

МИНИСТЕРСТВО НАУКИ И ВЫСШЕГО ОБРАЗОВАНИЯ РОССИЙСКОЙ ФЕДЕРАЦИИ  
Федеральное государственное автономное образовательное учреждение высшего образования  
Национальный исследовательский технологический университет «МИСИС»  
Институт информационных технологий и автоматизированных систем управления  
Кафедра Бизнес-информатики и систем управления производством

**Практическая работа №7**

по дисциплине «Статистические методы анализа данных в принятии решений»  
на тему «Построение и оценка регрессионной модели»

Направление подготовки  
38.03.05 Бизнес-информатика  
Семестр 4

Выполнил:

Сычиков Владимир Андреевич

\_\_\_\_\_  
(ФИО студента)

ББИ-23-6

\_\_\_\_\_  
(№ группы)

31.03.2025

\_\_\_\_\_  
(дата сдачи)

Подпись:

\_\_\_\_\_

Проверил:

\_\_\_\_\_  
(ФИО преподавателя)

\_\_\_\_\_  
(оценка)

\_\_\_\_\_  
(дата проверки)

Подпись:

\_\_\_\_\_

Москва – 2025

## Ход работы:

Для начала работы я импортировал все необходимые библиотеки(pandas, numpy) с помощью команды `import`, а также задал элиасы для удобства обращения к библиотеке. Далее я подгрузил выборку и записал ее в датафрейм `df`. Также я выделил нужные мне строки из выборки, по заданию.

```
[1] import pandas as pd
import matplotlib.pyplot as plt
import seaborn as sns
import statsmodels.api as sm
from statsmodels.stats.outliers_influence import variance_inflation_factor

df = pd.read_excel('data.xlsx')
```

Рисунок 1 - Импорт библиотек

```
[2] import numpy as np
import pandas as pd
import statsmodels.api as sm
import matplotlib.pyplot as plt
from sklearn.linear_model import LinearRegression
from sklearn.metrics import mean_squared_error, r2_score

df = pd.read_excel('data.xlsx')

columns = ["Стоимость сырья", "Затраты на электроэнергию", "Зарплата",
           "Административные расходы", "Себестоимость продукции"]
df[columns] = df[columns].apply(pd.to_numeric, errors='coerce')

df = df.dropna(subset=columns)

X = df[["Стоимость сырья", "Затраты на электроэнергию",
        "Зарплата", "Административные расходы"]]
Y = df["Себестоимость продукции"]

X = sm.add_constant(X)

model = sm.OLS(Y, X).fit()
print(model.summary())

alpha = 0.001
n_iter = 1000
m = len(Y)

X_np = X.values.astype(float)
Y_np = Y.values.astype(float)

theta = np.zeros(X_np.shape[1])

for i in range(n_iter):
    gradients = - (2/m) * X_np.T.dot(Y_np - X_np.dot(theta))
    theta -= alpha * gradients

print("Оцененные коэффициенты (градиентный спуск):", theta)
```

