

МИНИСТЕРСТВО ОБРАЗОВАНИЯ И НАУКИ РОССИЙСКОЙ ФЕДЕРАЦИИ  
ФЕДЕРАЛЬНОЕ ГОСУДАРСТВЕННОЕ АВТОНОМНОЕ ОБРАЗОВАТЕЛЬНОЕ УЧРЕЖДЕНИЕ ВЫСШЕГО ПРОФЕССИОНАЛЬНОГО  
ОБРАЗОВАНИЯ

**Национальный исследовательский ядерный университет «МИФИ»**

---



**Институт интеллектуальных кибернетических систем**

**КАФЕДРА КИБЕРНЕТИКИ**

**БДЗ**

**по курсу "Математическая статистика"**

**студента группы Б22-534**

**Когановского Григория**

**Вариант №7**

**Оценка:** \_\_\_\_\_

**Подпись:** \_\_\_\_\_

**2024 г.**

## 1. Описательные статистики

### 1.1. Выборочные характеристики

Анализируемый признак 1 – C9 (Number of alcoholic drinks consumed per week)

Анализируемый признак 2 – C10 (Cholesterol consumed (mg per day))

Анализируемый признак 3 – C11 (Dietary beta-carotene consumed (mcg per day))

*а) Привести формулы расчёта выборочных характеристик*

Выборочная хар-ка	Формула расчета
Объём выборки	$n$
Среднее	$\bar{x} = \frac{1}{n} \sum_{i=1}^n x_i$
Выборочная дисперсия	$D_X^* = \frac{1}{n} \sum_{i=1}^n (x_i - \bar{x})^2$
Выборочное среднеквадратическое отклонение	$\sigma_X^* = \sqrt{D_X^*} = \sqrt{\frac{1}{n} \sum_{i=1}^n (x_i - \bar{x})^2}$
Выборочный коэффициент асимметрии	$\gamma_X^* = \frac{\mu_{3,X}^*}{(\sigma_X^*)^3} = \frac{\frac{1}{n} \sum_{i=1}^n (x_i - \bar{x})^3}{\left(\frac{1}{n} \sum_{i=1}^n (x_i - \bar{x})^2\right)^{\frac{3}{2}}}$
Выборочный эксцесс	$\epsilon_X^* = \frac{\mu_{4,X}^*}{(\sigma_X^*)^4} = \frac{\frac{1}{n} \sum_{i=1}^n (x_i - \bar{x})^4}{\left(\frac{1}{n} \sum_{i=1}^n (x_i - \bar{x})^2\right)^2}$

*б) Рассчитать выборочные характеристики*

Выборочная хар-ка	Признак 1	Признак 2	Признак 3
Среднее	3.28	242.46	2185.60
Выборочная дисперсия	151.37	17366.48	2165445.23
Выборочное среднеквадратическое отклонение	12.30	131.78	1471.55
Выборочный коэффициент асимметрии	13.76	1.47	1.61
Выборочный эксцесс	217.82	3.34	3.40

*1.2. Группировка и гистограммы частот*

Анализируемый признак – C9 (Number of alcoholic drinks consumed per week)

Объём выборки –  $n = 315$

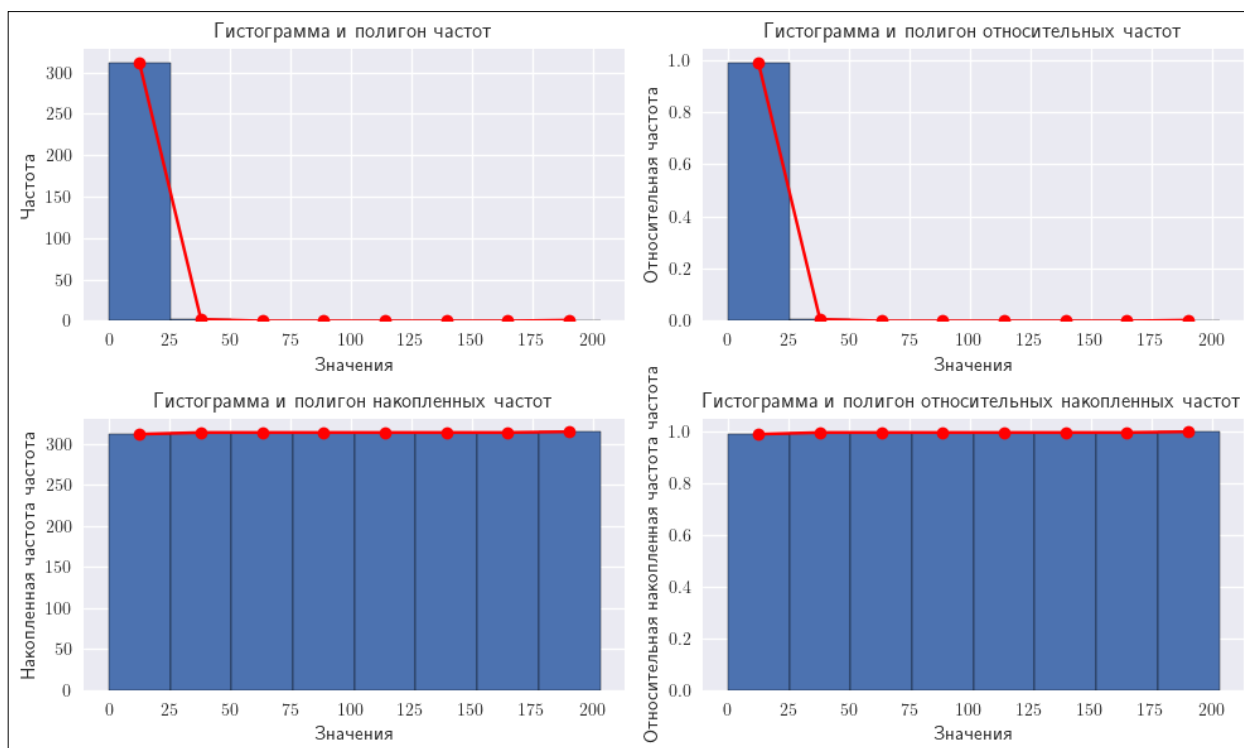
*а) Выбрать число групп*

Число групп	Обоснование выбора числа групп	Ширина интервалов
8	Формула Стерджесса: $k \approx 1 + 1,3 \ln n$	от 25.38 до 25.58

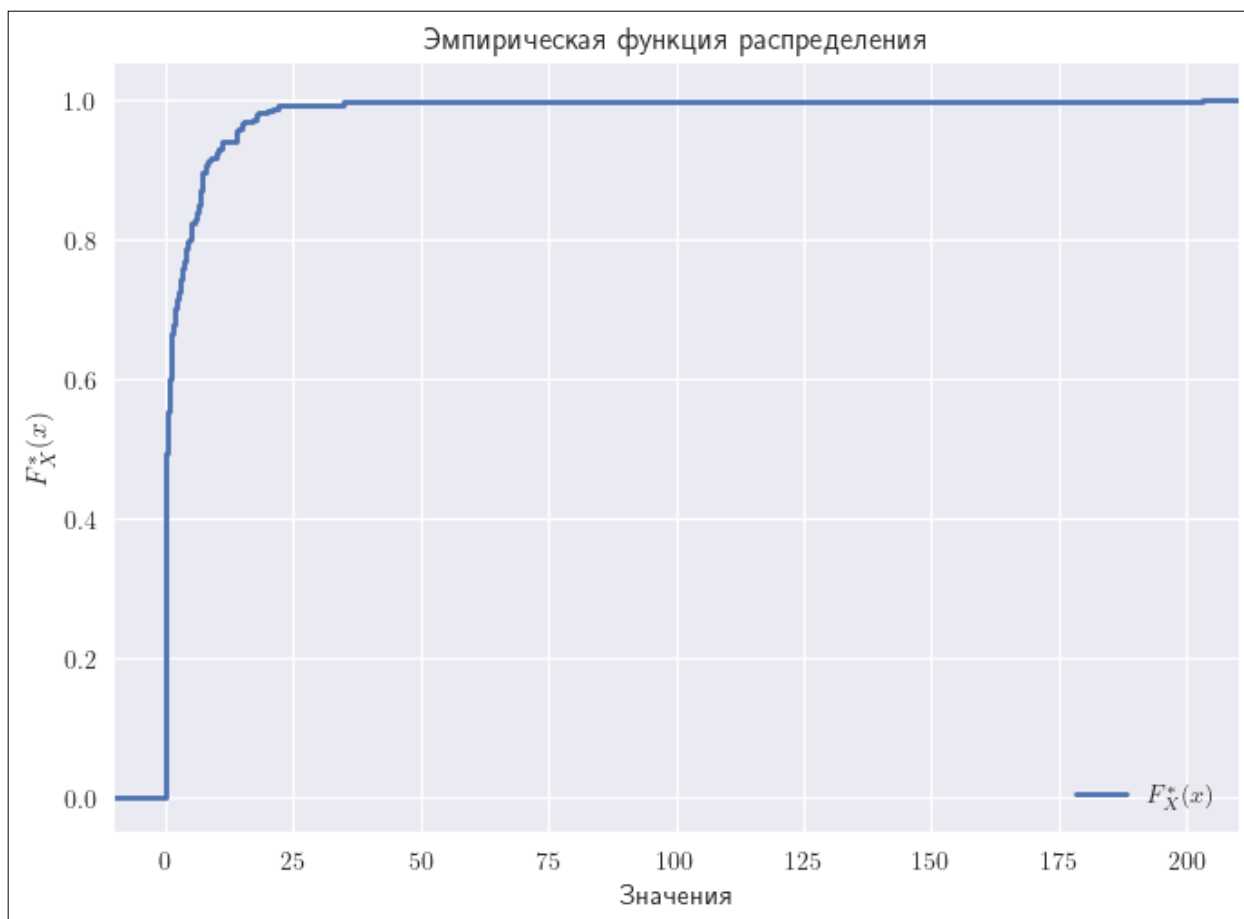
*б) Построить таблицу частот*

Номер интервала	Нижняя граница	Верхняя граница	Частота	Относит. частота	Накопл. частота	Относит. накопл. частота
1	-0.2	25.38	312	0.99	312	0.99
2	25.38	50.75	2	0.01	314	1.00
3	50.75	76.12	0	0.00	314	1.00
4	76.12	101.50	0	0.00	314	1.00
5	101.50	126.88	0	0.00	314	1.00
6	126.88	152.25	0	0.00	314	1.00
7	152.25	177.62	0	0.00	314	1.00
8	177.62	203.0	1	0.00	315	1.00

в) Построить гистограммы частот и полигоны частот



г) Построить график эмпирической функции распределения



## 2. Интервальные оценки

### 2.1. Доверительные интервалы для мат. ожидания

Анализируемый признак – C9 (Number of alcoholic drinks consumed per week)

