

Министерство науки и образования РФ  
Федеральное государственное бюджетное учреждение  
высшего образования  
“Тверской Государственный Технический Университет”  
(ТвГТУ)

Кафедра Программного обеспечения

Отчет по лабораторной работе №2  
По дисциплине: «Анализ больших данных»  
Тема: “Исследовательский анализ данных. Постановка гипотез  
Категориальные данные“

Выполнил:  
студент группы  
Б.ПИН.РИС-21.06  
Миронов М.В.

Проверила:  
старший преподаватель  
кафедры ПО  
Корнеева Е.И.

Тверь 2025 г.

# Содержание

1. Задача .....	1
2. Вариант задачи .....	2
3. Ссылка на код .....	2
4. Описание проделанной работы .....	3
4.1. Датасет seaborn.mpg .....	3
4.1.1. Данные по столбцам .....	3
4.1.2. Гипотезы .....	7
4.1.3. Корреляции .....	8
4.1.4. Стохастический и градиентный спуски .....	9
4.1.5. Вывод .....	9
4.2. Датасет Pulsar .....	10
4.2.1. Данные по столбцам .....	10
4.2.2. Гипотезы .....	14
4.2.3. Корреляции .....	16
4.2.4. Стохастический и градиентный спуски .....	16
4.2.5. Вывод .....	17
5. Вывод .....	17
Список Литературы .....	18

# 1. Задача

1. Посчитать количество строк и столбцов
2. Провести разведочный анализ
  - Для числовых столбцов
    - Доля пропусков
    - Мин/Макс значения
    - Среднее значение
    - Медиана
    - Дисперсия
    - Квантиль 0.1 и 0.9
    - Квартиль 1 и 3
  - Для категориальных столбцов
    - Доля пропусков
    - Количество уникальных значений
    - Мода
3. Сформулировать и проверить минимум 2 статистические гипотезы. Выбор критериев обосновать. Сделать выводы в терминах предметной области.
4. Категориальные данные необходимые для задачи, закодировать методами OneHotEncoding или LabelEncoding, в случае если это требуется.
5. Построить таблицу корреляции признаков и целевого столбца.
6. Реализовать стохастический и обычный градиентный спуск.

Выше описанное требуется реализовать для датасета `seaborn.mpg` и датасета данного по варианту.

Сложность “Well-done”

## **2. Вариант задачи**

Binary Classification with a Tabular Pulsar Dataset [1]

## **3. Ссылка на код**

Дамп таблиц не будет предоставлен, так как БД содержит слишком много данных для системы elearning. Взамен будет предоставлен DDL.

<https://github.com/NydusBorn/big-data>

## 4. Описание проделанной работы

Значительная часть требуемых данных автоматически получена штатными средствами rpycharm

