



Министерство науки и высшего образования Российской Федерации
Федеральное государственное бюджетное образовательное учреждение
высшего образования
«Московский государственный технический университет имени
Н.Э. Баумана
(национальный исследовательский университет)»
(МГТУ им. Н.Э. Баумана)

ФАКУЛЬТЕТ «Информатика и системы управления»

КАФЕДРА «Программное обеспечение ЭВМ и информационные технологии»

Отчёт по лабораторным работам №1-2 по дисциплине "Методы машинного обучения"

Тема Модель полиномиальной регрессии

Студент Варламова Е. А.

Группа ИУ7-23М

Оценка (баллы) _____

Преподаватели Солодовников Владимир Игоревич

СОДЕРЖАНИЕ

1	Теоретическая часть	3
1.1	Полиномиальная регрессия	3
1.2	Постановка задачи	3
1.2.1	Задача 1	3
1.2.2	Задача 2	4
1.3	Функционал эмпирического риска	4
1.4	Обобщающая способность	5
1.5	Описание алгоритма	5
2	Практическая часть	7
2.1	Выбор средств разработки	7
2.2	Исследование ПО	7
2.2.1	Задача 1	7
2.2.2	Задача 2	11
	СПИСОК ИСПОЛЬЗОВАННЫХ ИСТОЧНИКОВ	20

1 | Теоретическая часть

1.1 Полиномиальная регрессия

Полиномиальная регрессия – это метод восстановления зависимости между независимыми и зависимыми переменными при помощи полиномиальной функции. Он часто используется для приближения нелинейного поведения данных и улучшения качества предсказаний по сравнению с линейной регрессией. Полиномиальная регрессия позволяет уловить сложные взаимосвязи в данных и учитывать нелинейные зависимости.

Целью данной лабораторной работы является изучение модели полиномиальной регрессии.

Для этого необходимо решить следующие задачи:

- формализовать задачу;
- описать алгоритм работы ПО, решающего поставленную задачу;
- привести особенности реализации ПО, решающего поставленную задачу;
- провести исследование зависимости среднеквадратичной ошибки регрессии от степени полинома;
- провести исследование зависимости значения функционала эмпирического риска на обучающей и контрольной выборках от степени полинома.

1.2 Постановка задачи

1.2.1 Задача 1

Создать обучающую выборку с использованием функции

$$y(x) = \theta_1 x + \theta_2 \sin(x) + \theta_3 \quad (1.1)$$

с добавлением шума с нормальным распределением.

Построить модель полиномиальной регрессии, аппроксимирующей данные обучающей выборки. Исходить из того, что степень полинома (начальный закон генерации обучающей выборки) неизвестен. Обучение проводить методом наименьших квадратов.

1.2.2 Задача 2

Феномен Рунге – это эффект нежелательных осцилляций, возникающий при использовании полиномов высоких степеней для интерполяции.

Функция:

$$y(x) = \frac{1}{1 + 25x^2}, x \in [-2, 2] \quad (1.2)$$

Обучающая выборка:

$$S_l : x_i = \frac{4(i-1)}{l-1} - 2, i = 1, \dots, l \quad (1.3)$$

Контрольная выборка:

$$S_k : x_i = \frac{4(i-0.5)}{l-1} - 2, i = 1, \dots, l-1. \quad (1.4)$$

Рассчитать функционал эмпирического риска (функционал качества) для обучающей и контрольной выборок (вывести графики). Оценить обобщающую способность (generalization ability). Найти оптимальную степень полинома для аппроксимации.

1.3 Функционал эмпирического риска

Функционал эмпирического риска (empirical risk functional) используется в машинном обучении для измерения качества модели на обучающей выборке. Он представляет собой среднее значение функции потерь (loss function) на обучающих примерах.

Для задачи регрессии, наши данные состоят из пар (x_i, y_i) , где x_i - входное значение, а y_i - соответствующее целевое значение. Пусть $h(x)$ - модель, а $\ell(h(x), y)$ - функция потерь. Тогда эмпирический риск $R(h)$ может быть записан следующим образом:

$$R(h) = \frac{1}{N} \sum_{i=1}^N \ell(h(x_i), y_i)$$

Здесь N - количество обучающих примеров, и сумма берется по всем парам (x_i, y_i) . Функция потерь $\ell(h(x_i), y_i)$ оценивает разницу между предсказанным значением $h(x_i)$ и истинным значением y_i .

В данной работе используется квадратичная функция потерь, а, соответственно, функционал эмпирического риска равен среднеквадратичной ошибке.

1.4 Обобщающая способность

Обобщающая способность (generalization ability) модели является ее способностью хорошо предсказывать новые, невиданные ранее данные после обучения на имеющемся наборе обучающих данных. Обобщающая способность – это ключевой критерий эффективности модели и важна для того, чтобы избежать переобучения.

Обобщающая способность зависит от сбалансированности модели между точностью на обучающем наборе и способностью обобщаться на новые данные. Модель с хорошей обобщающей способностью сможет давать точные предсказания на новых данных, не привязываясь к особенностям обучающего набора.

Для оценки обобщающей способности модели после обучения ее на обучающем наборе, обычно используют разделение данных на обучающую и тестовую выборки, а также кросс-валидацию. Это помогает оценить, насколько модель способна хорошо предсказывать на новых данных.

1.5 Описание алгоритма

Схема алгоритма, вычисляющего оптимальную степень полинома по обучающей выборке, представлена на рисунке 1.1.

Данный алгоритм используется в обеих задачах, однако во второй задаче среднеквадратичная ошибка вычисляется не на обучающей, а на контрольной выборке.