Объём выборки –  $n = 315$

Оцениваемый параметр –  $m$

#### а) Привести формулы расчёта доверительных интервалов

Граница доверительного интервала	Формула расчета
Нижняя граница	$\bar{X} - \frac{S}{\sqrt{n}} \cdot t_{1-\frac{\alpha}{2},(n-1)}$
Верхняя граница	$\bar{X} + \frac{S}{\sqrt{n}} \cdot t_{1-\frac{\alpha}{2},(n-1)}$

#### б) Рассчитать доверительные интервалы

Граница доверительного интервала	$\alpha = 0.01$	$\alpha = 0.05$	$\alpha = 0.1$
Нижняя граница	1.48	1.91	2.13
Верхняя граница	5.08	4.65	4.42

### 2.2. Доверительные интервалы для дисперсии

Анализируемый признак – C9 (Number of alcoholic drinks consumed per week)

Объём выборки –  $n = 315$

Оцениваемый параметр –  $\sigma^2$

#### а) Привести формулы расчёта доверительных интервалов

Граница доверительного интервала	Формула расчета
Нижняя граница	$\frac{(n-1) \cdot S^2}{\chi^2_{1-\frac{\alpha}{2},(n-1)}}$
Верхняя граница	$\frac{(n-1) \cdot S^2}{\chi^2_{\frac{\alpha}{2},(n-1)}}$

*б) Рассчитать доверительные интервалы*

Граница доверительного интервала	$\alpha = 0.01$	$\alpha = 0.05$	$\alpha = 0.1$
Нижняя граница	124.72	130.64	133.82
Верхняя граница	188.31	178.72	174.05

*2.3. Доверительные интервалы для разности мат. ожиданий*

Анализируемый признак 1 – C11 (Dietary beta-carotene consumed (mcg per day))

Анализируемый признак 2 – C12 (Dietary retinol consumed (mcg per day))

Объёмы выборок –  $n_1 = n_2 = 315$

Оцениваемый параметр –  $m_1 - m_2$

*а) Привести формулы расчёта доверительных интервалов*

Граница доверительного интервала	Формула расчета
Нижняя граница	$(\bar{X}_1 - \bar{X}_2) - t_{1-\frac{\alpha}{2}, (n_1+n_2-2)} \cdot S \cdot \sqrt{\frac{1}{n_1} + \frac{1}{n_2}}$
Верхняя граница	$(\bar{X}_1 - \bar{X}_2) + t_{1-\frac{\alpha}{2}, (n_1+n_2-2)} \cdot S \cdot \sqrt{\frac{1}{n_1} + \frac{1}{n_2}}$
$S$	$\sqrt{\frac{(n_1 - 1)S_1^2 + (n_2 - 1)S_2^2}{n_1 + n_2 - 2}}$

*б) Рассчитать доверительные интервалы*

Граница доверительного интервала	$\alpha = 0.01$	$\alpha = 0.05$	$\alpha = 0.1$
Нижняя граница	1121.82	1177.26	1205.56
Верхняя граница	1583.96	1528.52	1500.21

#### 2.4. Доверительные интервалы для отношения дисперсий

Анализируемый признак 1 – C11 (Dietary beta-carotene consumed (mcg per day))

Анализируемый признак 2 – C12 (Dietary retinol consumed (mcg per day))

Объёмы выборок –  $n_1 = n_2 = 315$

Оцениваемый параметр –  $\frac{\sigma_1^2}{\sigma_2^2}$

а) Привести формулы расчёта доверительных интервалов

Граница доверительного интервала	Формула расчета
Нижняя граница	$\frac{S_1^2}{S_2^2} \cdot F_{\frac{\alpha}{2}, (n_1-1, n_2-1)}$
Верхняя граница	$\frac{S_1^2}{S_2^2} \cdot F_{1-\frac{\alpha}{2}, (n_1-1, n_2-1)}$

б) Рассчитать доверительные интервалы

Граница доверительного интервала	$\alpha = 0.01$	$\alpha = 0.05$	$\alpha = 0.1$
Нижняя граница	4.67	5.01	5.19
Верхняя граница	8.37	7.81	7.53

### 3. Проверка статистических гипотез о математических ожиданиях и дисперсиях

#### 3.1. Проверка статистических гипотез о математических ожиданиях

Анализируемый признак – C9 (Number of alcoholic drinks consumed per week)

Объём выборки –  $n = 315$

Статистическая гипотеза –  $H_0 : m = m_0$   
 $H' : m \neq m_0$

а) Указать формулы расчёта показателей, используемых при проверке статистических гипотез

	Выражение
Формула расчета статистики критерия	$Z = \frac{\bar{X} - m_0}{S/\sqrt{n}}$
Закон распределения статистики критерия при условии истинности основной гипотезы	$T(n - 1)$
Формулы расчета критических точек	$\pm t_{1-\frac{\alpha}{2}, n-1}$
Формула расчета $p$ -value	$2 \min(F_Z(z_{\text{выб}}   H_0), 1 - F_Z(z_{\text{выб}}   H_0))$

б) Выбрать произвольные значения  $m_0$  и проверить статистические гипотезы

$m_0$	Уровень значимости	Выборочное значение статистики критерия	$p$ -value	Статистическое решение	Вывод
0	0.1	4.72	0.00	$H_0$ отклоняется	$m \neq 0$
2	0.1	1.84	0.07	$H_0$ отклоняется	$m \neq 2$
4	0.1	-1.04	0.30	$H_0$ принимается	$m = 4$

#### 3.2. Проверка статистических гипотез о дисперсиях

Анализируемый признак – C9 (Number of alcoholic drinks consumed per week)

Объём выборки –  $n = 315$

Статистическая гипотеза –  $H_0 : \sigma = \sigma_0$   
 $H' : \sigma \neq \sigma_0$



а) Указать формулы расчёта показателей, используемых при проверке статистических гипотез

	Выражение
Формула расчета статистики критерия	$Z = \frac{(n-1)S^2}{\sigma_0^2}$
Закон распределения статистики критерия при условии истинности основной гипотезы	$\chi^2(n-1)$
Формулы расчета критических точек	$\chi_{\frac{\alpha}{2}, n-1}^2, \chi_{1-\frac{\alpha}{2}, n-1}^2$
Формула расчета $p$ -value	$2 \min(F_Z(z_{\text{выб}}   H_0), 1 - F_Z(z_{\text{выб}}   H_0))$

б) Выбрать произвольные значения  $\sigma_0$  и проверить статистические гипотезы

$\sigma_0$	Уровень значимости	Выборочное значение статистики критерия	$p$ -value	Статистическое решение	Вывод
10	0.1	476.82	0.00	$H_0$ отклоняется	$\sigma \neq 10$
12.5	0.1	305.16	0.74	$H_0$ принимается	$\sigma = 12.5$
15	0.1	211.92	0.00	$H_0$ отклоняется	$\sigma \neq 15$

### 3.3. Проверка статистических гипотез о равенстве математических ожиданий

Анализируемый признак 1 – C11 (Dietary beta-carotene consumed (mcg per day))

Анализируемый признак 2 – C12 (Dietary retinol consumed (mcg per day))

Объёмы выборок –  $n_1 = n_2 = 315$

Статистическая гипотеза –  $H_0: m_1 = m_2$   
 $H': m_1 \neq m_2$

а) Указать формулы расчёта показателей, используемых при проверке статистических гипотез

	Выражение
Формула расчета статистики критерия	$Z = \frac{\bar{X}_1 - \bar{X}_2}{S \sqrt{\frac{1}{n_1} + \frac{1}{n_2}}}, \text{ где } S^2 = \frac{(n_1 - 1)S_1^2 + (n_2 - 1)S_2^2}{n_1 + n_2 - 2}$
Закон распределения статистики критерия при условии истинности основной гипотезы	$T(n_1 + n_2 - 2)$
Формулы расчета критических точек	$\pm t_{1-\frac{\alpha}{2}, n_1+n_2-2}$
Формула расчета $p$ -value	$2 \min(F_Z(z_{\text{выб}}   H_0), 1 - F_Z(z_{\text{выб}}   H_0))$

*б) Проверить статистические гипотезы*

Уровень значимости	Выборочное значение статистики критерия	$p$ -value	Статистическое решение	Вывод
0.01	15.13	0.00	$H_0$ отклоняется	$m_1 \neq m_2$
0.05			$H_0$ отклоняется	$m_1 \neq m_2$
0.1			$H_0$ отклоняется	$m_1 \neq m_2$

*3.4. Проверка статистических гипотез о равенстве дисперсий*

Анализируемый признак 1 – C11 (Dietary beta-carotene consumed (mcg per day))

Анализируемый признак 2 – C12 (Dietary retinol consumed (mcg per day))

Объёмы выборок –  $n_1 = n_2 = 315$

Статистическая гипотеза –  $H_0 : \sigma_1 = \sigma_2$   
 $H' : \sigma_1 \neq \sigma_2$

а) Указать формулы расчёта показателей, используемых при проверке статистических гипотез

	Выражение
Формула расчета статистики критерия	$Z = \frac{S_1^2}{S_2^2}$
Закон распределения статистики критерия при условии истинности основной гипотезы	$F(n_1 - 1, n_2 - 1)$
Формулы расчета критических точек	$F_{\frac{\alpha}{2}, n_1-1, n_2-1}, F_{1-\frac{\alpha}{2}, n_1-1, n_2-1}$
Формула расчета $p$ -value	$2 \min(F_Z(z_{\text{выб}}   H_0), 1 - F_Z(z_{\text{выб}}   H_0))$

б) Проверить статистические гипотезы

Уровень значимости	Выборочное значение статистики критерия	$p$ -value	Статистическое решение	Вывод
0.01	6.26	0.00	$H_0$ отклоняется	$\sigma_1 \neq \sigma_2$
0.05			$H_0$ отклоняется	$\sigma_1 \neq \sigma_2$
0.1			$H_0$ отклоняется	$\sigma_1 \neq \sigma_2$

#### 4. Критерии согласия

Анализируемый признак – C9 (Number of alcoholic drinks consumed per week)

Объём выборки –  $n = 315$

##### 4.1. Критерий хи-квадрат

Теоретическое распределение – нормальное.