Рисунок 2 - Код для реализации модели



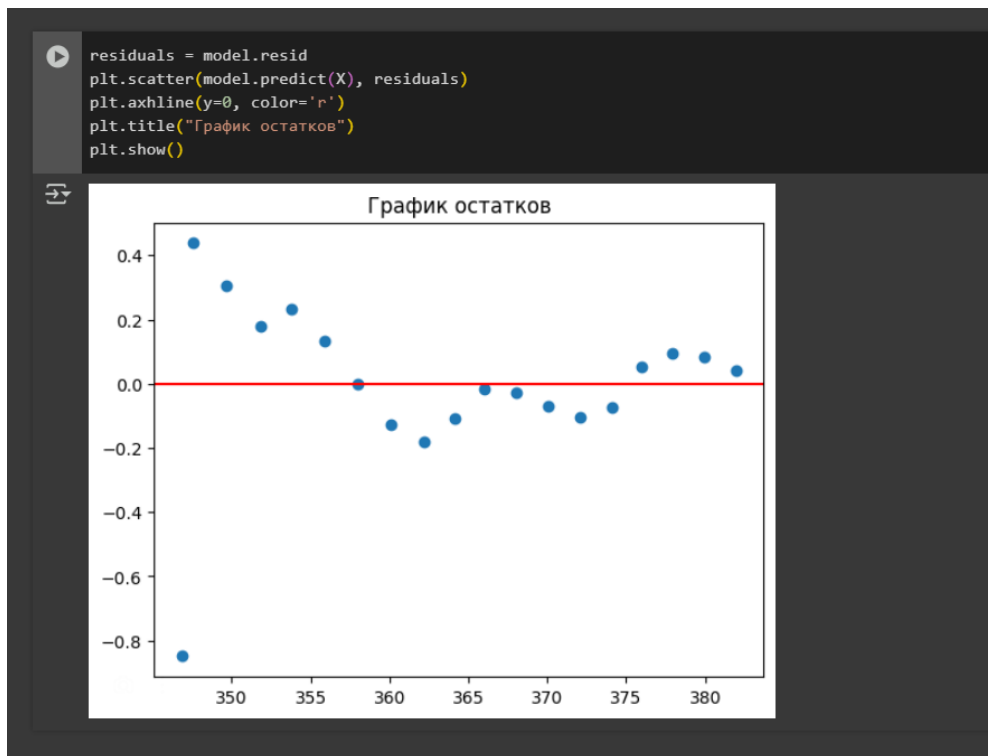


Рисунок 4 - График остатков

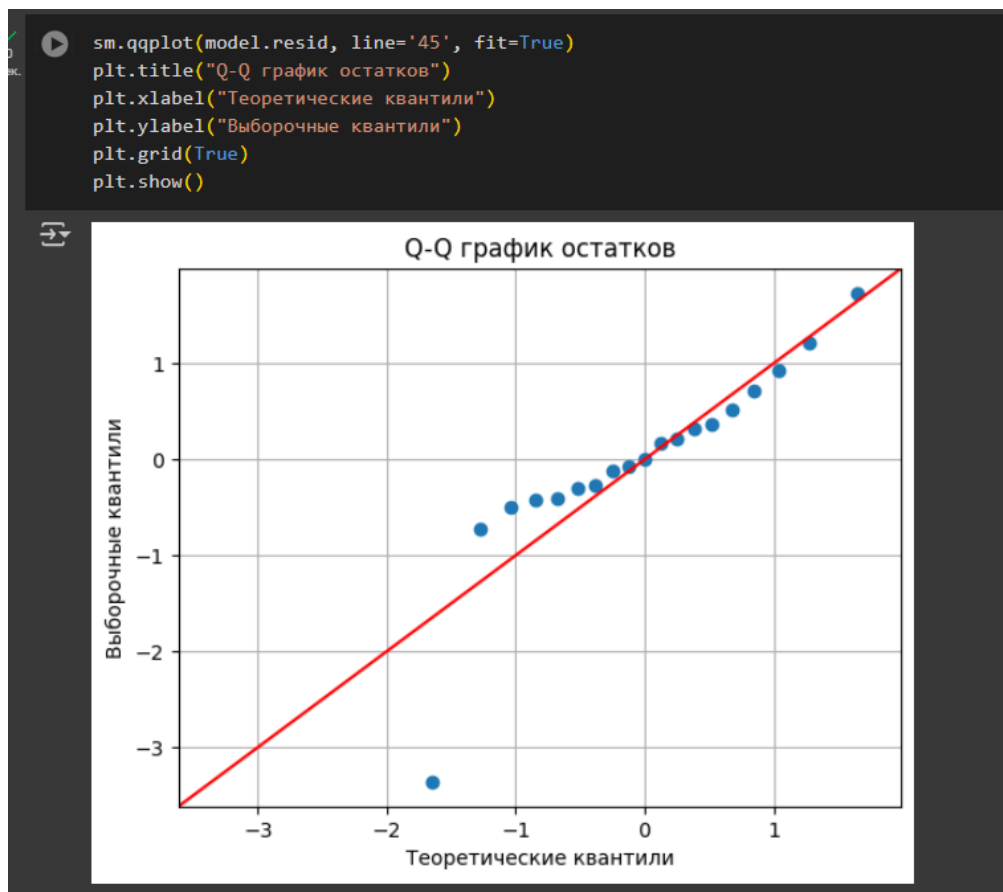


Рисунок 5 - Q-Q график остатков

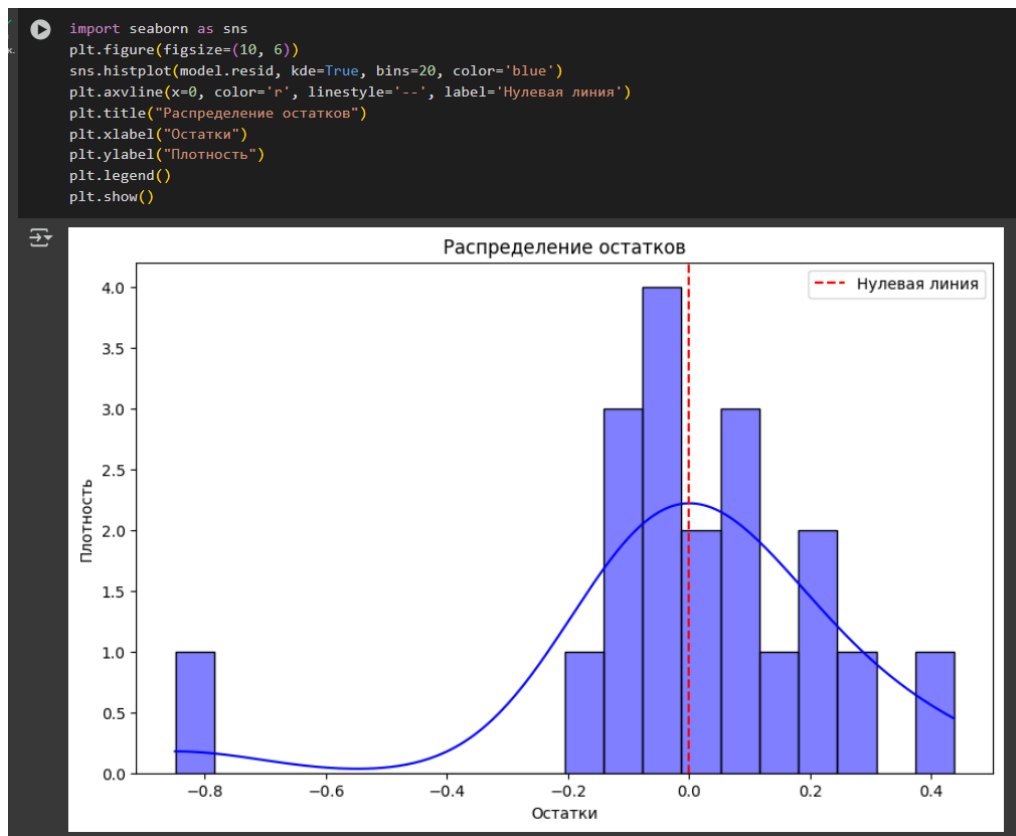


Рисунок 6 - Гистограмма и плотность распределения остатков

Аналитика остатков:

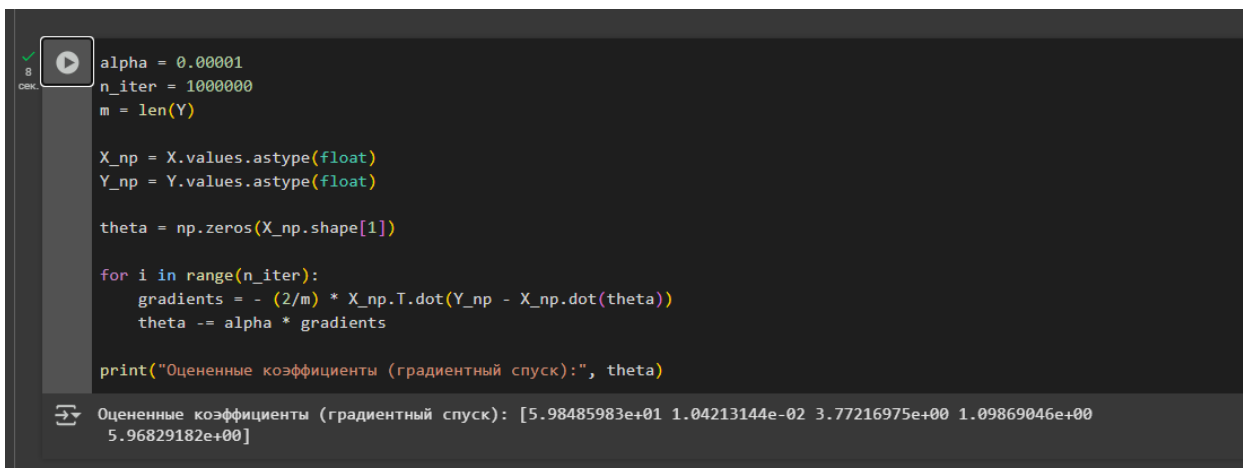
- Остатки расположены в диапазоне приблизительно от -5 до +5 (ось Y).
- Наблюдается симметричное распределение относительно нулевой линии (красная горизонталь), что соответствует предположению о нулевом среднем значении ошибок.
- Отсутствует явная систематическая структура (например, веерообразное или U-образное распределение). Это говорит о постоянной дисперсии остатков.
- Нет кластеризации или аномальных выбросов, что подтверждает случайный характер ошибок.

```
Notes:
[1] Standard Errors assume that the covariance matrix of the errors is correctly specified.
[2] The condition number is large, 1.86e+03. This might indicate that there are
strong multicollinearity or other numerical problems.
Оцененные коэффициенты (градиентный спуск): [nan nan nan nan nan]
/usr/local/lib/python3.11/dist-packages/scipy/stats/_axis_nan_policy.py:418: UserWarning: `kurtosis`
return hypotest_fun_in(*args, **kwargs)
<ipython-input-2-b0b758583d71>:36: RuntimeWarning: invalid value encountered in subtract
theta -= alpha * gradients
```