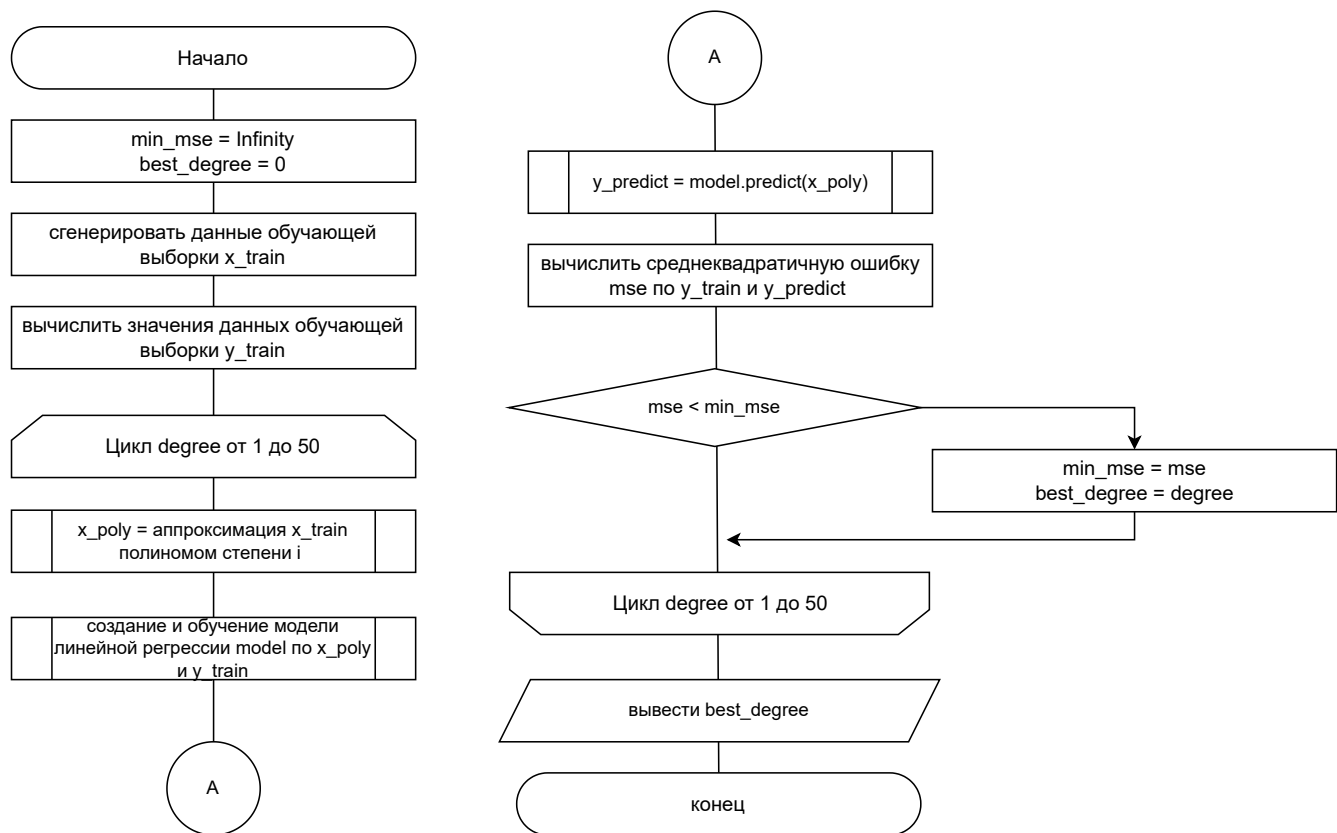


Рис. 1.1: Схема работы алгоритма

2 | Практическая часть

2.1 Выбор средств разработки

В качестве языка программирования был использован язык Python, поскольку этот язык кроссплатформенный и для него разработано огромное количество библиотек и модулей, решающих разнообразные задачи.

В частности, имеются библиотеки, включающие в себя алгоритмы аппроксимации полиномом и линейной регрессии в библиотеке [1].

Для создания графиков была выбрана библиотека matplotlib [2], доступная на языке Python, так как она предоставляет удобный интерфейс для работы с данными и их визуализации.

2.2 Исследование ПО

2.2.1 Задача 1

В листинге 2.1 представлен код, вычисляющий оптимальную степень полинома по обучающей выборке и рисует зависимость среднеквадратичной ошибки модели по обучающей выборке от степени полинома.

Листинг 2.1: код

```
1 import numpy as np
2 from sklearn.preprocessing import PolynomialFeatures
3 from sklearn.linear_model import LinearRegression
4 from sklearn.metrics import mean_squared_error
5 import matplotlib.pyplot as plt
6
7 theta_1 = 2
8 theta_2 = 1
9 theta_3 = 0.5
10
11 np.random.seed(0)
12 X_train = np.linspace(0, 10, 100)
13 y_train = theta_1 * X_train + theta_2 * np.sin(X_train) + theta_3 + np.
    random.normal(0, 0.5, 100)
14
```

```

15 min_mse = float('inf')
16 best_degree = 0
17 degrees = list(range(1, 50))
18 errors = []
19
20 for degree in degrees:
21     poly_features = PolynomialFeatures(degree=degree)
22     X_poly = poly_features.fit_transform(X_train.reshape(-1, 1))
23     model = LinearRegression()
24     model.fit(X_poly, y_train)
25
26     y_pred = model.predict(X_poly)
27
28     mse = mean_squared_error(y_train, y_pred)
29     errors.append(mse)
30     if mse < min_mse:
31         min_mse = mse
32         best_degree = degree
33
34 print(f"Optimal degree of poly: {best_degree}")
35
36 best_poly_features = PolynomialFeatures(degree=best_degree)
37 X_poly_best = best_poly_features.fit_transform(X_train.reshape(-1, 1))
38
39 best_model = LinearRegression()
40 best_model.fit(X_poly_best, y_train)
41
42 plt.plot(degrees, errors)
43 plt.xlabel('Degree of poly')
44 plt.ylabel('Value error')
45 plt.title('Dependence of value error on degree of poly')
46 plt.show()

```

На рисунках 2.1-2.2 показан график тренировочных данных и аппроксимирующего их полинома 13 и 20 степеней соответственно. Видим, что полином 13-ой степени лучше аппроксимирует тренировочные данные, чем 20-ой.