Статистическая гипотеза –  $H_0 : X \sim N$   
 $H' : X \not\sim N$

а) Указать формулы расчета показателей, используемых при проверке статистических гипотез

	Выражение	Пояснение использованных обозначений
Формула расчета статистики критерия	$Z = \sum_{i=1}^k \frac{(n_i - n \cdot p_i)^2}{n \cdot p_i}$	$k$ - число интервалов в группированном статистическом ряду. $n_i$ - частота попадания случайной величины $X$ в интервал $\Delta_i$ . $p_i$ - вероятность попадания случайной величины $X$ в интервал $\Delta_i$ в условиях $H_0$ , то есть, если $\Delta_i = (a_{i-1}, a_i]$ , то $p_i = \int_{a_{i-1}}^{a_i} g(x)dx = G(a_i) - G(a_{i-1}).$
Закон распределения статистики критерия при условии истинности основной гипотезы	$\chi^2(k - r - 1)$	$r$ - количество оцениваемых параметров у предполагаемого распределения $G$ .
Формула расчета критической точки	$\chi^2_{1-\alpha, k-r-1}$	Малые значения $Z$ нам также подходят, поэтому критическая область выбирается правосторонней
Формула расчета $p$ -value	$1 - F_Z(z_{\text{выб}}   H_0)$	

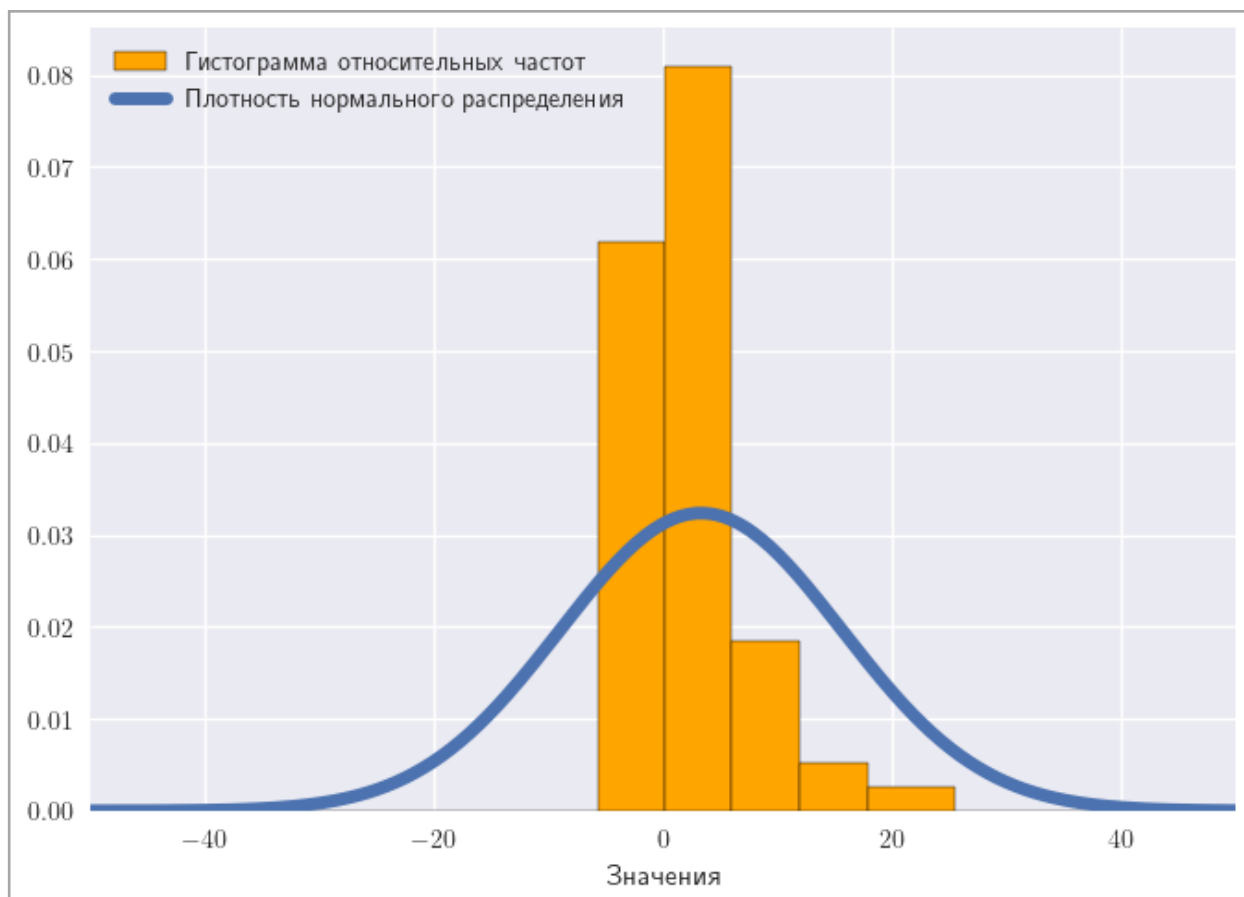
б) Выбрать число групп

Число групп	Обоснование выбора числа групп	Ширина интервалов
8	$k \approx 1 + 1.3 \ln n$ - формула Стерджесса $n \cdot p_i \gtrsim 5$ - поправка на чувствительность критерия	От 5.69 до $\infty$

в) Построить таблицу частот

Номер интервала	Нижняя граница	Верхняя граница	Частота	Относит. частота	Вероятность попадания в интервал при условии истинности основной гипотезы
1	$-\infty$	-11.83	0	0	0.11
2	-11.83	-5.65	0	0	0.12
3	-5.65	0.04	111	0,35	0.16
4	0.04	5.92	150	0,48	0.19
5	5.92	11.91	35	0,11	0.17
6	11.91	17.91	10	0,03	0.12
7	17.91	25.38	6	0,02	0.08
8	25.38	$+\infty$	3	0,01	0.04

г) Построить гистограмму относительных частот и функцию плотности теоретического распределения на одном графике



д) Проверить статистические гипотезы

Уровень значимости	Выборочное значение статистики критерия	$p$ -value	Статистическое решение	Вывод
0.01	332.48	0.00	$H_0$ отклоняется	$X \sim N$
0.05			$H_0$ отклоняется	$X \sim N$
0.1			$H_0$ отклоняется	$X \sim N$

4.2. Проверка гипотезы о нормальности на основе коэффициента асимметрии и эксцесса (критерий Харке-Бера)

Статистическая гипотеза –  $H_0 : X \sim N$   
 $H' : X \not\sim N$

а) Указать формулы расчёта показателей, используемых при проверке статистических гипотез

	Выражение	Пояснение использованных обозначений
Формула расчета статистики критерия	$Z = Z_1^2 + Z_2^2, \text{ где}$ $Z_1 = \frac{\gamma_X^*}{\sqrt{\frac{6}{n}}};$ $Z_2 = \frac{\epsilon_X^*}{\sqrt{\frac{24}{n}}}$	$\gamma_X^*$ - выборочный коэффициент асимметрии. $\epsilon_X^*$ - выборочный эксцесс. $n$ - объём выборки.
Закон распределения статистики критерия при условии истинности основной гипотезы	$\chi^2(2)$	
Формула расчета критической точки	$\chi_{1-\alpha, 2}^2$	Правосторонняя критическая область
Формула расчета $p$ -value	$1 - F_Z(z_{\text{выб}}   H_0)$	

б) Проверить статистические гипотезы

Уровень значимости	Выборочное значение статистики критерия	<i>p-value</i>	Статистическое решение	Вывод
0.01	632630.74	0.00	$H_0$ отклоняется	$X \approx N$
0.05			$H_0$ отклоняется	$X \approx N$
0.1			$H_0$ отклоняется	$X \approx N$

Вывод (в терминах предметной области)

В результате проведённого в п.4 статистического анализа обнаружено, что оба критерия согласия отвергают гипотезу о нормальности распределения величины С9, поэтому можно сделать вывод, что С9 не имеет нормального распределения.

## 5. Проверка однородности выборок

Анализируемый признак 1 – C11 (Dietary beta-carotene consumed (mcg per day))

Анализируемый признак 2 – C12 (Dietary retinol consumed (mcg per day))

Объёмы выборок –  $n_1 = n_2 = 315$

### 5.1 Критерий знаков

Статистическая гипотеза –  
 $H_0 : F_X(x) = F_Y(x)$   
 $H' : F_X(x) \neq F_Y(x)$

а) Указать формулы расчёта показателей, используемых при проверке статистических гипотез

	Выражение	Пояснение использованных обозначений
Формула расчета статистики критерия	$Z = \frac{K_+ - \frac{n}{2}}{\sqrt{n}/2}$	$K_+$ - число знаков '+' в выборке $z_1, \dots, z_n = x_1 - y_1, \dots, x_n - y_n$
Закон распределения статистики критерия при условии истинности основной гипотезы	$N(0,1)$	
Формула расчета критической точки	$\pm u_{1-\frac{\alpha}{2}}$	Двусторонняя критическая область
Формула расчета $p$ -value	$2 \min(F_Z(z_{\text{выб}}   H_0), 1 - F_Z(z_{\text{выб}}   H_0))$	

б) Проверить статистические гипотезы

Уровень значимости	Выборочное значение статистики критерия	$p$ -value	Статистическое решение	Вывод
0.01	14.03	0.00	$H_0$ отклоняется	$F_X(x) \neq F_Y(x)$
0.05			$H_0$ отклоняется	$F_X(x) \neq F_Y(x)$
0.1			$H_0$ отклоняется	$F_X(x) \neq F_Y(x)$



## 5.2. Критерий хи-квадрат

Статистическая гипотеза –  $H_0 : F_X(x) = F_Y(x)$   
 $H' : F_X(x) \neq F_Y(x)$

а) Указать формулы расчёта показателей, используемых при проверке статистических гипотез

	Выражение	Пояснение использованных обозначений
Формула расчета статистики критерия	$Z = nm \sum_{i=1}^k \frac{1}{n_i + m_i} \left( \frac{n_i}{n} + \frac{m_i}{m} \right)$	<p><math>k</math> - число интервалов в группированном статистическом ряду.</p> <p><math>n_i</math> - частота попадания случайной величины <math>X</math> в интервал <math>\Delta_i</math>.</p> <p><math>m_i</math> - частота попадания случайной величины <math>Y</math> в интервал <math>\Delta_i</math>.</p> <p><math>n</math> - объём выборки <math>x_1, \dots, x_n</math></p> <p><math>m</math> - объём выборки <math>y_1, \dots, y_m</math></p>
Закон распределения статистики критерия при условии истинности основной гипотезы	$\chi^2(k - 1)$	
Формула расчета критической точки	$\chi^2_{1-\alpha, k-1}$	Правосторонняя критическая область
Формула расчета $p$ -value	$1 - F_Z(z_{\text{выб}}   H_0)$	