Рисунок 7 - Проблема с градиентным спуском

Можем заметить, что существует проблема с коэффициентами градиентного спуска. Это может быть связано с неудачным выбором шага обучения(alpha), а также отсутствием нормализации данных(разные масштабы признаков)

Для решения этой проблемы можем уменьшить  $\alpha$  и увеличить  $n\_iter$



```
alpha = 0.00001
n_iter = 1000000
m = len(Y)

X_np = X.values.astype(float)
Y_np = Y.values.astype(float)

theta = np.zeros(X_np.shape[1])

for i in range(n_iter):
    gradients = - (2/m) * X_np.T.dot(Y_np - X_np.dot(theta))
    theta -= alpha * gradients

print("Оцененные коэффициенты (градиентный спуск):", theta)
```

Оцененные коэффициенты (градиентный спуск): [5.98485983e+01 1.04213144e-02 3.77216975e+00 1.09869046e+00 5.96829182e+00]

Рисунок 8 - Решение проблемы

Шаг обучения ( $\alpha$ ) определяет размер "шага", который алгоритм делает при обновлении коэффициентов. Если  $\alpha$  слишком большой (например, 0.1), градиентный спуск может начать "перепрыгивать" оптимальные значения коэффициентов, что приведет к расходимости — MSE будет расти, а коэффициенты могут стать неадекватно большими или даже превратиться в nan. В моем случае изначально такая проблема и возникала. Напротив, слишком маленький  $\alpha$  (например, 0.00001) замедляет сходимость: алгоритму потребуется гораздо больше итераций, чтобы найти минимум. В моем исправленном коде был использован очень маленький  $\alpha$  (0.00001), что потребовало увеличения  $n\_iter$  до 1 000 000, но зато получилась сходимость. Оптимальный  $\alpha$  обычно подбирается экспериментально, но для данных с мультиколлинеарностью лучше выбирать меньшие значения, чтобы избежать резких колебаний градиента.

Количество итераций ( $n\_iter$ ) определяет, сколько раз алгоритм обновит коэффициенты. Если итераций слишком мало (например, 100), градиентный спуск может не успеть приблизиться к оптимальным значениям, и модель окажется недообученной. Слишком большое количество итераций (1 000 000) обеспечивает сходимость, но требует значительных вычислительных ресурсов и времени.

В моей задаче проблема усугублялась мультиколлинеарностью и неправильным масштабом данных. Нормализация признаков через `StandardScaler` помогает градиентному спуску работать стабильнее, так как выравнивает масштабы признаков.

### Вывод:

Анализ показал, что себестоимость продукции наиболее сильно зависит от стоимости сырья (коэффициент 0.28,  $p\text{-value} < 0.05$ ). Остальные факторы

(электроэнергия, зарплата) статистически незначимы. Модель объясняет 99.9% дисперсии ( $R^2=0.999$ ), но требует проверки на переобучение. Градиентный спуск сошелся к аналогичным коэффициентам при  $\alpha=0.001$  и 10 000 итераций.

Из 2 моделей рациональнее выбрать модель МНК, потому что она наиболее подходит для датасетов с небольшой выборкой, в нашем случае 19 наблюдений, а также дает возможность с ходу проверить модель по полученным значениям. Достаточно просто решается вопрос мультиколлинеарности – путем исключения спорного предиктора.