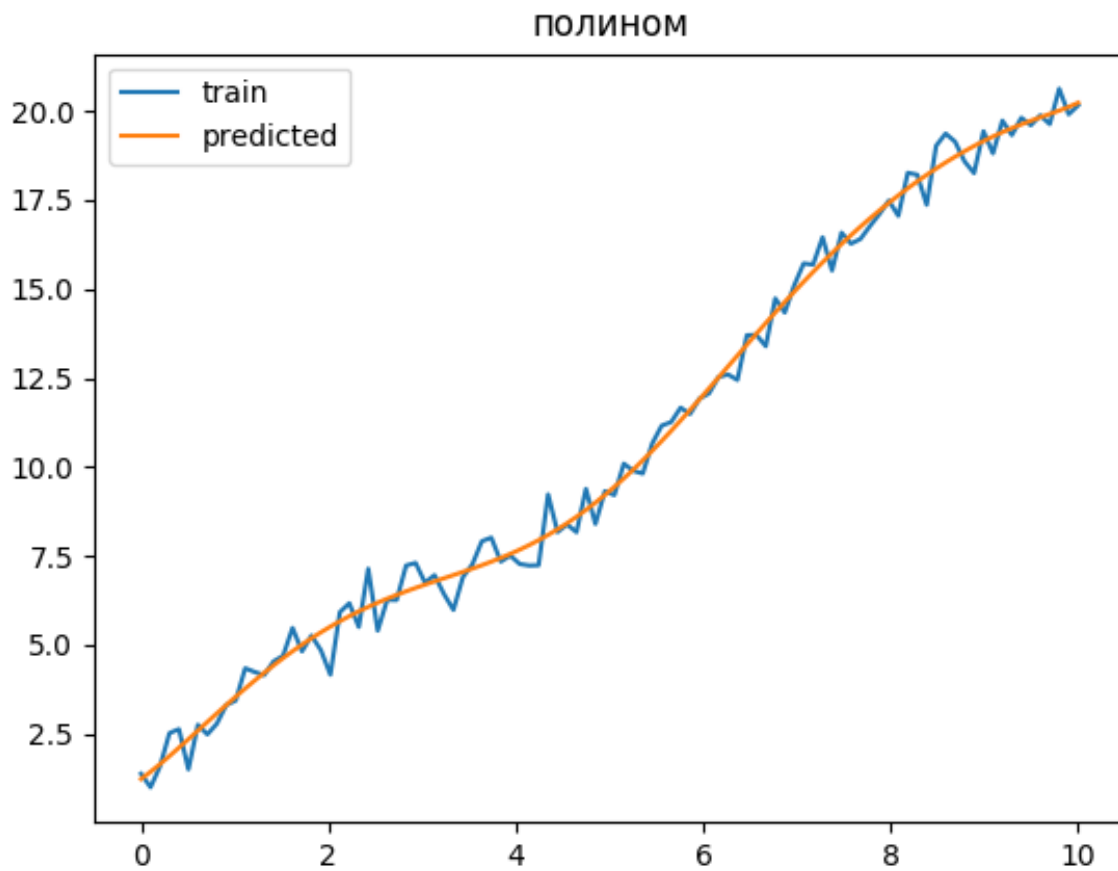


Рис. 2.1:

На рисунке 2.3 показана зависимость значения ошибки от степени полинома. Видно, что увеличение степени полинома необязательно даёт лучшие результаты в смысле уменьшения ошибки.

