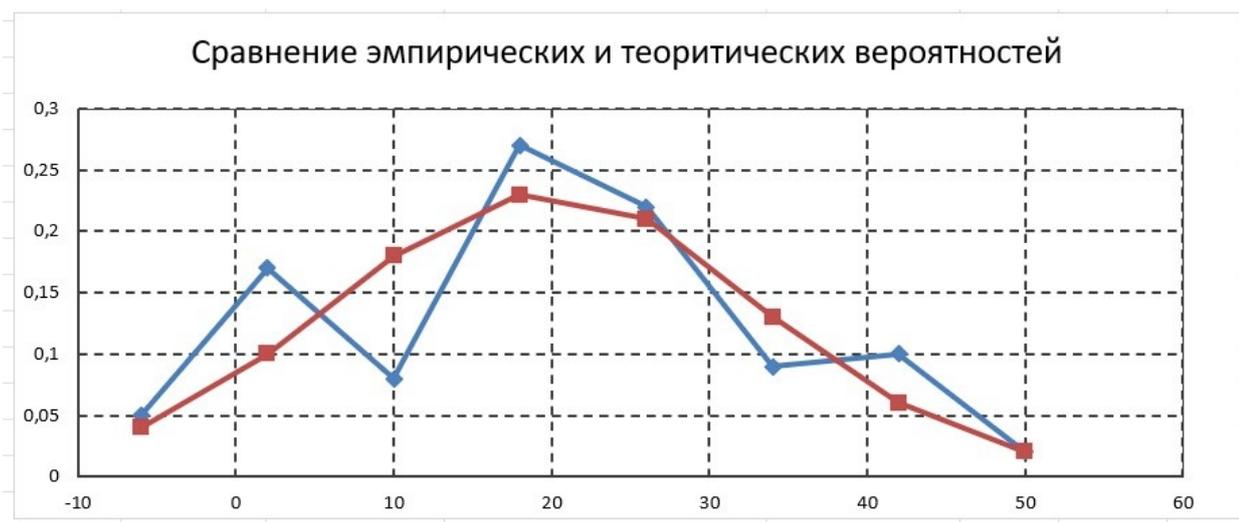
Таким образом,  $\chi_{набл}^2 = 14,101$ . По таблице критических точек распределения  $\chi^2$  найдем  $\chi_{критич}$  при числе степеней свободы  $n = k - 3 = 6 - 3 = 3$  и уровне значимости  $\alpha = 0,05$ :

$$\chi_{критич}^2 = 7,815$$

Сравнивая найденное значение  $\chi_{набл}^2 = 14,101$  с критическим (14,101 > 7,815), определяем, что  $\chi_{набл}^2$  лежит в критической области

Таким образом, при уровне значимости  $\alpha = 0,05$  можно утверждать, что распределение рассматриваемой генеральной совокупности  $\xi$  **сильно отличается от нормальной**

Построим сравнительный график эмпирических (ломаная линия с маркерами) и теоретических (сглаженная линия без маркеров) частот (используя точечную диаграмму).



### Задача 3. Корреляционный и регрессионный анализ данных.

**Условие задачи:**

1) Установить наличие или отсутствие связи между случайными величинами  $X$  и  $Y$ , вычислив выборочный коэффициент корреляции.

2) Найти выборочные регрессии  $Y$  на  $X$  и  $X$  на  $Y$ , предполагая, что они линейные

3) Построить линии регрессий и точки условных средних на одном чертеже (точки условных средних  $X$  и регрессия  $X$  и  $Y$  изображаются одним цветом, а точки условных средних  $Y$  и регрессия  $Y$  на  $X$  – другим цветом).

### Исходные данные:

945									
946	<b>Вар. 106</b>	$X$							
947		<b>11</b>	<b>19</b>	<b>27</b>	<b>35</b>	<b>43</b>	<b>51</b>	<b>59</b>	
948	$Y$	<b>17</b>	109	79	63	47	31	3	1
949		<b>27</b>	78	108	80	64	48	32	16
950		<b>37</b>	62	81	107	82	65	49	33
951		<b>47</b>	46	66	82	114	79	66	65
952		<b>57</b>	16	51	66	71	111	115	82
953		<b>67</b>	2	17	50	54	84	84	131

### Решение:

<b>Вар. 106</b>		$X$							$n_j$	$\sum_{i=1}^7 x_i n_{ji}$	$\bar{x}_{y_j} = \frac{1}{n_j} \sum_{i=1}^7 x_i n_{ji}$
		<b>11</b>	<b>19</b>	<b>27</b>	<b>35</b>	<b>43</b>	<b>51</b>	<b>59</b>			
$Y$	<b>17</b>	109	79	63	47	31	3	1	333	7591	22,80
	<b>27</b>	78	108	80	64	48	32	16	426	11950	28,05
	<b>37</b>	62	81	107	82	65	49	33	479	15221	31,78
	<b>47</b>	46	66	82	114	79	66	65	518	18562	35,83
	<b>57</b>	16	51	66	71	111	115	82	512	20888	40,80
	<b>67</b>	2	17	50	54	84	84	131	422	19210	45,52
$n_i$		313	402	448	432	418	349	328	2690	93422	34,73
$\sum_{j=1}^6 y_j n_{ji}$		9461	14404	18156	18584	19896	18013	18176	116690		
$\bar{y}_{x_i} = \frac{1}{n_i} \sum_{j=1}^6 y_j n_{ji}$		30,23	35,83	40,53	43,02	47,60	51,61	55,41	43,38		