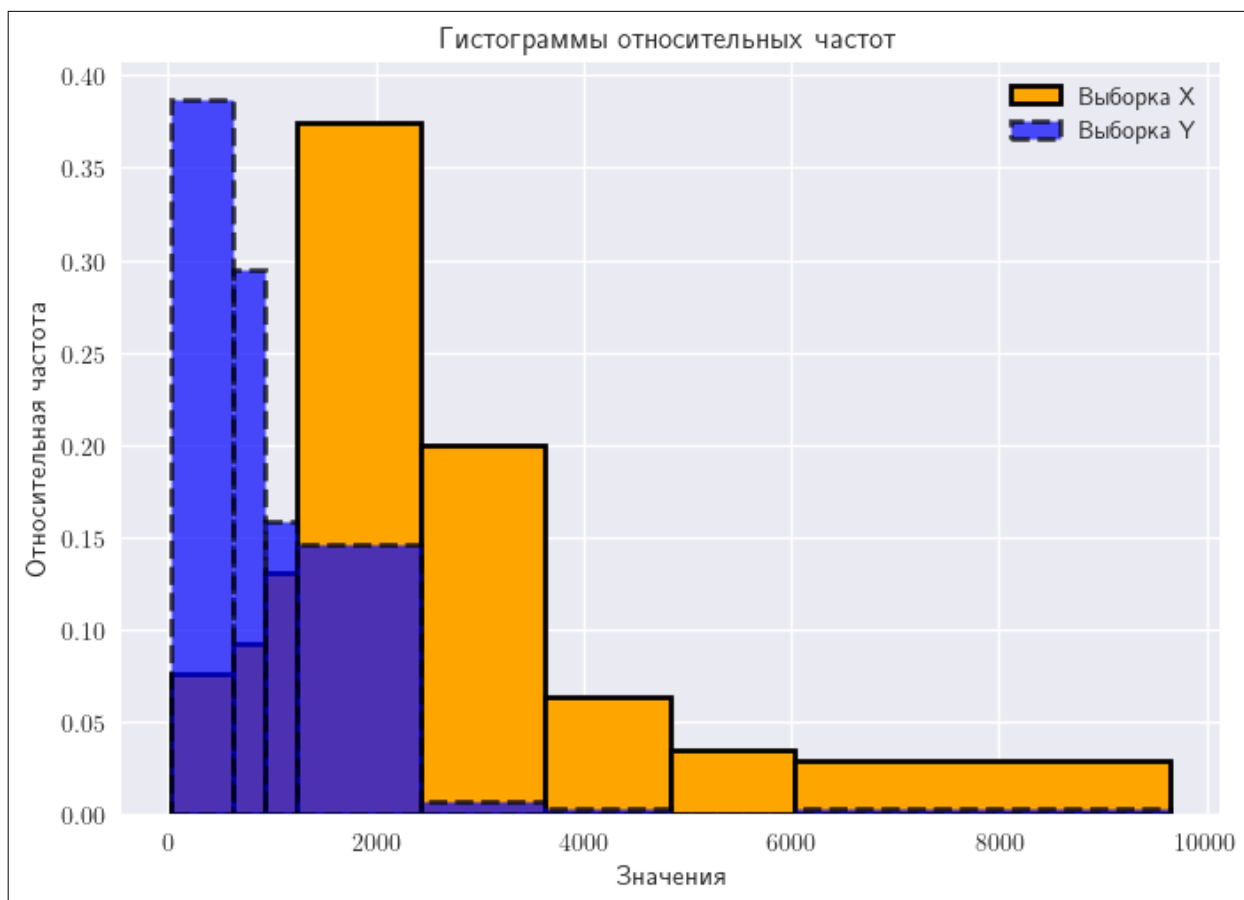
б) Выбрать число групп

Число групп	Обоснование выбора числа групп	Ширина интервалов
8	$k \approx 1 + 1.3 \ln \max(n, m)$ - формула Стерджесса $n_i + m_i \gtrsim 5$ - поправка на чувствительность критерия	от 300.38 до 3604.50

в) Построить таблицу частот

Номер интервала	Нижняя граница	Верхняя граница	Частота признака 1	Частота признака 2	Относит. частота признака 1	Относит. частота признака 2
1	30.00	630.75	24	122	0,08	0,39
2	630.75	931.12	29	93	0,09	0,3
3	931.12	1231.50	41	50	0,13	0,16
4	1231.50	2433.0	118	46	0,37	0,15
5	2433.0	3634.50	63	2	0,2	0,01
6	3634.50	4836.00	20	1	0,06	0
7	4836.00	6037.50	11	0	0,03	0
8	6037.50	9642.00	9	1	0,03	0

г) Построить гистограммы относительных частот на одном графике



д) Проверить статистические гипотезы

Уровень значимости	Выборочное значение статистики критерия	$p$ -value	Статистическое решение	Вывод
0.01	223.69	0.00	$H_0$ отклоняется	$F_X(x) \neq F_Y(x)$
0.05			$H_0$ отклоняется	$F_X(x) \neq F_Y(x)$
0.1			$H_0$ отклоняется	$F_X(x) \neq F_Y(x)$

Вывод (в терминах предметной области)

В результате проведённого в п.5 статистического анализа обнаружено, что выборки С11 и С12 неоднородны.

## 6. Таблицы сопряжённости

Факторный признак  $x$  – C2 (Sex)

Результативный признак  $y$  – C5 (Vitamin Use)

Объёмы выборок –  $n_1 = n_2 = n = 315$

Статистическая гипотеза –  $H_0 : F_Y(y |_{X=x^{(1)}}) = F_Y(y |_{X=x^{(2)}}) = \dots = F_Y(y |_{X=x^{(k_1)}}) = F_Y(y)$   
 $H' : \exists i, j : F_Y(y |_{X=x^{(i)}}) \neq F_Y(y |_{X=x^{(j)}})$

*а) Указать формулы расчёта показателей, используемых при проверке статистических гипотез*

	Выражение	Пояснение использованных обозначений
Формула расчета статистики критерия	$Z = \sum_{i=1}^{k_1} \sum_{j=1}^{k_2} \frac{(n_{ij} - m_{ij})^2}{m_{ij}}$	$(x_1, y_1), \dots, (x_n, y_n)$ - наблюдения случайного вектора $(X, Y)$ , где $X, Y$ - случайные величины дискретного типа $x^{(1)}, \dots, x^{(k_1)}$ - варианты признака $X$ $y^{(1)}, \dots, y^{(k_2)}$ - варианты признака $Y$ $n_{ij}$ - выборочная частота варианта $(x^{(i)}, y^{(j)})$ в выборке $(x_1, y_1), \dots, (x_n, y_n)$ $m_{ij}$ - теоретическая частота варианта $(x^{(i)}, y^{(j)})$ в выборке $(x_1, y_1), \dots, (x_n, y_n)$ при условии истинности $H_0$
Закон распределения статистики критерия при условии истинности основной гипотезы	$\chi^2((k_1 - 1)(k_2 - 1))$	
Формула расчета критической точки	$\chi^2_{1-\alpha, (k_1-1)(k_2-1)}$	Правосторонняя критическая область
Формула расчета $p$ -value	$1 - F_Z(z_{\text{выб}}   H_0)$	

б) Построить эмпирическую таблицу сопряжённости

Вариант Y Вариант X	$y^{(1)}$	$y^{(2)}$	$y^{(3)}$	$\Sigma$
$x^{(1)}$	87	77	109	273
$x^{(2)}$	24	5	13	42
$\Sigma$	111	82	122	315

в) Построить теоретическую таблицу сопряжённости

Вариант Y Вариант X	$y^{(1)}$	$y^{(2)}$	$y^{(3)}$	$\Sigma$
$x^{(1)}$	96,20	71,07	105,73	273
$x^{(2)}$	14,80	10,93	16,27	42
$\Sigma$	111,00	82	122	315

г) Проверить статистические гипотезы

Уровень значимости	Выборочное значение статистики критерия	$p$ -value	Статистическое решение	Вывод
0.01	11.07	0.00	$H_0$ отклоняется	$\exists i, j : F_Y(y  _{X=x^{(i)}}) \neq F_Y(y  _{X=x^{(j)}})$
0.05			$H_0$ отклоняется	$\exists i, j : F_Y(y  _{X=x^{(i)}}) \neq F_Y(y  _{X=x^{(j)}})$
0.1			$H_0$ отклоняется	$\exists i, j : F_Y(y  _{X=x^{(i)}}) \neq F_Y(y  _{X=x^{(j)}})$

Вывод (в терминах предметной области)

В результате проведённого в п.6 статистического анализа обнаружено, что между факторным признаком С2 и результативным признаком С5 присутствует статистическая связь. Под действием С2 оказывается влияние на распределение С5.

## 7. Дисперсионный анализ

Факторный признак  $x$  – C5 (Vitamin Use)

Результативный признак  $y$  – C4 (Quetelet (weight/height^2))

Число вариантов факторного признака –  $k = 3$

Объёмы выборок –  $n_1 = n_2 = n = 315$

Статистическая гипотеза –  $H_0 : F_Y(y |_{X=x_1}) = F_Y(y |_{X=x_2}) = \dots = F_Y(y |_{X=x_k}) = F_Y(y)$   
 $H' : \exists i, j : F_Y(y |_{X=x_i}) \neq F_Y(y |_{X=x_j})$

*а) Рассчитать групповые выборочные характеристики*

№ п/п	Вариант факторного признака	Объём выборки	Групповые средние	Групповые дисперсии
1	No	111	26.53	34.21
2	Not often	82	26.63	41.24
3	Often	122	25.51	33.53

*б) Привести формулы расчёта показателей вариации, используемых в дисперсионном анализе*

Источник вариации	Показатель вариации	Число степеней свободы	Несмещенная оценка
Факторный признак	$\tilde{D}_{\text{межгр}} = \frac{1}{n} \sum_{i=1}^k n_i (\bar{y}_i - \bar{y})^2$	$k - 1$	$\frac{n}{k - 1} \tilde{D}_{\text{межгр}}$
Остаточные признаки	$\tilde{D}_{\text{внутр}} = \frac{1}{n} \sum_{i=1}^k n_i \tilde{\sigma}_i^2$ , где $\tilde{\sigma}_i^2 = \frac{1}{n_i} \sum_{j=1}^{n_i} (y_{ij} - \bar{y}_i)^2$ — групповая дисперсия	$n - k$	$\frac{n}{n - k} \tilde{D}_{\text{внутр}}$
Все признаки	$\tilde{D}_{\text{общ}} = \frac{1}{n} \sum_{i=1}^k \sum_{j=1}^{n_i} (y_{ij} - \bar{y})^2$	$n - 1$	$\frac{n}{n - 1} \tilde{D}_{\text{общ}}$

в) Рассчитать показатели вариации, используемые в дисперсионном анализе

Источник вариации	Показатель вариации	Число степеней свободы	Несмещенная оценка
Факторный признак	$\tilde{D}_{\text{межгр}} = 0.27$	$k - 1 = 2$	$\frac{n}{k - 1} \tilde{D}_{\text{межгр}} = 42.52$
Остаточные признаки	$\tilde{D}_{\text{внутр}} = 35.78$	$n - k = 312$	$\frac{n}{n - k} \tilde{D}_{\text{внутр}} = 36.12$
Все признаки	$\tilde{D}_{\text{общ}} = 36.05$	$n - 1 = 314$	$\frac{n}{n - 1} \tilde{D}_{\text{общ}} = 36.16$

г) Проверить правило сложения дисперсий

Показатель	$\tilde{D}_{\text{межгр}}$	$\tilde{D}_{\text{внутр}}$	$\tilde{D}_{\text{межгр}}$	$\tilde{D}_{\text{межгр}} + \tilde{D}_{\text{внутр}}$
Значение	0.27	35.78	36.05	36.05

д) Рассчитать показатели тесноты связи между факторным и результативным признаками

Показатель	Формула расчета	Значение
Эмпирический коэффициент детерминации	$\tilde{\eta}^2 = \frac{\tilde{D}_{\text{межгр}}}{\tilde{D}_{\text{общ}}}$	0.01
Эмпирическое корреляционное отношение	$\tilde{\eta} = \sqrt{\tilde{\eta}^2} = \sqrt{\frac{\tilde{D}_{\text{межгр}}}{\tilde{D}_{\text{общ}}}}$	0.09

е) Охарактеризовать тип связи между факторным и результативным признаками

По шкале Чеддока наблюдается отсутствие статистической связи между факторным признаком С5 и результативным признаком С4.