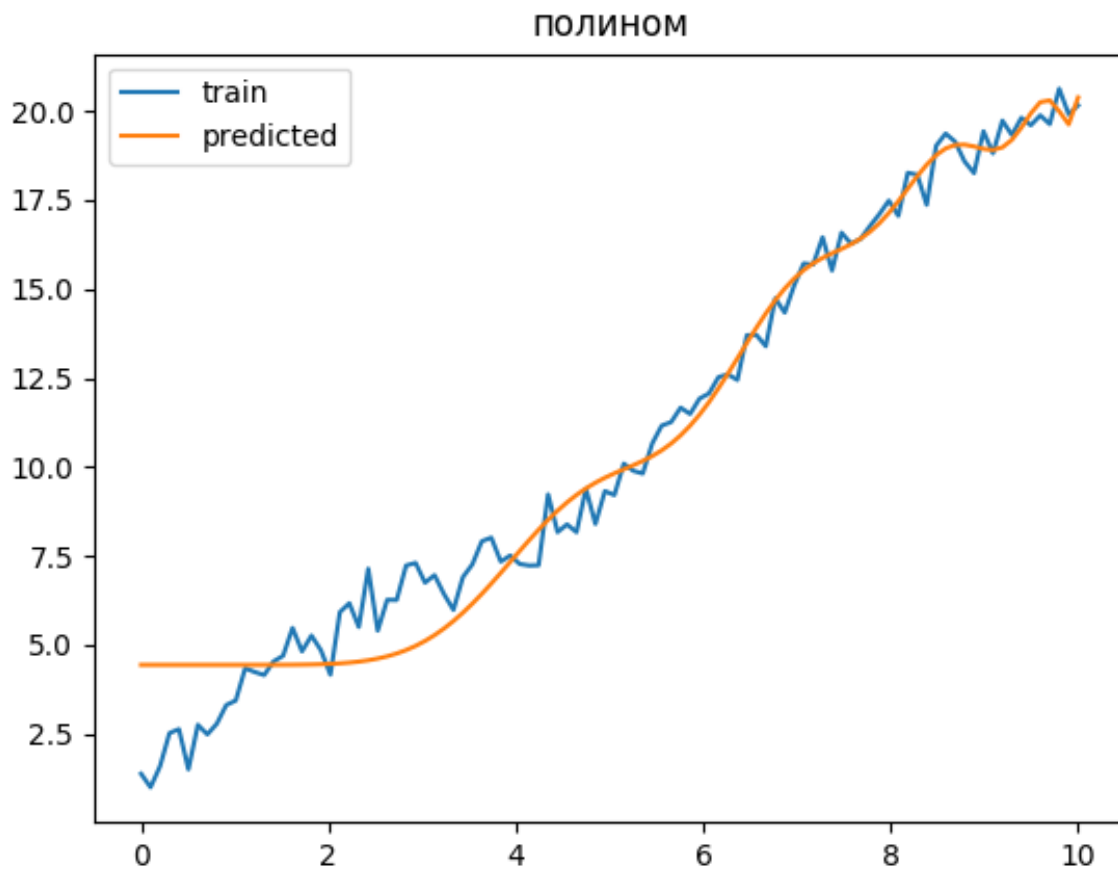


Рис. 2.2:

Результат работы ПО – было вычислено, что оптимальная степень полинома равна 13.



Рис. 2.3:

2.2.2 Задача 2

В листинге 2.2 представлен код, вычисляющий оптимальную степень полинома по контрольной выборке и рисует зависимость среднеквадратичной ошибки модели (функционала эмпирического риска) для обучающей и контрольной выборок от степени полинома.

Листинг 2.2: код

```
1 import numpy as np
2 import matplotlib.pyplot as plt
3 from sklearn.preprocessing import PolynomialFeatures
4 from sklearn.linear_model import LinearRegression
5 from sklearn.metrics import mean_squared_error
6
7 def true_function(x):
8     return 1 / (1 + 25 * x**2)
9
10 l = 21
```

```

11 X_train = np.array([4 * (i - 1) / (l - 1) - 2 for i in range(1, l + 1)]).
    reshape(-1, 1)
12 X_control = np.array([4 * (i - 0.5) / (l - 1) - 2 for i in range(1, l)]).
    reshape(-1, 1)
13 y_train = true_function(X_train)
14 y_control = true_function(X_control)
15
16 def fit_polynomial_regression(X, y, degree):
17     poly_features = PolynomialFeatures(degree=degree)
18     X_poly = poly_features.fit_transform(X)
19     model = LinearRegression()
20     model.fit(X_poly, y)
21     return model, poly_features
22
23
24 def calculate_error(model, poly_features, X, y):
25     X_poly = poly_features.transform(X)
26     y_pred = model.predict(X_poly)
27     return mean_squared_error(y, y_pred)
28
29
30 degrees = np.arange(1, 50)
31 train_errors = []
32 control_errors = []
33
34 for degree in degrees:
35     model, poly_features = fit_polynomial_regression(X_train, y_train,
36     degree)
37     train_error = calculate_error(model, poly_features, X_train, y_train)
38     control_error = calculate_error(model, poly_features, X_control,
39     y_control)
40     train_errors.append(train_error)
41     control_errors.append(control_error)
42
43 plt.plot(degrees, train_errors, label='Train Error')
44 plt.plot(degrees, control_errors, label='Control Error')
45 plt.xlabel('Degree of Polynomial')
46 plt.ylabel('Mean Squared Error')
47 plt.title('Error vs Polynomial Degree')
48 plt.legend()
49 plt.show()
50
51 optimal_degree = degrees[np.argmin(control_errors)]
52 print(f'Optimal polynomial degree for approximation: {optimal_degree}')

```

На рисунках 2.4-2.7 показан график тренировочных и контрольных данных и аппроксимирующих их полиномов 5, 10, 16 и 20 степеней соответственно.