### 4.1. Датасет seaborn.mpg

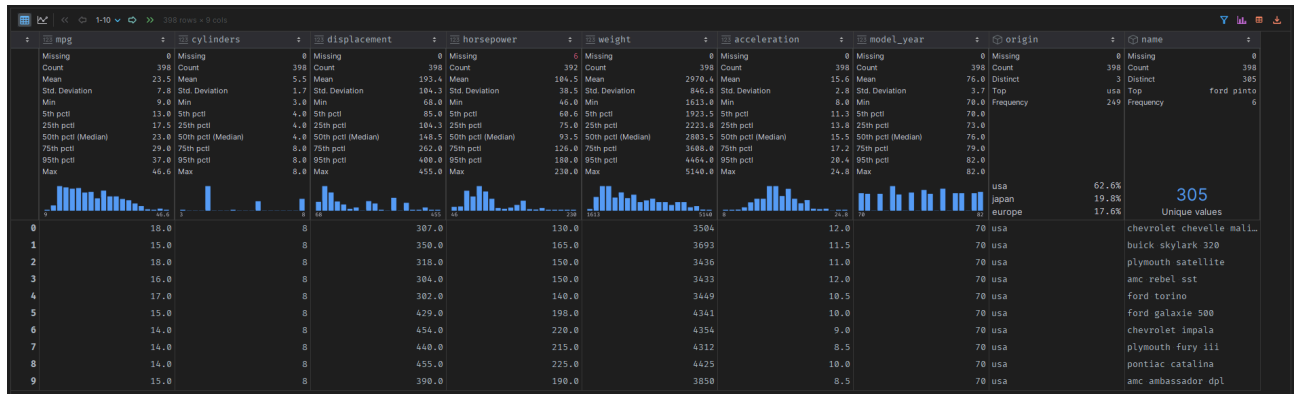


Figure 1: Предварительные данные

Соответственно здесь можно обнаружить следующие параметры:

- Количество строк и столбцов: 398 строк и 9 столбцов
- Количество пропусков
- Среднее для числовых
- Мин/Макс для числовых
- Медиана для числовых
- Количество уникальных значений и мода для категориальных

Остальные данные (дисперсия, квантили и квартили посчитаем сами)

#### 4.1.1. Данные по столбцам

##### 1. mpg

- Пропуски: 0
- Мин: 9
- Макс: 46.6

- Среднее: 23.5
- Медиана: 23
- Дисперсия: 61.1
- Квантиль 0.1: 14
- Квантиль 0.9: 34.3
- Квартиль 1: 17.5
- Квартиль 3: 29

## 2. cylinders

- Пропуски: 0
- Мин: 3
- Макс: 8
- Среднее: 5.5
- Медиана: 4
- Дисперсия: 2.9
- Квантиль 0.1: 4
- Квантиль 0.9: 8
- Квартиль 1: 4
- Квартиль 3: 8

## 3. displacement

- Пропуски: 0
- Мин: 68
- Макс: 455
- Среднее: 193.4
- Медиана: 148.5
- Дисперсия: 10872.2
- Квантиль 0.1: 90
- Квантиль 0.9: 350

- Квартиль 1: 104.25
- Квартиль 3: 262

#### 4. horsepower

- Пропуски: 6 (1.5%)
- Мин: 46
- Макс: 230
- Среднее: 104.5
- Медиана: 93.5
- Дисперсия: 1481.5
- Квантиль 0.1: 67
- Квантиль 0.9: 157.7
- Квартиль 1: 75
- Квартиль 3: 126

#### 5. weight

- Пропуски: 0
- Мин: 1613
- Макс: 5140
- Среднее: 2970.4
- Медиана: 2803.5
- Дисперсия: 717141
- Квантиль 0.1: 1988.5
- Квантиль 0.9: 4275.2
- Квартиль 1: 2223.75
- Квартиль 3: 3608

#### 6. acceleration

- Пропуски: 0
- Мин: 8

- Макс: 24.8
- Среднее: 15.6
- Медиана: 15.5
- Дисперсия: 7.6
- Квантиль 0.1: 12
- Квантиль 0.9: 19
- Квартиль 1: 13.8
- Квартиль 3: 17.2

#### 7. model\_year

- Пропуски: 0
- Мин: 70
- Макс: 82
- Среднее: 76
- Медиана: 76
- Дисперсия: 13.6
- Квантиль 0.1: 71
- Квантиль 0.9: 81
- Квартиль 1: 73
- Квартиль 3: 79

#### 8. origin

- Пропуски: 0
- Уникальных: 3
- Мода: usa

#### 8. Name

- Пропуски: 0
- Уникальных: 305
- Мода: ford pinto



### 4.1.2. Гипотезы

Основным интересным для нас параметром предположим количество лошадиных сил, соответственно остальные данные будем считать дополнительными.

1. Предположим нулевую гипотезу: Год выпуска не имеет влияния на количество лошадиных сил.