### Математические характеристики $X$ :

1) Математическое ожидание (выборочное среднее)

$$M(X) = \bar{x}^i = \frac{1}{N} \sum_{i=1}^7 x_i n_i = 34,73$$

2) Дисперсия (выборочная дисперсия)

$$D(X) = D^i(x) = \frac{1}{N} \sum_{i=1}^7 x_i^2 n_i - (\bar{x}^i)^2 = 229$$

3) Выборочное среднее квадратическое отклонение

$$\sigma(X) = \sigma_x^i = \sqrt{D^i(x)} = 15,1$$

4) Условные средние, вычисленные в последнем столбце по формуле

$$\bar{x}_{y_j} = \frac{1}{m_j} \sum_{i=1}^4 x_i n_{ji}$$

**Математические характеристики Y:**

1) Математическое ожидание (выборочное среднее)

$$M(Y) = \bar{y}^i = \frac{1}{N} \sum_{i=1}^6 y_i m_i = 43,3792$$

2) Дисперсия (выборочная дисперсия)

$$D(Y) = D^i(y) = \frac{1}{N} \sum_{i=1}^6 y_i^2 m_i - (\bar{y}^i)^2 = 261,24$$

3) Выборочное среднее квадратическое отклонение

$$\sigma(Y) = \sigma_y^i = \sqrt{D^i(y)} = 16,163$$

4) Условные средние, вычисленные в последнем столбце по формуле

$$\bar{y}_{x_j} = \frac{1}{n_j} \sum_{i=1}^6 y_j n_{ji}$$

**Математическое ожидание совместного появления (произведения) значений X и Y:**

$$M(XY) = \frac{1}{N} \sum_{i=1}^6 \sum_{j=1}^7 x_i y_j n_{ji} = \frac{1}{N} \sum_{i=1}^6 \left( x_i \sum_{j=1}^7 y_j n_{ji} \right) = 1622,67$$

**Выборочный коэффициент корреляции:**

$$r_{xy}^i = \frac{M(XY) - \bar{x}^i \bar{y}^i}{\sigma_x^i \sigma_y^i} = 0,4746$$

**Регрессия Y на X** — функция, показывающая зависимость условных средних  $\bar{y}_{x_i}$  от значений  $x_i$   $y = \bar{y}^i + r_{xy}^i \frac{\sigma_y^i}{\sigma_x^i} (x - \bar{x}^i) \Rightarrow$

$$y = 43,3792 + 0,4746 \cdot \frac{16,163}{15,1} (x - 34,73) \Rightarrow y = 43,3792 + 0,51x - 17,64 \Rightarrow$$

$$y_{\text{рег}} = 25,739 + 0,51x$$

Составим таблицу значений регрессии Y на X (для известных X)

$x_i$	11	19	27	35	43	51	59
$y_{пер}$	31,36	35,41	39,46	43,52	47,57	51,62	55,67

Уравнение линейной регрессии X на Y:  $x = \bar{x}^i + r_{xy}^i \frac{\sigma_x^i}{\sigma_y^i} (y - \bar{y}^i) \Rightarrow$

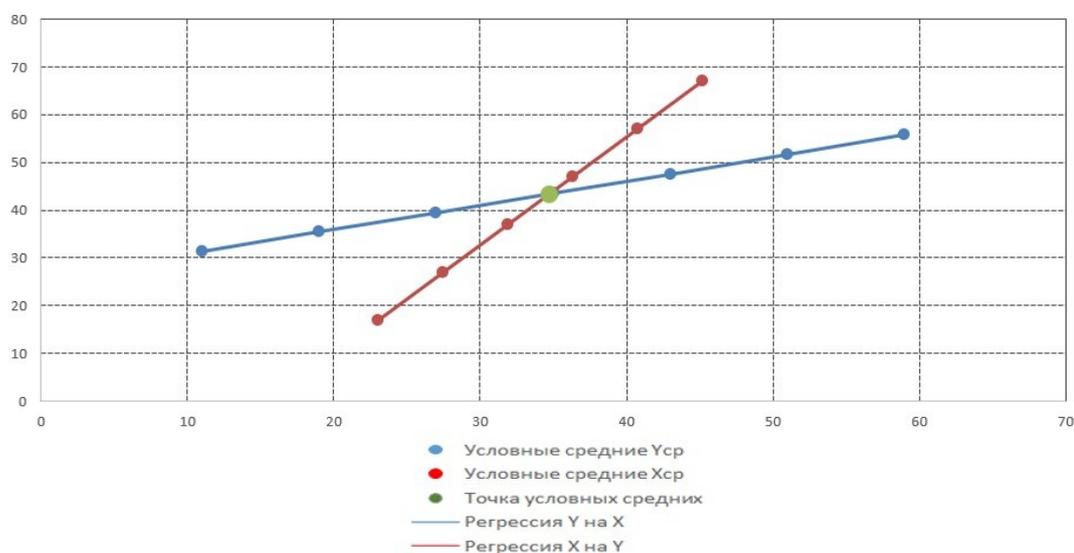
$$x = 34,73 + 0,4746 \cdot \frac{15,1}{16,163} (y - 43,3792) \Rightarrow x_{пер} = 34,73 + 0,443 y - 19,234 \Rightarrow$$

$$x_{пер} = 15,496 + 0,443 y$$

Составим таблицу значений регрессии X на Y (для известных Y)

$y_j$	17	27	37	47	57	67
$x_{пер}$	23,00	27,45	31,89	36,34	40,78	45,2

Сравнение регрессий:



**Анализ корреляции:**  $r_{xy}^i \approx 0,4746$  — коэфф-т корреляции  $\neq 0$  (существенно больше нуля)  $\Rightarrow$  связь между X и Y есть, и существенная (величины X и Y коррелированы), при этом значение  $r_{xy}$  существенно меньше 1  $\Rightarrow$  связь нелинейная (существенно нелинейная).