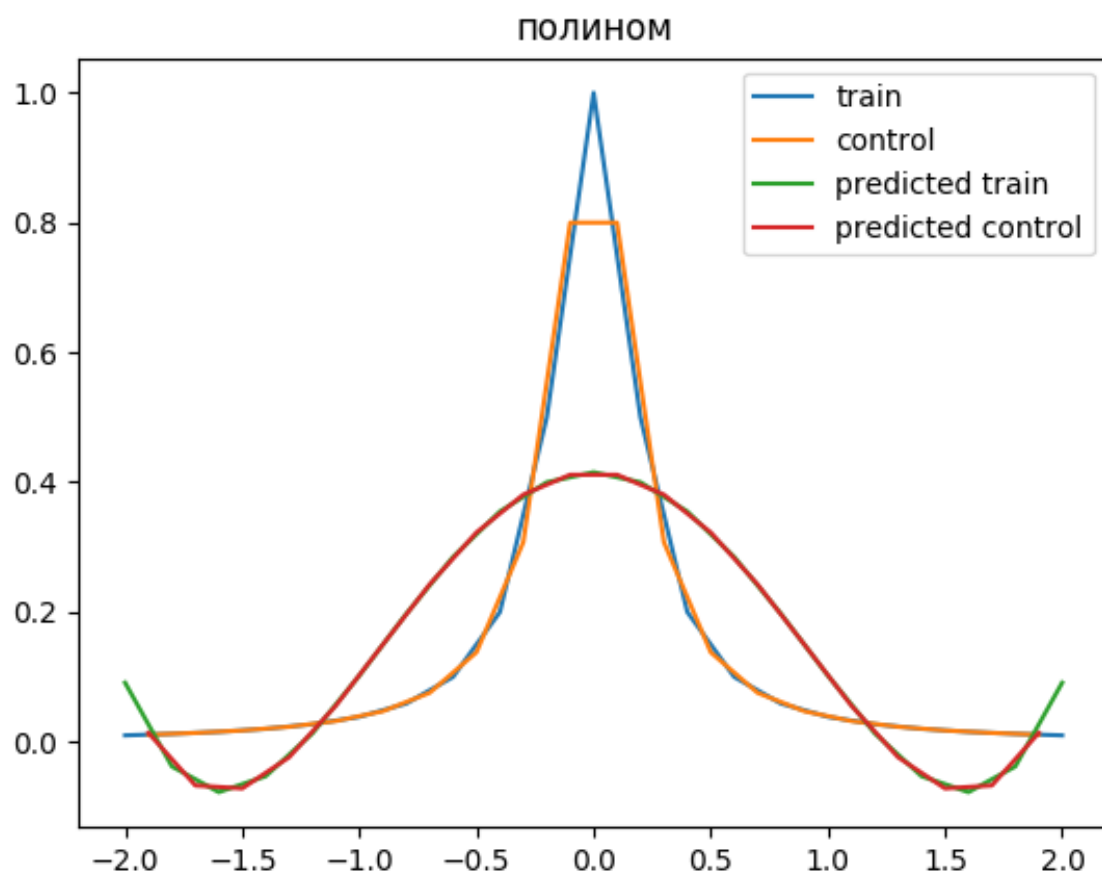


Рис. 2.4:

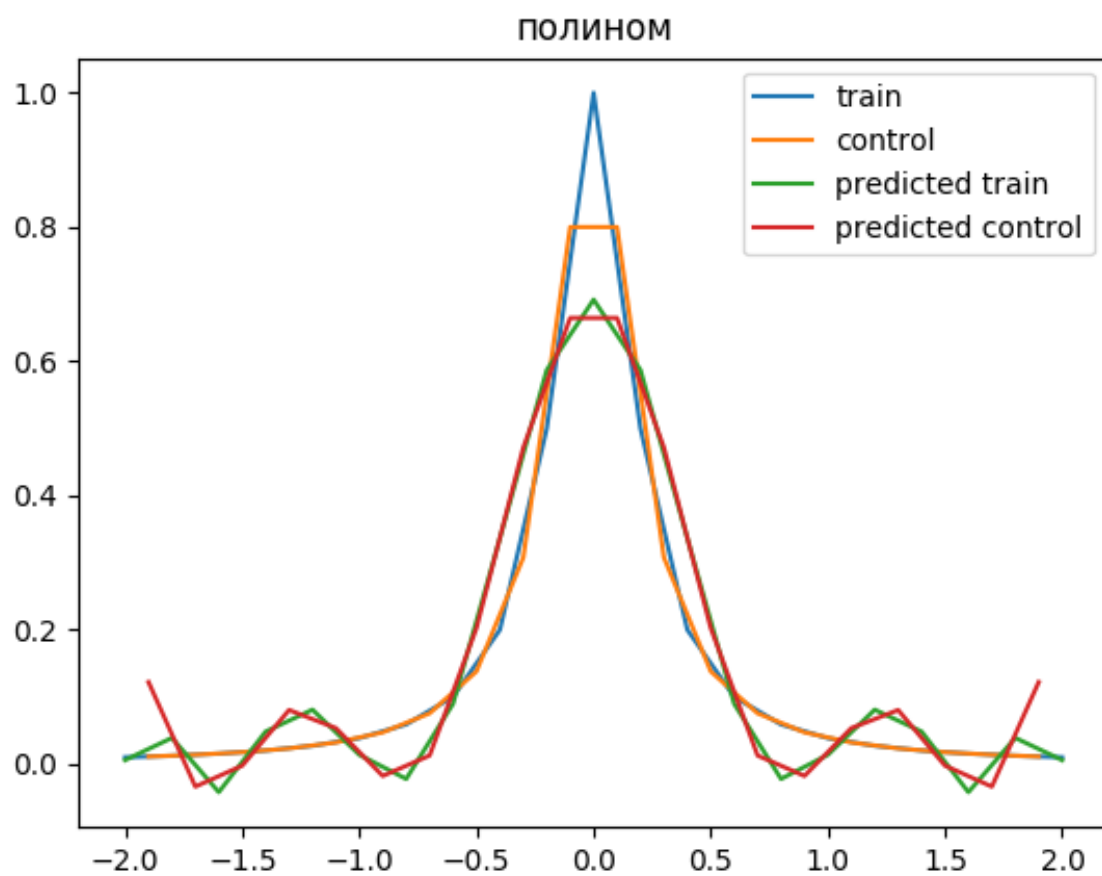


Рис. 2.5:

