

Отчет

Проект по дисциплине «Машинное обучение»

Международный университет Астана
Факультет информационных технологий
Кафедра Data Science DS-23B

Датасет «winequality-white.csv »

Студенты:
Абу Олжас
Муратов Бекнур

Преподаватель: Әбдіқалық Г.Д.

Астана — 2025

Введение

Цель проекта — изучить методы машинного обучения на реальном датасете, реализовать линейную и логистическую регрессию с нуля, провести классификацию, оценить качество моделей и визуализировать результаты.

Мы использовали **датасет качества вина** (Wine Quality Dataset, Kaggle) с числовыми признаками, включающий химические характеристики вина (кислотность, сахар, pH, спирт и др.) и целевую переменную «quality».

Задачи проекта:

- ⑩ Линейная регрессия для предсказания качества вина.
- ⑩ Логистическая регрессия для классификации (например, хорошее/плохое качество).
- ⑩ Сравнение методов классификации (логистическая регрессия vs. решающее дерево).
- ⑩ Анализ метрик (accuracy, precision, recall, F1, ROC AUC).
- ⑩ Демонстрация работы с интерактивными виджетами.

Данные

Источник датасета: Kaggle, Wine Quality Dataset

Ссылка: <https://www.kaggle.com/datasets/uciml/red-wine-quality-cortez-et-al-2009>

Описание признаков:

- ⑩ fixed acidity — постоянная кислотность
- ⑩ volatile acidity — летучая кислотность
- ⑩ citric acid — лимонная кислота
- ⑩ residual sugar — остаточный сахар
- ⑩ chlorides — хлориды
- ⑩ free sulfur dioxide — свободный диоксид серы
- ⑩ total sulfur dioxide — общий диоксид серы
- ⑩ density — плотность
- ⑩ pH — показатель pH
- ⑩ sulphates — сульфаты
- ⑩ alcohol — содержание алкоголя
- ⑩ quality — оценка качества вина (целевой признак)

Подготовка данных:

- ⑩ Проверка на пропуски — отсутствуют.
- ⑩ Дубликаты — удалены.
- ⑩ Стандартизация числовых признаков (MinMaxScaler или StandardScaler).
- ⑩ Разделение на train/test (70/30).

Линейная регрессия (с нуля)

Метод: Реализация градиентного спуска с помощью Numpy.

- ⑩ Функция потерь: MSE (среднеквадратичная ошибка)

- ⑩ Параметры: скорость обучения (learning rate), количество эпох, batch size

Результаты:

- ⑩ График функции потерь по эпохам (сходимость).
- ⑩ Найденные коэффициенты и интерсепт.
- ⑩ Scatter plot: связь признаков (например, alcohol, residual sugar) с quality, линия регрессии, доверительный интервал

Вывод: Линейная регрессия показывает базовую зависимость между химическими свойствами вина и оценкой качества, но может недооценивать нелинейные эффекты.

Логистическая регрессия (с нуля)

Метод:

- ⑩ Сигмоида + log-loss
- ⑩ Градиентный спуск (batch)
- ⑩ Возможная регуляризация L2

Результаты:

- ⑩ Коэффициенты модели
- ⑩ График функции потерь
- ⑩ Проблема: пустые/NaN коэффициенты при исходных данных с категориальными признаками → исправлено сменой датасета на полностью числовой.

Вывод: Логистическая регрессия хорошо работает для бинарной классификации (например, «качество ≥ 7 » — хорошее, <7 — плохое).

Классификация (два подхода)

Подходы:

1. Логистическая регрессия (с нуля)
2. Decision Tree (sklearn)

Метрики:

- ⑩ Accuracy, Precision, Recall, F1-score, Confusion Matrix, ROC AUC
- ⑩ Сравнение моделей: дерево лучше справляется с нелинейными зависимостями, логистическая регрессия быстрее обучается и проще интерпретируется.

Вывод: Для данного датасета с числовыми признаками Decision Tree показывает более высокое качество классификации, однако логистическая регрессия остаётся хорошим базовым методом.

6. Эксперименты и анализ метрик

- ⑩ Изменяли learning rate, epochs и batch size для градиентного спуска.
- ⑩ Наблюдали влияние на сходимость MSE и log-loss.
- ⑩ Decision Tree обучалась без гиперпараметров (базовая версия), можно улучшить с помощью max_depth, min_samples_split.
- ⑩ Анализ метрик показал, что на малом learning rate сходимость медленная, на большом — возможны колебания функции потерь.

7. Интерфейс / демонстрация

- ⑩ Использованы виджеты для:
 - ⑩ загрузки данных
 - ⑩ выбора модели (Линейная регрессия / Логистическая регрессия / Decision Tree)
 - ⑩ задания гиперпараметров (learning rate, epochs, batch size)
- ⑩ Графики потерь, метрик и scatter plot обновляются интерактивно.
- ⑩ Титульный лист отображён в отдельной ячейке (PDF или изображение).

8. Выводы

- ⑩ Мы реализовали линейную и логистическую регрессию с нуля.
- ⑩ Провели классификацию с двумя подходами, сравнили метрики.

- ⑩ На основе экспериментов показали влияние гиперпараметров на качество моделей.
- ⑩ Понимание данных и предварительная подготовка (удаление дубликатов, стандартизация) критично важны.
- ⑩ Проект демонстрирует базовые навыки работы с Python, Numpy, sklearn и визуализацией данных.