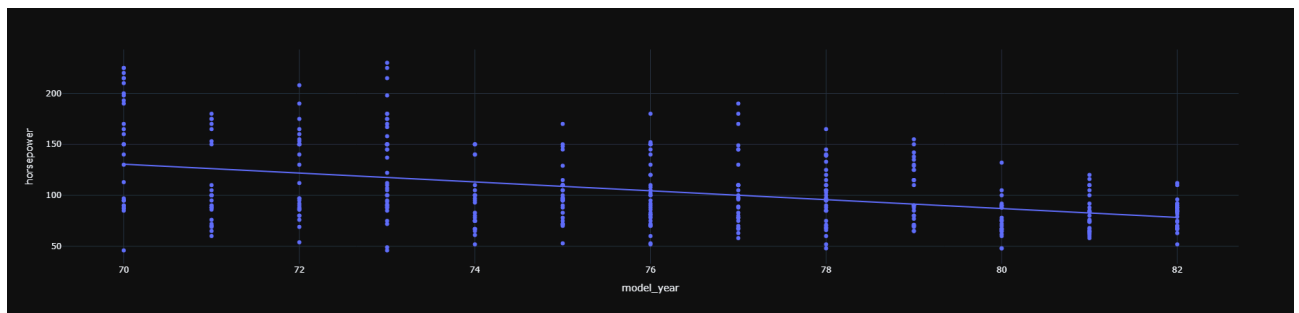


Figure 2: x - год выпуска, y - количество лошадиных сил

Как мы видим наша нулевая гипотеза не подтвердилась, можно утверждать что с ростом года выпуска количество лошадиных сил уменьшается.

2. Предположим альтернативную гипотезу: Количество цилиндров имеет положительное влияние на количество лошадиных сил.

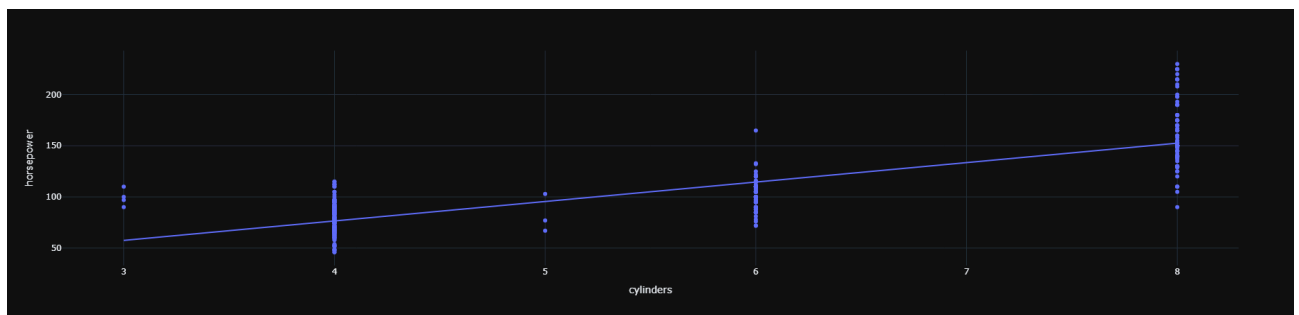


Figure 3: x - количество цилиндров, y - количество лошадиных сил

Как мы видим наша гипотеза подтвердилась.

### 4.1.3. Корреляции

Для этого дополнительно закодируем страны в OneHotEncoding.