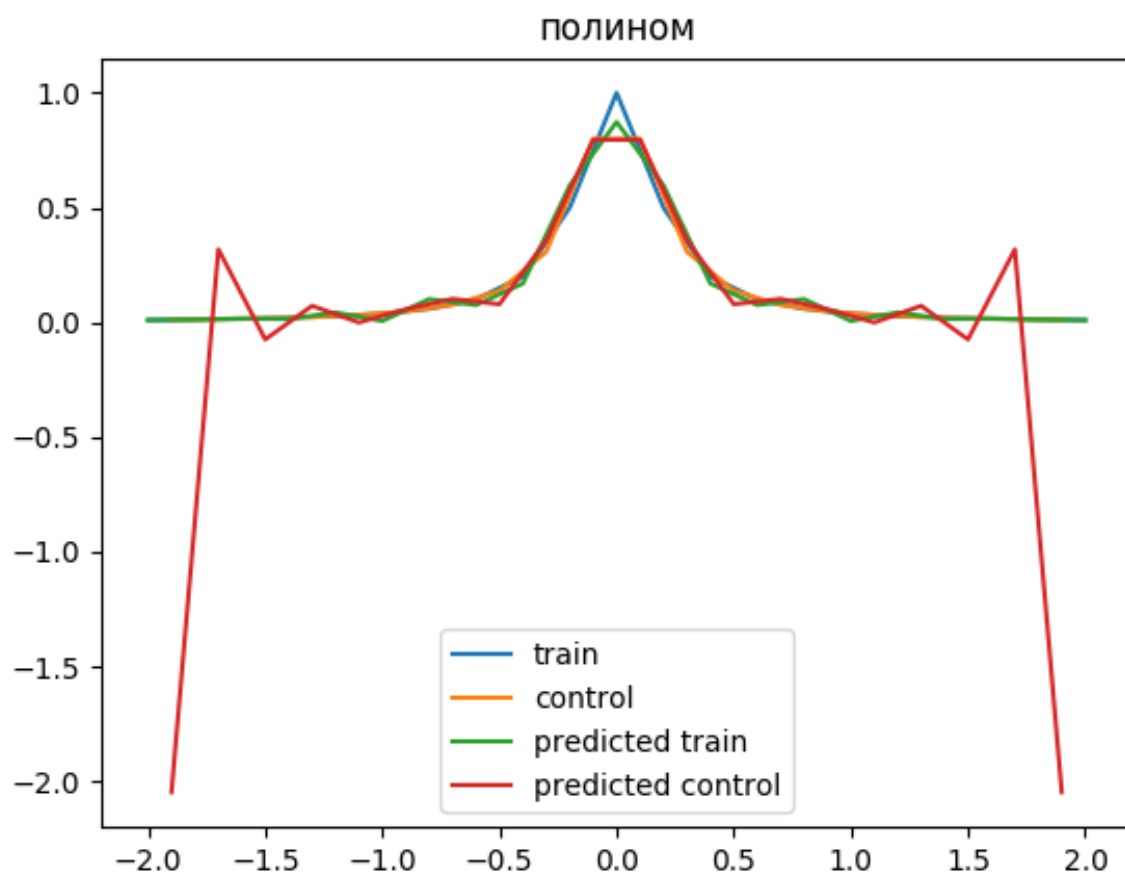


Рис. 2.6:

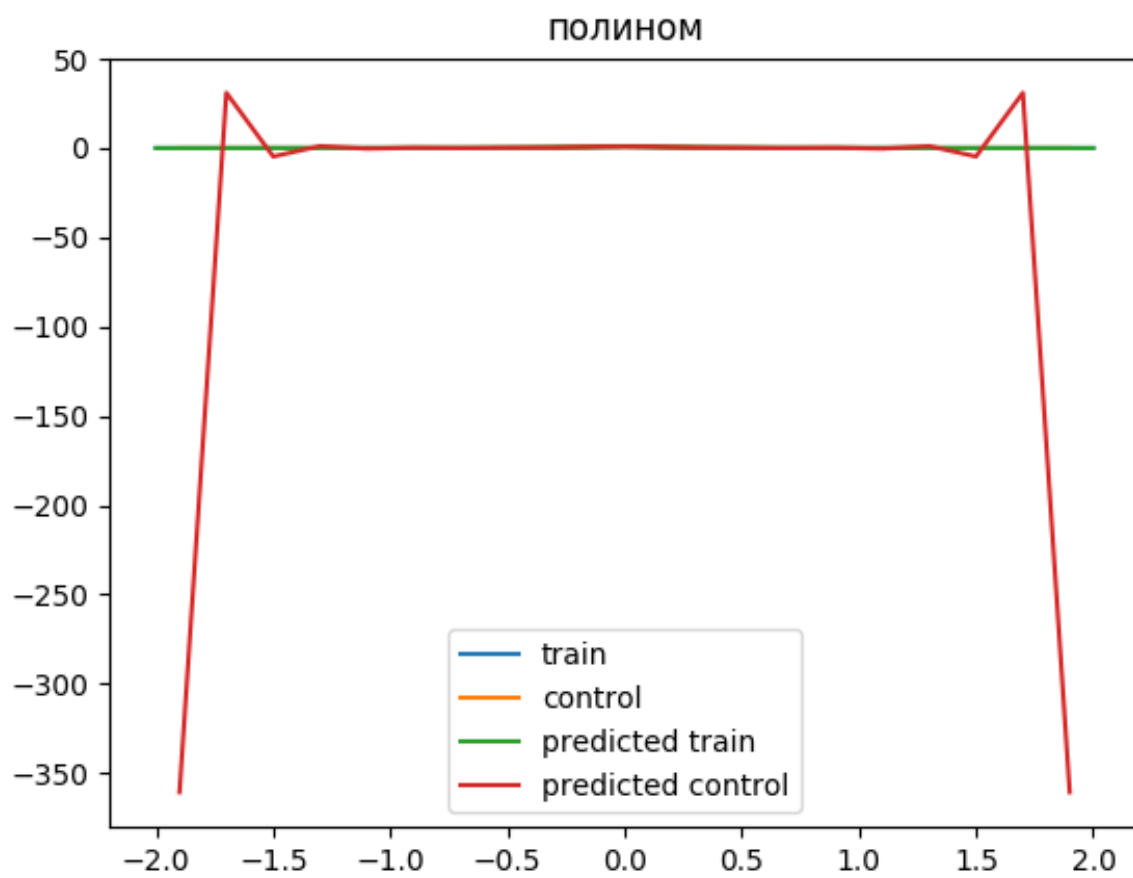


Рис. 2.7:

На рисунке 2.8 показана зависимость значения ошибки от степени полинома. Как и в предыдущей задаче, видно, что увеличение степени полинома необязательно даёт лучшие результаты в смысле уменьшения ошибки.