ж) Указать формулы расчёта показателей, используемых при проверке статистической гипотезы дисперсионного анализа

	Выражение	Пояснение использованных обозначений
Формула расчета статистики критерия	$Z = \frac{n - k}{k - 1} \frac{\tilde{D}_{\text{межгр}}}{\tilde{D}_{\text{внутр}}}$	$k$ — число групп $n$ — объём выборки $y_1, \dots, y_n$
Закон распределения статистики критерия при условии истинности основной гипотезы	$F(k - 1, n - k)$	
Формула расчета критической точки	$F_{1-\alpha, k-1, n-k}$	Правосторонняя критическая область
Формула расчета $p$ -value	$1 - F_Z(z_{\text{выб}}   H_0)$	

з) Проверить статистическую гипотезу дисперсионного анализа

Уровень значимости	Выборочное значение статистики критерия	$p$ -value	Статистическое решение	Вывод
0.01	1.18	0.31	$H_0$ принимается	$F_Y(y   x=x_1) = F_Y(y   x=x_2) = \dots = F_Y(y   x=x_k) = F_Y(y)$
0.05			$H_0$ принимается	$F_Y(y   x=x_1) = F_Y(y   x=x_2) = \dots = F_Y(y   x=x_k) = F_Y(y)$
0.1			$H_0$ принимается	$F_Y(y   x=x_1) = F_Y(y   x=x_2) = \dots = F_Y(y   x=x_k) = F_Y(y)$

Вывод (в терминах предметной области)

В результате проведённого в п.7 статистического анализа обнаружено, что между факторным признаком С5 и результативным признаком С4 отсутствует статистическая связь. Под действием С5 не оказывается влияние на распределение С4.



## 8. Корреляционный анализ

### 8.1. Расчёт парных коэффициентов корреляции

Анализируемый признак 1 – C11 (Dietary beta-carotene consumed (mcg per day))

Анализируемый признак 2 – C12 (Dietary retinol consumed (mcg per day))

Объёмы выборок –  $n_1 = n_2 = n = 315$

#### а) Рассчитать точечные оценки коэффициентов корреляции

	Формула расчёта	Значение
Линейный коэффициент корреляции	$\tilde{\rho}_{XY} = \frac{\text{cov}(X, Y)}{\tilde{\sigma}_X \cdot \tilde{\sigma}_Y} = \frac{\sum_{i=1}^n (x_i - \bar{x})(y_i - \bar{y})}{\sqrt{\sum_{i=1}^n (x_i - \bar{x})^2 \cdot \sum_{i=1}^n (y_i - \bar{y})^2}}$	0.05
Ранговый коэффициент корреляции по Спирмену	<p><math>\tilde{\rho}_{XY}^{(sp)} = \tilde{\rho}_{RS}</math>, где <math>R</math> и <math>S</math> — ранги для выборок <math>X</math> и <math>Y</math> соответственно.</p> <p>Можно показать, что <math>\tilde{\rho}_{XY}^{(sp)} = 1 - \frac{6S}{n(n^2 - 1)}</math>, где</p> $S = \sum_{i=1}^n (r_i - s_i)^2$	0.20
Ранговый коэффициент корреляции по Кендаллу	<p><math>\tilde{\tau}_{XY} = \frac{N_+ - N_-}{n(n-1)/2}</math>, где</p> <p><math>N_+</math>, <math>N_-</math> — количество пар точек <math>(x_i, y_i)</math> таких, что <math>(x_i - x_j)(y_i - y_j) &gt; 0</math>, <math>(x_i - x_j)(y_i - y_j) &lt; 0</math> соответственно.</p> <p>Другой вариант расчёта: <math>\tilde{\tau}_{XY} = \frac{4R}{n(n-1)} - 1</math>, где</p> $R = \sum_{i=1}^n \sum_{j=i+1}^n 1[s_j > s_i]$ <p>— число инверсий в выборке <math>(r_1, s_1), \dots, (r_n, s_n)</math>, предварительно отсортированной по возрастанию <math>r_i</math>.</p>	0.13

б) Привести формулы расчёта доверительного интервала для линейного коэффициента корреляции

Для небольших объёмов выборок:  $n < 500$

Граница доверительного интервала	Формула расчета
Нижняя граница	$th\left(\frac{1}{2} \ln \frac{1 + \tilde{\rho}}{1 - \tilde{\rho}} + \frac{\tilde{\rho}}{2(n-1)} - \frac{u_{1-\frac{\alpha}{2}}}{\sqrt{n-3}}\right)$
Верхняя граница	$th\left(\frac{1}{2} \ln \frac{1 + \tilde{\rho}}{1 - \tilde{\rho}} + \frac{\tilde{\rho}}{2(n-1)} + \frac{u_{1-\frac{\alpha}{2}}}{\sqrt{n-3}}\right)$

в) Рассчитать доверительные интервалы для линейного коэффициента корреляции

Граница доверительного интервала	$\alpha = 0.01$	$\alpha = 0.05$	$\alpha = 0.1$
Нижняя граница	-0.09	-0.06	-0.04
Верхняя граница	0.2	0.16	0.15

г) Указать формулы расчёта показателей, используемых при проверке значимости коэффициентов корреляции

Статистическая гипотеза	Формула расчета статистики критерия	Закон распределения статистики критерия при условии истинности основной
$H_0 : \rho_{XY} = 0$ $H' : \rho_{XY} \neq 0$	$Z = \frac{\tilde{\rho}_{XY} \cdot \sqrt{n-2}}{\sqrt{1 - \tilde{\rho}_{XY}^2}}$	$T(n-2)$
$H_0 : \rho_{XY}^{(sp)} = 0$ $H' : \rho_{XY}^{(sp)} \neq 0$	$Z = \frac{\tilde{\rho}_{XY}^{(sp)} \cdot \sqrt{n-2}}{\sqrt{1 - \left(\tilde{\rho}_{XY}^{(sp)}\right)^2}}$	$T(n-2)$
$H_0 : \tau_{XY} = 0$ $H' : \tau_{XY} \neq 0$	$\sqrt{\frac{9n(n+1)}{2(2n+5)}} \cdot \tilde{\tau}_{XY}$	$N(0,1)$

д) Проверить значимость коэффициентов корреляции

Статистическая гипотеза	Уровень значимости	Выборочное значение статистики критерия	$p$ -value	Статистическое решение	Вывод
$H_0 : \rho_{XY} = 0$ $H' : \rho_{XY} \neq 0$	0.1	0.94	0.35	$H_0$ принимается	$\rho_{XY} = 0$
$H_0 : \rho_{XY}^{(sp)} = 0$ $H' : \rho_{XY}^{(sp)} \neq 0$	0.1	3.54	0.00	$H_0$ отклоняется	$\rho_{XY}^{(sp)} \neq 0$
$H_0 : \tau_{XY} = 0$ $H' : \tau_{XY} \neq 0$	0.1	3.53	0.00	$H_0$ отклоняется	$\tau_{XY} \neq 0$

## 8.2. Расчёт множественных коэффициентов корреляции

Анализируемый признак 1 – C9 (Number of alcoholic drinks consumed per week)

Анализируемый признак 2 – C10 (Cholesterol consumed (mg per day))

Анализируемый признак 3 – C11 (Dietary beta-carotene consumed (mcg per day))

Объёмы выборок –  $n_1 = n_2 = n_3 = n = 315$

а) Рассчитать матрицу ранговых коэффициентов корреляции по Кендаллу

Признак \ Признак	C9	C10	C11
C9	1.00	0.07	0.04
C10	0.07	1.00	0.13
C11	0.04	0.13	1.00

б) Рассчитать матрицу значений  $p$ -value для ранговых коэффициентов корреляции по Кендаллу

Статистическая гипотеза:  $H_0 : \tau = 0$   
 $H' : \tau \neq 0$

Признак \ Признак	C9	C10	C11
C9	—	0.10	0.27
C10	0.10	—	0.00
C11	0.27	0.00	—

в) Рассчитать точечную оценку коэффициента конкордации

	Формула расчета	Значение
Коэффициент конкордации	$\tilde{W} = \frac{12}{k^2(n^3 - n)} \cdot \sum_{i=1}^n \left( \sum_{j=1}^k r_{ij} - \frac{k(n+1)}{2} \right)^2$ , где $r_{ij}$ — ранг $i$ -ого объекта в $j$ -ой выборке. $k$ — количество выборок. $n$ — объём выборок.	0.41

г) Указать формулы расчёта показателей, используемых при проверке значимости коэффициента конкордации

Статистический критерий:  $H_0 : W = 0$   
 $H' : W \neq 0$

	Выражение	Пояснение использованных обозначений
Формула расчета статистики критерия	$Z = n(k - 1)\tilde{W}$	$\tilde{W}$ — точечная оценка коэффициента конкордации. $k$ — количество выборок. $n$ — объём каждой выборки.
Закон распределения статистики критерия при условии истинности основной гипотезы	$\chi^2(n - 1)$	
Формула расчета критической точки	$\chi^2_{1-\alpha, n-1}$	Правосторонняя критическая область

Формула расчета <i>p-value</i>	$1 - F_Z(z_{\text{выб}}   H_0)$
--------------------------------	---------------------------------

д) Проверить значимость коэффициента конкордации

Уровень значимости	Выборочное значение статистики критерия	<i>p-value</i>	Статистическое решение	Вывод
0.01	255.27	0.99	$H_0$ принимается	$W = 0$
0.05			$H_0$ принимается	$W = 0$
0.1			$H_0$ принимается	$W = 0$

Вывод (в терминах предметной области)

В результате проведенного в п.8 статистического анализа обнаружено, что между признаками С11 и С12 нет линейной корреляционной связи, однако существует слабая, но статистически значимая монотонная корреляционная связь. Между признаками С9, С10, С11 наблюдается отсутствие значимой монотонной связи между их парами. Наибольшая монотонная связь между С10 и С11, притом она статистически значимая. Между остальными парами можно утверждать о том, что монотонная связь между ними отсутствует. Также общая согласованность рангов наблюдается слабая, статистически незначимая.