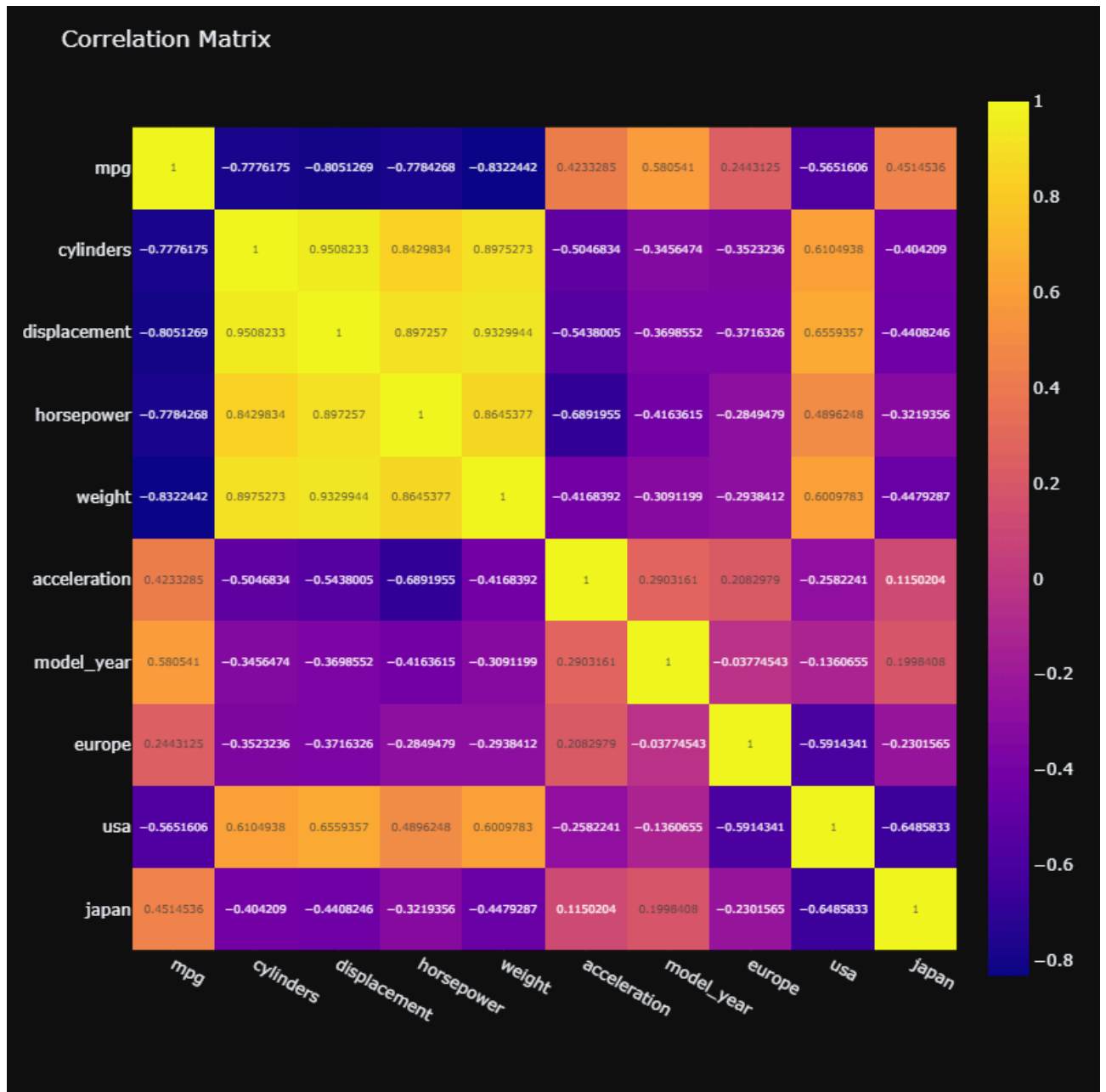


Figure 4: Корреляционная матрица

#### 4.1.4. Стохастический и градиентный спуски

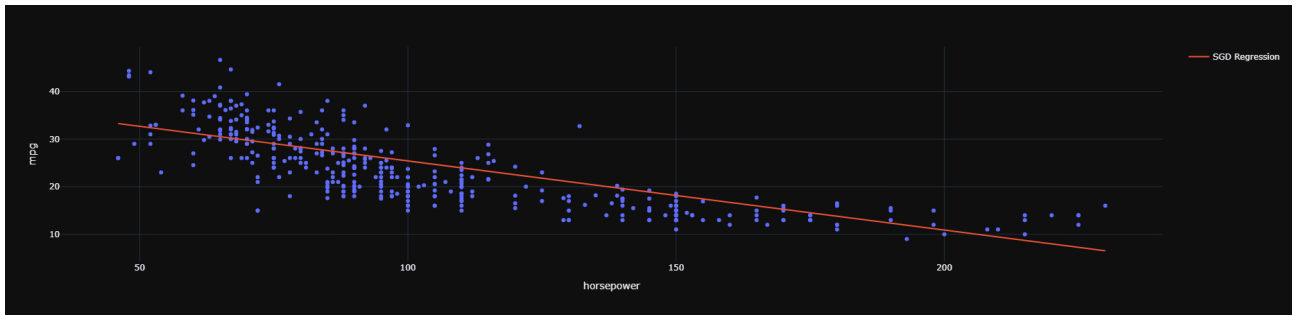


Figure 5: Стохастический градиентный спуск

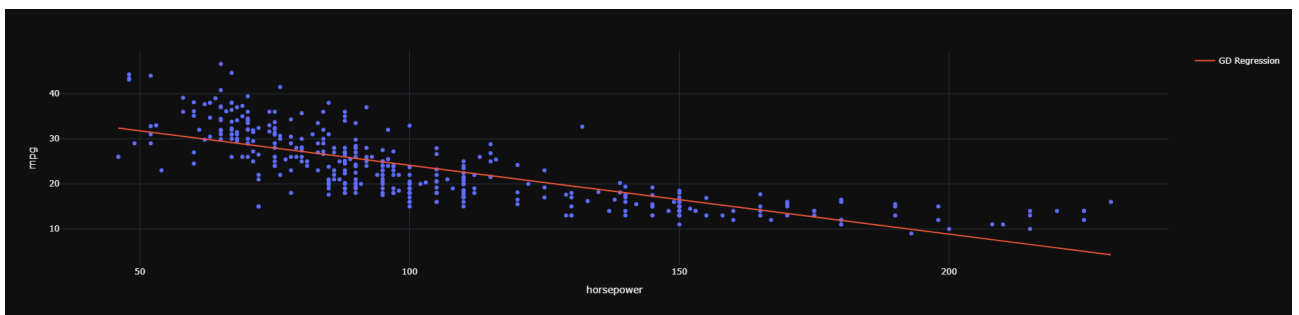


Figure 6: Обычный градиентный спуск

#### 4.1.5. Вывод

Мы провели статистический анализ датасета `seaborn.mpg`. Как мы выяснили количество лошадиных сил выше в США, в моделях выпущенных раньше, при этом количество цилиндров увеличивается и объем двигателя растет, также есть положительная корреляция с весом (более тяжелые автомобили имеют больше лошадиных сил), но автомобили с высоким количеством лошадиных сил имеют низкое ускорение.