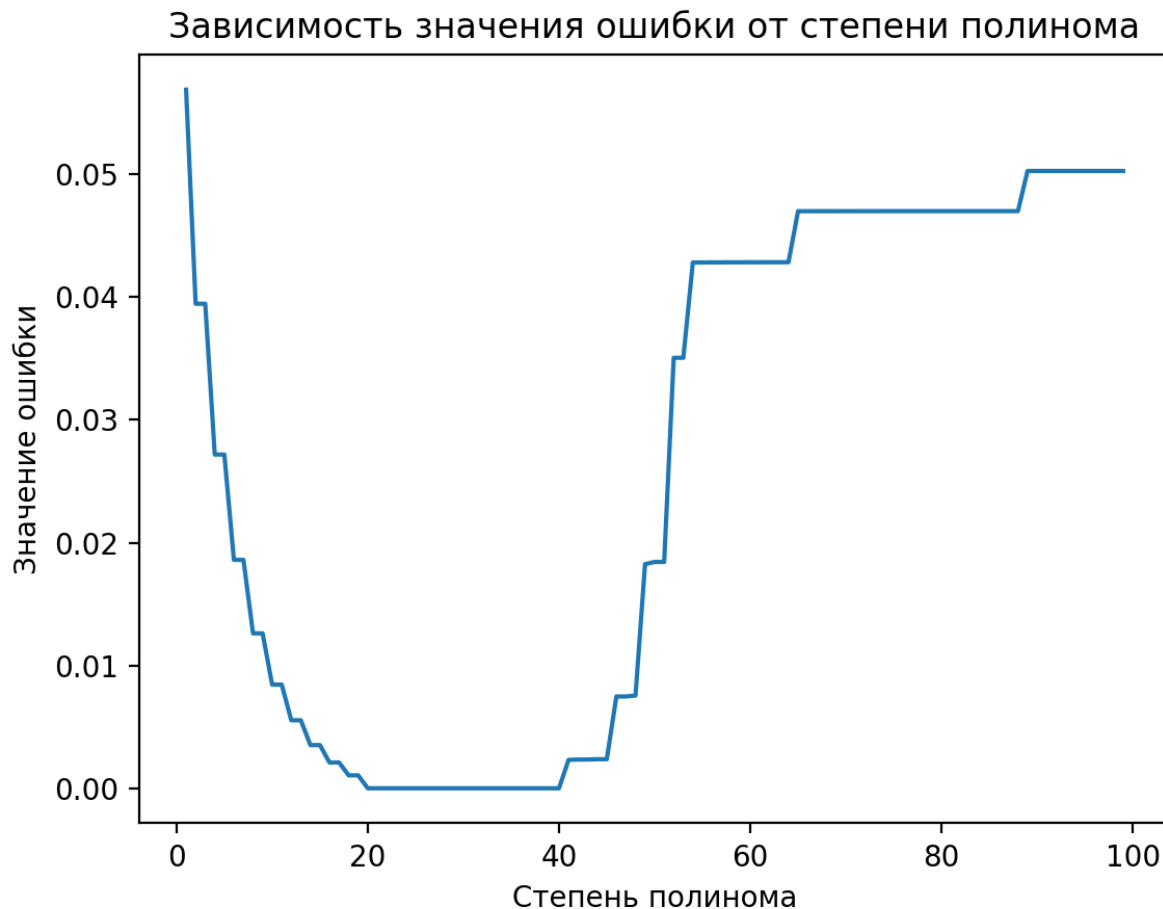


Рис. 2.8:

На рисунке 2.9 показана зависимость значения ошибки от степени полинома для обучающей и контрольной выборок. Видим, что для контрольной выборки с определённого значения степени полинома ошибка стремительно растёт, что демонстрирует эффект Рунге – эффект нежелательных осцилляций или колебаний вблизи крайних точек интерполяции при использовании полиномов высоких степеней.