## 9. Регрессионный анализ

### 9.1 Простейшая линейная регрессионная модель

Факторный признак  $x$  – C11 (Dietary beta-carotene consumed (mcg per day))

Результативный признак  $y$  – C13 (Plasma beta-carotene (ng/ml))

Уравнение регрессии –  $f(x) = \beta_0 + \beta_1 x$

#### 9.1.1. Точечные оценки линейной регрессионной модели

а) Рассчитать точечные оценки параметров линейной регрессионной модели

Параметр	Формула расчета	Значение
$\beta_0$	$\bar{y} - \tilde{\rho}_{XY} \cdot \frac{\tilde{\sigma}_Y}{\tilde{\sigma}_X} \cdot \bar{x}$	128.89
$\beta_1$	$\tilde{\rho}_{XY} \cdot \frac{\tilde{\sigma}_Y}{\tilde{\sigma}_X}$	0.03

б) Записать точечную оценку уравнения регрессии

$$\tilde{f}(x) = \tilde{\beta}_0 + \tilde{\beta}_1 x = 128.89 + 0.03x$$

в) Привести формулы расчёта показателей вариации, используемых в регрессионном анализе

Источник вариации	Показатель вариации	Число степеней свободы	Несмещенная оценка
Регрессия	$\tilde{D}_{Y X} = \frac{1}{n} \sum_{i=1}^n (\tilde{f}(x_i) - \bar{y})^2$	$k - 1$	$\frac{n}{k - 1} \tilde{D}_{Y X}$
Остаточные признаки	$\tilde{D}_{Y \text{ ост}} = \frac{1}{n} \sum_{i=1}^n (\tilde{f}(x_i) - y_i)^2$	$n - k$	$\frac{n}{n - k} \tilde{D}_{Y \text{ ост}}$
Все признаки	$\tilde{D}_{Y \text{ общ}} = \frac{1}{n} \sum_{i=1}^n (y_i - \bar{y})^2$	$n - 1$	$\frac{n}{n - 1} \tilde{D}_{Y \text{ общ}}$

$k$  — число оцениваемых параметров функции регрессии  $f(x)$ .

г) Рассчитать показатели вариации, используемые в регрессионном анализе

Источник вариации	Показатель вариации	Число степеней свободы	Несмещенная оценка
Регрессия	$\tilde{D}_{Y X} = 1686.70$	$k - 1 = 1$	$\frac{n}{k - 1} \tilde{D}_{Y X} = 531311.35$
Остаточные признаки	$\tilde{D}_{Y \text{ ост}} = 31696.28$	$n - k = 313$	$\frac{n}{n - k} \tilde{D}_{Y \text{ ост}} = 31898.81$
Все признаки	$\tilde{D}_{Y \text{ общ}} = 33382.98$	$n - 1 = 314$	$\frac{n}{n - 1} \tilde{D}_{Y \text{ общ}} = 33489.29$

д) Проверить правило сложения дисперсий

Показатель	$\tilde{D}_{Y X}$	$\tilde{D}_{Y \text{ ост}}$	$\tilde{D}_{Y \text{ общ}}$	$\tilde{D}_{Y X} + \tilde{D}_{Y \text{ ост}}$
Значение	1686.70	31696.28	33382.98	33382.98

е) Рассчитать показатели тесноты связи между факторным и результативным признаками

Показатель	Формула расчета	Значение
Коэффициент детерминации	$\tilde{R}^2 = \frac{\tilde{D}_{Y X}}{\tilde{D}_{Y \text{ общ}}}$	0.05
Корреляционное отношение	$\tilde{R} = \sqrt{\tilde{R}^2} = \sqrt{\frac{\tilde{D}_{Y X}}{\tilde{D}_{Y \text{ общ}}}}$	0.22

ж) Охарактеризовать тип связи между факторным и результативным признаками, определяемой рассчитанной линейной регрессией

Наблюдается очень слабая корреляционная связь между факторным признаком С11 и результативным признаком С13.

#### 9.1.2. Интервальные оценки линейной регрессионной модели

а) Привести формулы расчёта доверительных интервалов для параметров линейной регрессионной модели

Параметр	Границы доверительного интервала	Формула расчета
$\beta_0$	Нижняя граница	$\tilde{\beta}_0 - t_{1-\frac{\alpha}{2}, (n-2)} \cdot \sqrt{\tilde{D}_{Y \text{ ост}}^{\text{несмещ}}} \cdot \sqrt{\frac{\sum_{i=1}^n x_i^2}{n^2 \tilde{D}_X}}$
	Верхняя граница	$\tilde{\beta}_0 + t_{1-\frac{\alpha}{2}, (n-2)} \cdot \sqrt{\tilde{D}_{Y \text{ ост}}^{\text{несмещ}}} \cdot \sqrt{\frac{\sum_{i=1}^n x_i^2}{n^2 \tilde{D}_X}}$
$\beta_1$	Нижняя граница	$\tilde{\beta}_1 - t_{1-\frac{\alpha}{2}, (n-2)} \cdot \sqrt{\tilde{D}_{Y \text{ ост}}^{\text{несмещ}}} \cdot \sqrt{\frac{1}{n \tilde{D}_X}}$
	Верхняя граница	$\tilde{\beta}_1 + t_{1-\frac{\alpha}{2}, (n-2)} \cdot \sqrt{\tilde{D}_{Y \text{ ост}}^{\text{несмещ}}} \cdot \sqrt{\frac{1}{n \tilde{D}_X}}$

б) Рассчитать доверительные интервалы для параметров линейной регрессионной модели

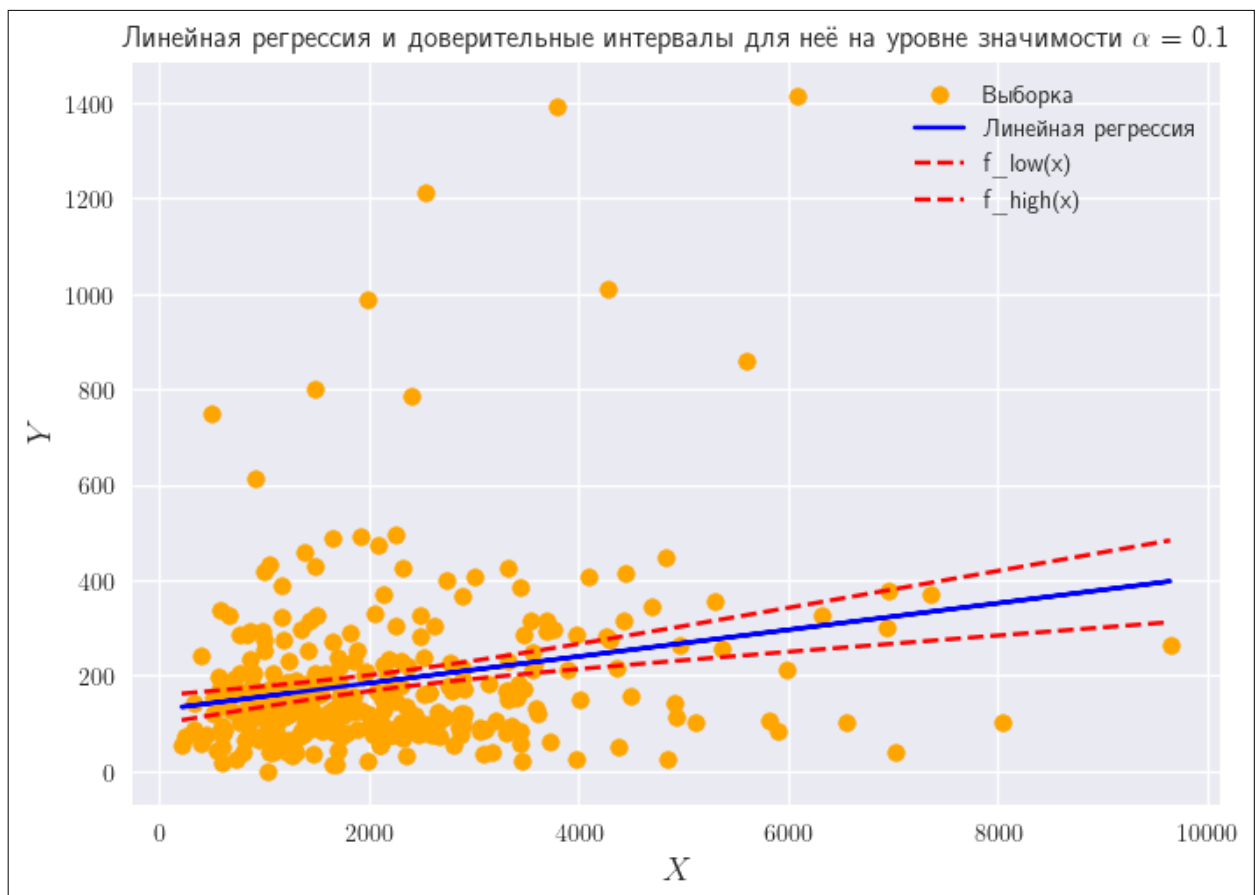
Параметр	Границы доверительного интервала	$\alpha = 0.01$	$\alpha = 0.05$	$\alpha = 0.1$
$\beta_0$	Нижняя граница	82.20	93.44	99.17
	Верхняя граница	175.59	164.35	158.62
$\beta_1$	Нижняя граница	0.01	0.01	0.02
	Верхняя граница	0.05	0.04	0.04



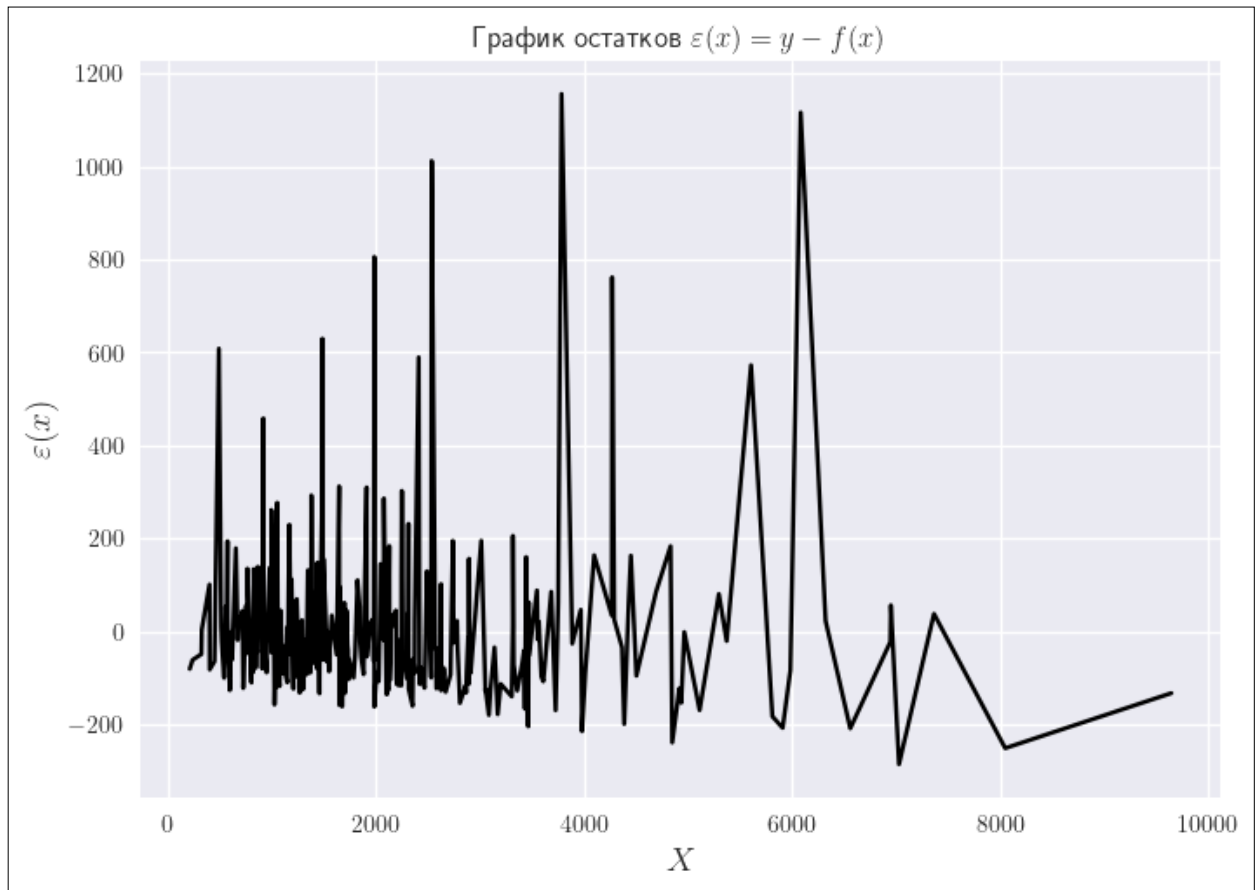
в) Привести формулы расчёта доверительного интервала для значений регрессии  $f(x)$

Границы доверительного интервала	Формула расчета
Нижняя граница $f_{low}(x)$	$(\tilde{\beta}_0 + \tilde{\beta}_1 x) - t_{1-\frac{\alpha}{2}, (n-2)} \cdot \sqrt{\tilde{D}_{Y \text{ ост}}} \cdot \sqrt{\frac{1}{n} + \frac{(x - \bar{x})^2}{n \tilde{D}_X}}$
Верхняя граница $f_{high}(x)$	$(\tilde{\beta}_0 + \tilde{\beta}_1 x) + t_{1-\frac{\alpha}{2}, (n-2)} \cdot \sqrt{\tilde{D}_{Y \text{ ост}}} \cdot \sqrt{\frac{1}{n} + \frac{(x - \bar{x})^2}{n \tilde{D}_X}}$

г) Построить диаграмму рассеяния признаков  $x$  и  $y$ . Нанести на диаграмму функцию регрессии  $f(x)$ , а также нижние и верхние границы линии регрессии  $f_{low}(x)$  и  $f_{high}(x)$  на уровне значимости  $\alpha = 0.1$



д) Построить график остатков  $\varepsilon(x) = y - f(x)$



### 9.1.3. Проверка значимости линейной регрессионной модели

Статистическая гипотеза –  $H_0 : \beta_1 = 0$   
 $H' : \beta_1 \neq 0$

а) Указать формулы расчёта показателей, используемых при проверке значимости линейной регрессионной модели

	Выражение	Пояснение использованных обозначений
Формула расчета статистики критерия	$Z = \frac{\tilde{D}_{Y X}}{\tilde{D}_{Y \text{ ост}}/(n-2)}$	$n$ — объём выборки $(x_1, y_1), \dots, (x_n, y_n)$
Закон распределения статистики критерия при условии истинности основной гипотезы	$F(1, n-2)$	
Формула расчета критической точки	$F_{1-\alpha, 1, n-2}$	Правосторонняя критическая область
Формула расчета $p$ -value	$1 - F_Z(z_{\text{выб}}   H_0)$	

б) Проверить значимость линейной регрессионной модели

Уровень значимости	Выборочное значение статистики критерия	$p$ -value	Статистическое решение	Вывод
0.01	16.66	0.00	$H_0$ отклоняется	$\beta_1 \neq 0$
0.05			$H_0$ отклоняется	$\beta_1 \neq 0$
0.1			$H_0$ отклоняется	$\beta_1 \neq 0$

## 9.2 Линейная регрессионная модель общего вида

Факторный признак  $x$  – C11 (Dietary beta-carotene consumed (mcg per day))

Результативный признак  $y$  – C13 (Plasma beta-carotene (ng/ml))

Уравнение регрессии – квадратичное по  $x$ :  $f(x) = \beta_0 + \beta_1 x + \beta_2 x^2$

### 9.2.1. Точечные оценки линейной регрессионной модели

а) Рассчитать точечные оценки параметров линейной регрессионной модели

Параметр	Формула расчета	Значение
$\tilde{\beta}_{\downarrow} = \begin{pmatrix} \tilde{\beta}_0 \\ \tilde{\beta}_1 \\ \tilde{\beta}_2 \end{pmatrix}$	$\tilde{\beta}_{\downarrow} = (F^T F)^{-1} F^T y_{\downarrow}$ , где $F$ — регрессионная матрица, $y_{\downarrow}$ — вектор значений результативного признака.	$\begin{pmatrix} 117.46 \\ 0.04 \\ -0.00 \end{pmatrix}$

б) Записать точечную оценку уравнения регрессии

$$\tilde{f}(x) = \tilde{\beta}_0 + \tilde{\beta}_1 x + \tilde{\beta}_2 x^2 = 117.46 + 0.04x - 0.00x^2$$

в) Рассчитать показатели вариации, используемые в регрессионном анализе

Источник вариации	Показатель вариации	Число степеней свободы	Несмещенная оценка
Регрессия	$\tilde{D}_{Y X} = 1715.27$	$k - 1 = 2$	$\frac{n}{k - 1} \tilde{D}_{Y X} = 270154.30$
Остаточные признаки	$\tilde{D}_{Y \text{ ост}} = 31667.71$	$n - k = 312$	$\frac{n}{n - k} \tilde{D}_{Y \text{ ост}} = 31972.21$
Все признаки	$\tilde{D}_{Y \text{ общ}} = 33382.98$	$n - 1 = 314$	$\frac{n}{n - 1} \tilde{D}_{Y \text{ общ}} = 33489.29$

$k$  — число оцениваемых параметров функции регрессии  $f(x)$ .

г) Проверить правило сложения дисперсий

Показатель	$\tilde{D}_{Y X}$	$\tilde{D}_{Y \text{ ост}}$	$\tilde{D}_{Y \text{ общ}}$	$\tilde{D}_{Y X} + \tilde{D}_{Y \text{ ост}}$
Значение	1715.27	31667.71	33382.98	33382.98

д) Рассчитать показатели тесноты связи между факторным и результативным признаками

Показатель	Формула расчета	Значение
Коэффициент детерминации	$\tilde{R}^2 = \frac{\tilde{D}_{Y X}}{\tilde{D}_{Y \text{ общ}}}$	0.05
Корреляционное отношение	$\tilde{R} = \sqrt{\tilde{R}^2} = \sqrt{\frac{\tilde{D}_{Y X}}{\tilde{D}_{Y \text{ общ}}}}$	0.23

е) Охарактеризовать тип связи между факторным и результативным признаками, определяемой рассчитанной линейной регрессией

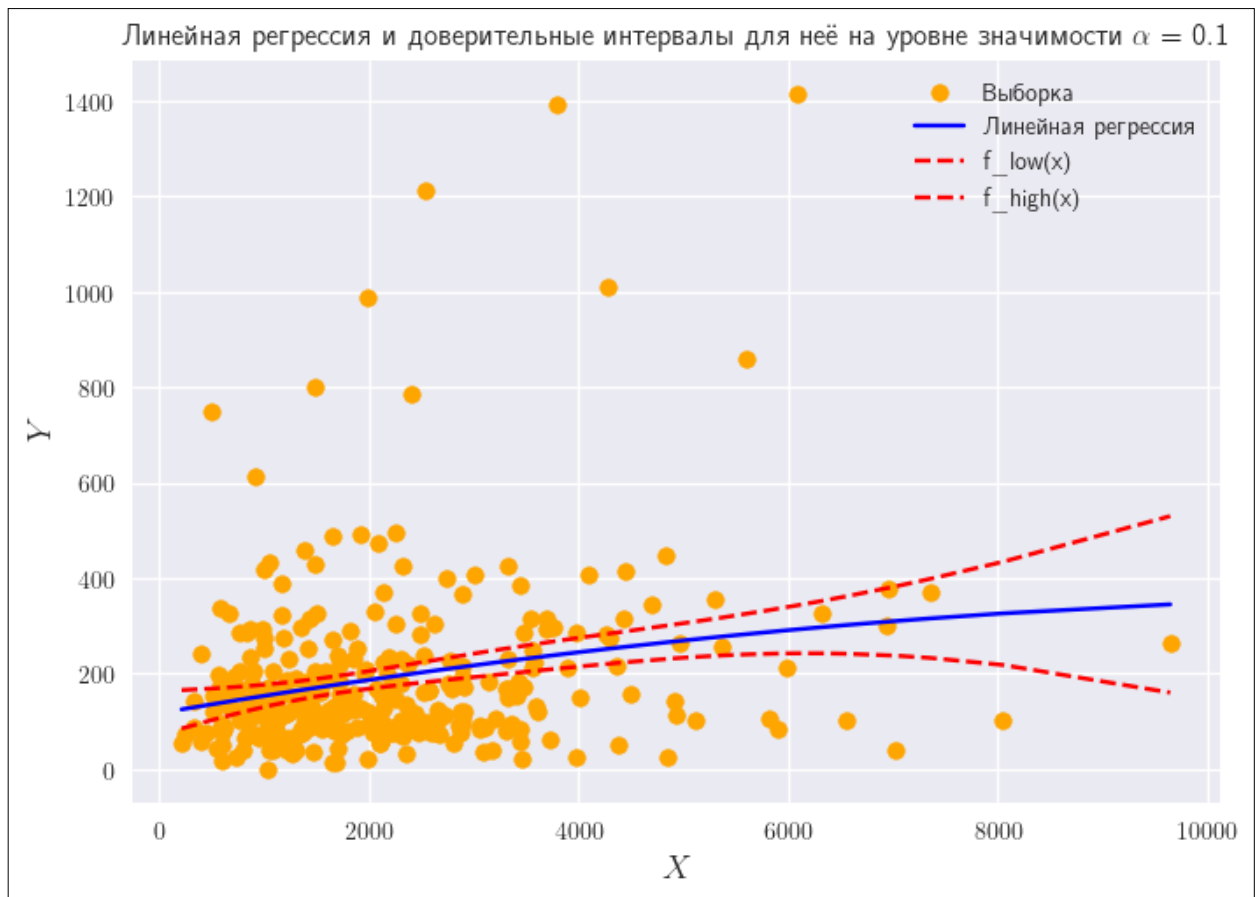
Наблюдается очень слабая корреляционная связь между факторным признаком С11 и результативным признаком С13.