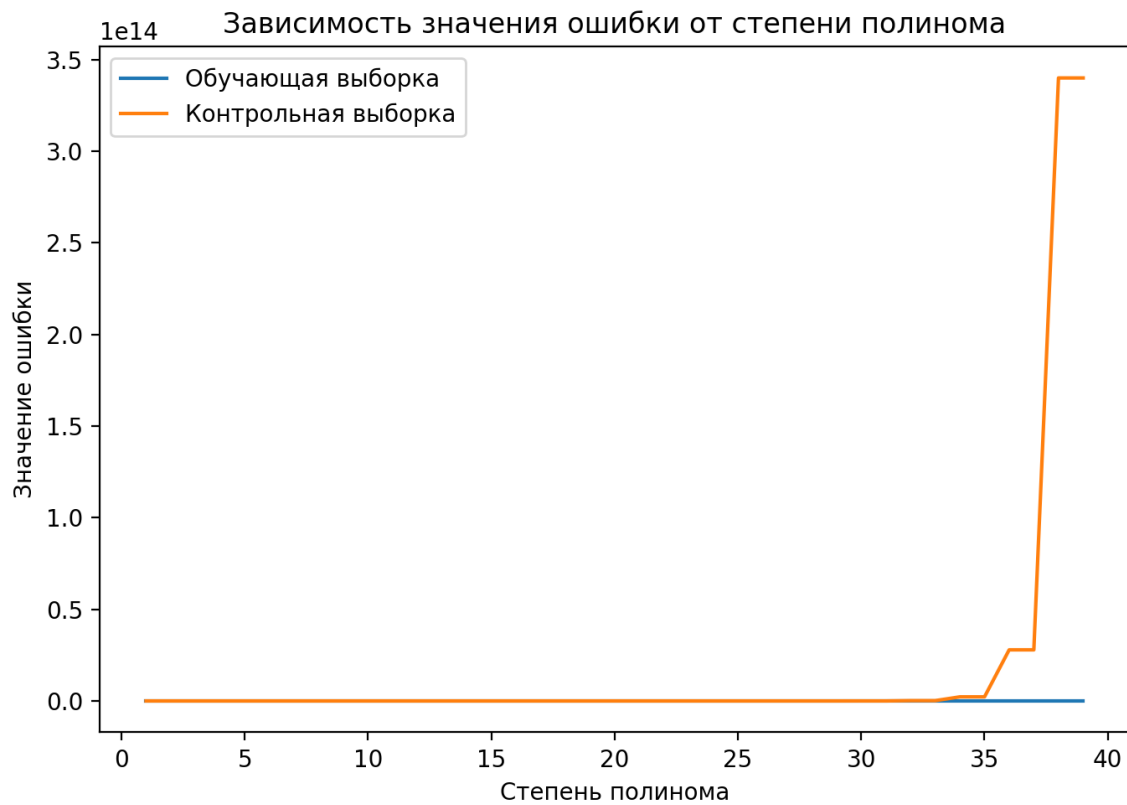


Рис. 2.9:

Рассмотрим интервал низких степеней полинома для получения оптимальной степени полинома на рисунке 2.10.

Видим, что оптимальная степень полинома равна 10.

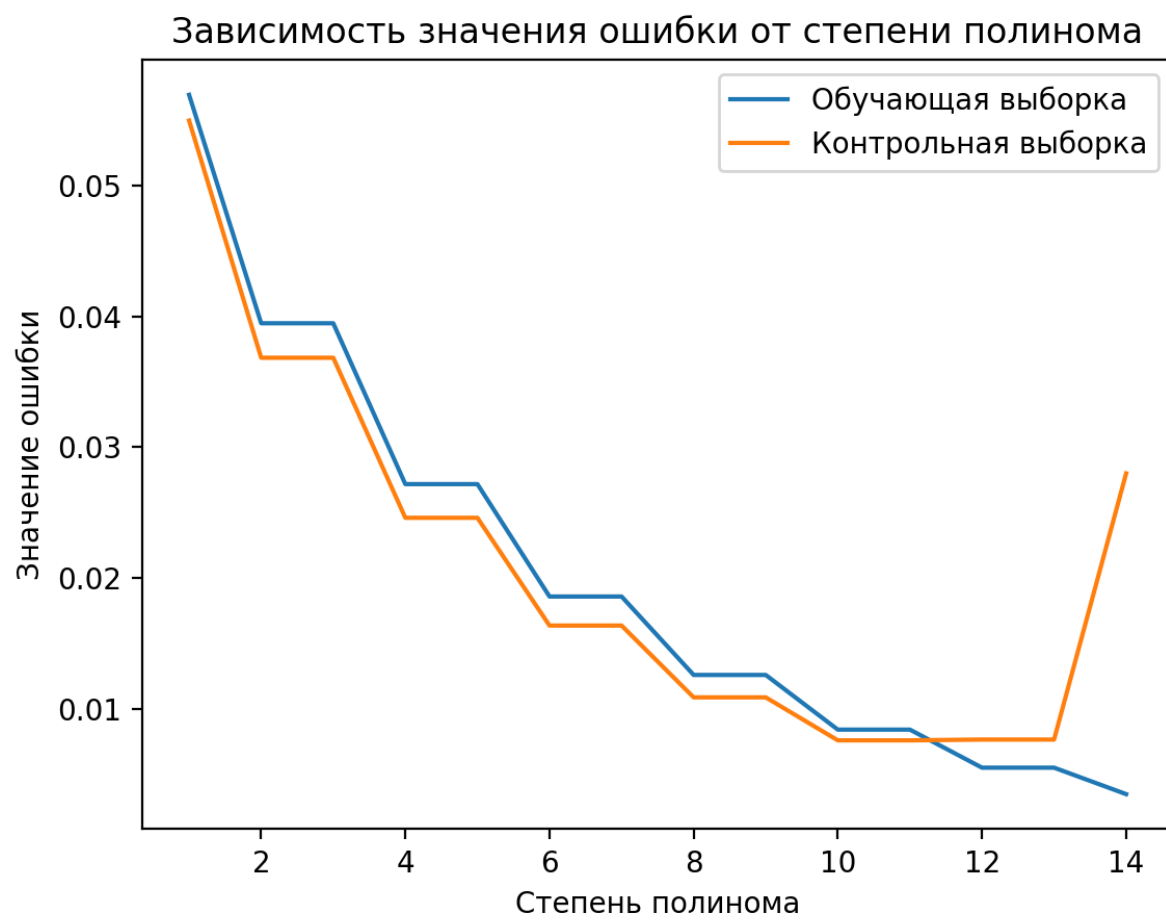


Рис. 2.10:

СПИСОК ИСПОЛЬЗОВАННЫХ ИСТОЧНИКОВ

1. Scikit-learn: Machine learning in Python / F. Pedregosa [и др.]. — 2011.
2. Библиотека визуализации данных matplotlib [Электронный ресурс]. — Режим доступа: URL: <https://matplotlib.org> (дата обращения: 13.12.2023).