### 9.2.2. Интервальные оценки линейной регрессионной модели

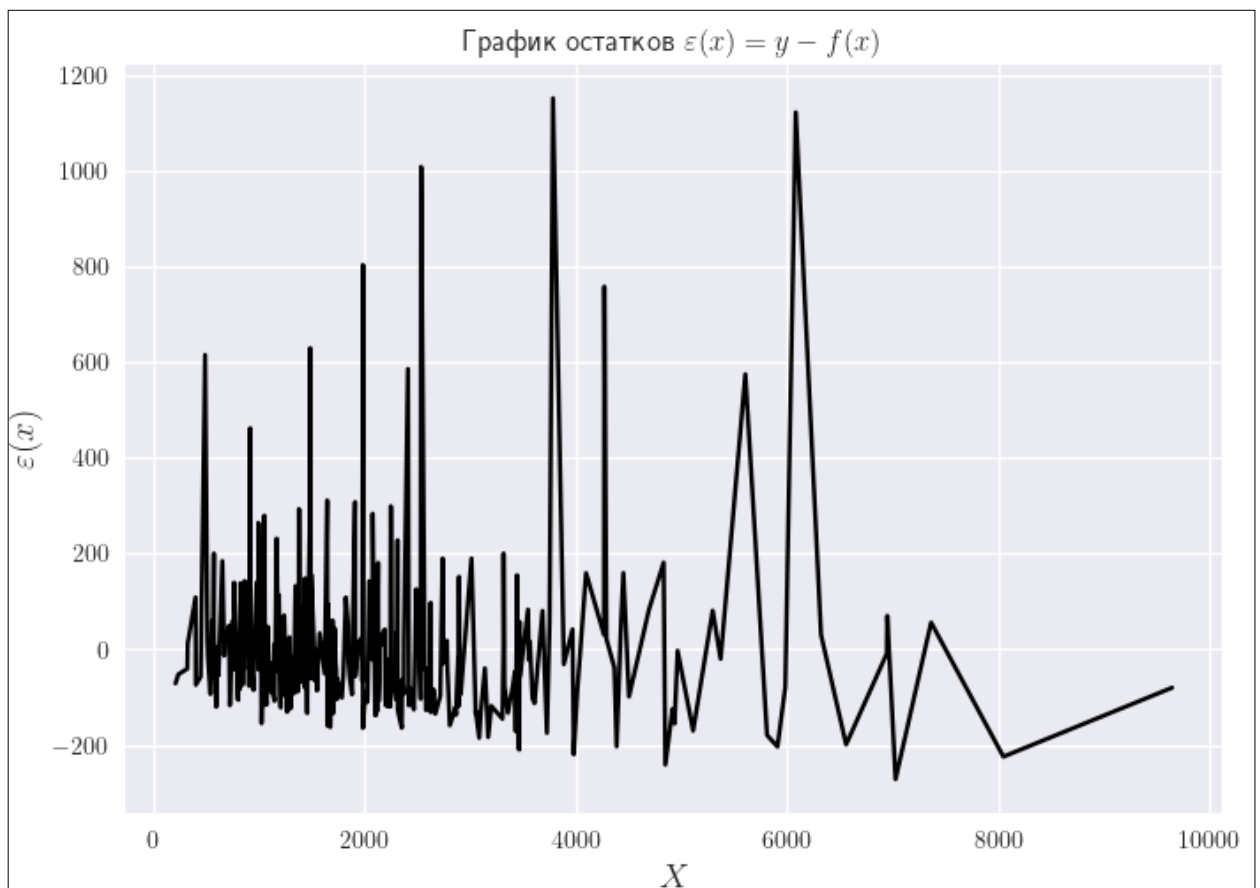
а) Привести формулы расчёта доверительного интервала для значений регрессии  $f(x)$

Границы доверительного интервала	Формула расчета
Нижняя граница $f_{low}(x)$	$\tilde{f}(x) - t_{1-\frac{\alpha}{2}, (n-k)} \cdot \sqrt{\tilde{D}_{Y \text{ ост}}^{\text{несмещ}}} \cdot \sqrt{\vec{\varphi}(x)(F^T F)^{-1} \varphi_{\downarrow}(x)}$
Верхняя граница $f_{high}(x)$	$\tilde{f}(x) + t_{1-\frac{\alpha}{2}, (n-k)} \cdot \sqrt{\tilde{D}_{Y \text{ ост}}^{\text{несмещ}}} \cdot \sqrt{\vec{\varphi}(x)(F^T F)^{-1} \varphi_{\downarrow}(x)}$

б) Построить диаграмму рассеяния признаков  $x$  и  $y$ . Нанести на диаграмму функцию регрессии  $f(x)$ , а также нижние и верхние границы линии регрессии  $f_{low}(x)$  и  $f_{high}(x)$  на уровне значимости  $\alpha = 0.1$



в) Построить график остатков  $\varepsilon(x) = y - f(x)$



### 9.2.3. Проверка значимости линейной регрессионной модели

Статистическая гипотеза –  $H_0 : \beta_1 = \beta_2 = 0$   
 $H' : \beta_1^2 + \beta_2^2 > 0$

а) Указать формулы расчёта показателей, используемых при проверке значимости линейной регрессионной модели

	Выражение	Пояснение использованных обозначений
Формула расчета статистики критерия	$Z = \frac{\tilde{D}_{Y X}^{\text{несмещ}}}{\tilde{D}_{Y \text{ ост}}^{\text{несмещ}}}$	$n$ — объём выборки $(x_1, y_1), \dots, (x_n, y_n)$
Закон распределения статистики критерия при условии истинности основной гипотезы	$F(k - 1, n - 2)$	
Формула расчета критической точки	$F_{1-\alpha, k-1, n-2}$	Правосторонняя критическая область
Формула расчета $p$ -value	$1 - F_Z(z_{\text{Выб}}   H_0)$	

б) Проверить значимость линейной регрессионной модели

Уровень значимости	Выборочное значение статистики критерия	$p$ -value	Статистическое решение	Вывод
0.01	8.45	0.00	$H_0$ отклоняется	$\beta_1^2 + \beta_2^2 > 0$
0.05			$H_0$ отклоняется	$\beta_1^2 + \beta_2^2 > 0$
0.1			$H_0$ отклоняется	$\beta_1^2 + \beta_2^2 > 0$

### 9.3 Множественная линейная регрессионная модель

Факторный признак 1  $x_1$  – C11 (Dietary beta-carotene consumed (mcg per day))

Факторный признак 2  $x_2$  – C4 (Quetelet (weight/(height^2)))

Результативный признак  $y$  – C13 (Plasma beta-carotene (ng/ml))

Уравнение регрессии –  $f(x) = \beta_0 + \beta_1 x_1 + \beta_2 x_2$

а) Рассчитать точечные оценки параметров линейной регрессионной модели

Параметр	Формула расчета	Значение
$\tilde{\beta}_{\downarrow} = \begin{pmatrix} \tilde{\beta}_0 \\ \tilde{\beta}_1 \\ \tilde{\beta}_2 \end{pmatrix}$	$\tilde{\beta}_{\downarrow} = (F^T F)^{-1} F^T y_{\downarrow}$ , где $F$ — регрессионная матрица, $y_{\downarrow}$ — вектор значений результативного признака.	$\begin{pmatrix} 310.72 \\ 0.03 \\ -6.94 \end{pmatrix}$

б) Записать точечную оценку уравнения регрессии

$\tilde{f}(x) = \tilde{\beta}_0 + \tilde{\beta}_1 x_1 + \tilde{\beta}_2 x_2 = 310.72 + 0.03x_1 - 6.94x_2$
---

в) Рассчитать показатели вариации, используемые в регрессионном анализе

Источник вариации	Показатель вариации	Число степеней свободы	Несмещенная оценка
Регрессия	$\tilde{D}_{Y X_1, X_2} = 3420.69$	$k - 1 = 2$	$\frac{n}{k - 1} \tilde{D}_{Y X_1, X_2} = 538757.89$
Остаточные признаки	$\tilde{D}_{Y \text{ ост}} = 29962.29$	$n - k = 312$	$\frac{n}{n - k} \tilde{D}_{Y \text{ ост}} = 30250.39$
Все признаки	$\tilde{D}_{Y \text{ общ}} = 33382.98$	$n - 1 = 314$	$\frac{n}{n - 1} \tilde{D}_{Y \text{ общ}} = 33489.29$

$k$  — число оцениваемых параметров функции регрессии  $f(x)$ .

г) Проверить правило сложения дисперсий

Показатель	$\tilde{D}_{Y X_1, X_2}$	$\tilde{D}_{Y \text{ ост}}$	$\tilde{D}_{Y \text{ общ}}$	$\tilde{D}_{Y X_1, X_2} + \tilde{D}_{Y \text{ ост}}$
Значение	3420.69	29962.29	33382.98	33382.98

д) Рассчитать показатели тесноты связи между факторным и результативным признаками

Показатель	Формула расчета	Значение
Коэффициент детерминации	$\tilde{R}^2 = \frac{\tilde{D}_{Y X_1, X_2}}{\tilde{D}_{Y \text{ общ}}}$	0.10
Корреляционное отношение	$\tilde{R} = \sqrt{\tilde{R}^2} = \sqrt{\frac{\tilde{D}_{Y X_1, X_2}}{\tilde{D}_{Y \text{ общ}}}}$	0.32

е) Охарактеризовать тип связи между факторным и результативным признаками, определяемой рассчитанной линейной регрессией

Наблюдается слабая корреляционная связь между факторным признаками С11 и С4 и результативным признаком С13.

#### 9.4. Выводы

а) Сводная таблица показателей вариации для различных регрессионных моделей

Источник вариации	Простейшая линейная модель	Линейная модель с квадратичным членом	Множественная линейная модель
Регрессия	$\tilde{D}_{Y X} = 1686.70$	$\tilde{D}_{Y X} = 1715.27$	$\tilde{D}_{Y X_1, X_2} = 3420.69$
Остаточные признаки	$\tilde{D}_{Y \text{ ост}} = 31696.28$	$\tilde{D}_{Y \text{ ост}} = 31667.71$	$\tilde{D}_{Y \text{ ост}} = 29962.29$
Все признаки	$\tilde{D}_{Y \text{ общ}} = 33382.98$	$\tilde{D}_{Y \text{ общ}} = 33382.98$	$\tilde{D}_{Y \text{ общ}} = 33382.98$

б) Сводная таблица свойств различных регрессионных моделей

Свойство	Простейшая линейная модель	Линейная модель с квадратичным членом	Множественная линейная модель
Точность	Неточная	Неточная	Неточная
Значимость	Значимая	Значимая	Значимая
Адекватность	Неадекватная	Неадекватная	Неадекватная
Степень тесноты связи	Очень слабая	Очень слабая	Слабая



*Вывод (в терминах предметной области)*

В результате проведённого в п.9 статистического анализа обнаружено, что результативный признак С13 имеет слабую статистическую связь как с факторным признаком С11 отдельно, так и с двумя факторными признаками С11 и С4 вместе.