## 4.2. Датасет Pulsar

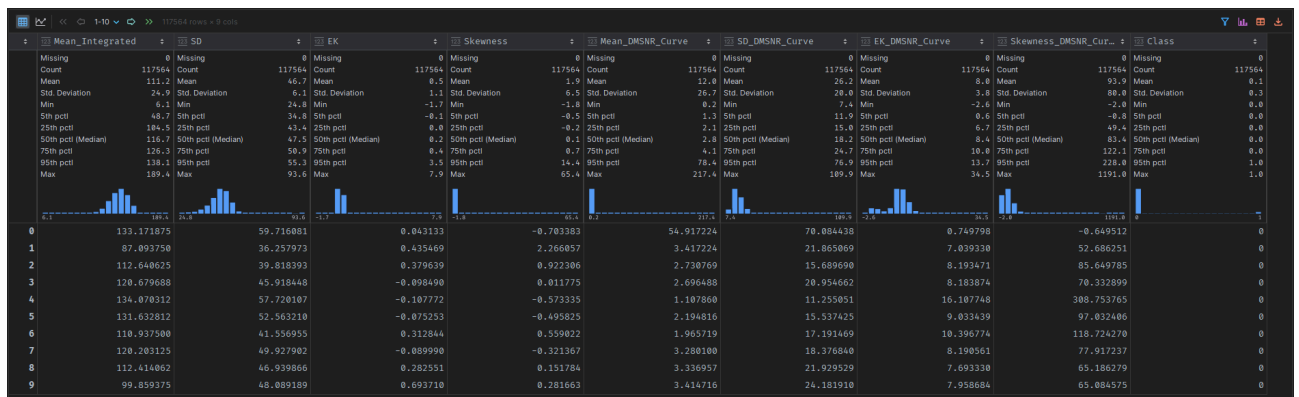


Figure 7: Предварительные данные

Количество строк и столбцов: 117564 строк и 9 столбцов

### 4.2.1. Данные по столбцам

#### 1. Mean\_Integrated

- Пропуски: 0
- Мин: 6.1
- Макс: 189.4
- Среднее: 111.2
- Медиана: 116.7
- Дисперсия: 620.3
- Квантиль 0.1: 85.6
- Квантиль 0.9: 134.2
- Квартиль 1: 104.5
- Квартиль 3: 126.3

#### 2. SD

- Пропуски: 0
- Мин: 24.8
- Макс: 93.6

- Среднее: 46.7
- Медиана: 47.5
- Дисперсия: 37.2
- Квантиль 0.1: 38.2
- Квантиль 0.9: 53.6
- Квартиль 1: 43.4
- Квартиль 3: 50.8

### 3. ЕК

- Пропуски: 0
- Мин: -1.7
- Макс: 7.9
- Среднее: 0.5
- Медиана: 0.2
- Дисперсия: 1.27
- Квантиль 0.1: -0.07
- Квантиль 0.9: 0.79
- Квартиль 1: 0.04
- Квартиль 3: 0.39

### 4. Skewness

- Пропуски: 0
- Мин: -1.8
- Макс: 65.4
- Среднее: 1.9
- Медиана: 0.1
- Дисперсия: 42.5
- Квантиль 0.1: -0.39
- Квантиль 0.9: 2.47

- Квартиль 1: -0.18
- Квартиль 3: 0.69

#### 5. Mean\_DMSNR\_Curve

- Пропуски: 0
- Мин: 0.2
- Макс: 217.4
- Среднее: 12
- Медиана: 2.8
- Дисперсия: 713.9
- Квантиль 0.1: 1.62
- Квантиль 0.9: 34.1
- Квартиль 1: 2.09
- Квартиль 3: 4.12

#### 6. SD\_DMSNR\_Curve

- Пропуски: 0
- Мин: 7.4
- Макс: 109.9
- Среднее: 26.2
- Медиана: 18.2
- Дисперсия: 401.6
- Квантиль 0.1: 12.9
- Квантиль 0.9: 62.7
- Квартиль 1: 14.9
- Квартиль 3: 24.7

#### 7. EK\_DMSNR\_Curve

- Пропуски: 0
- Мин: -2.6

- Макс: 34.5
- Среднее: 8
- Медиана: 8.4
- Дисперсия: 14.7
- Квантиль 0.1: 1.82
- Квантиль 0.9: 11.8
- Квартиль 1: 6.74
- Квартиль 3: 10.0

#### 8. Skewness\_DMSNR\_Curve

- Пропуски: 0
- Мин: -2
- Макс: 1191
- Среднее: 93.9
- Медиана: 83.4
- Дисперсия: 6393.9
- Квантиль 0.1: 2.24
- Квантиль 0.9: 174.3
- Квартиль 1: 49.4
- Квартиль 3: 122

#### 9. Class

- Пропуски: 0
- Мин: 0
- Макс: 1
- Среднее: 0.1
- Медиана: 0
- Дисперсия: 0.08
- Квантиль 0.1: 0

- Квантиль 0.9: 0
- Квартиль 1: 0
- Квартиль 3: 0

Технически это категориальный или бинарный признак, с 2 уникальными значениями и модой 0

#### 4.2.2. Гипотезы

Основным интересным для нас параметром предположим Skewness (наклон пульсара), соответственно остальные данные будем считать дополнительными.

1. Предположим нулевую гипотезу: Mean\_Integrated не имеет влияния на Skewness.

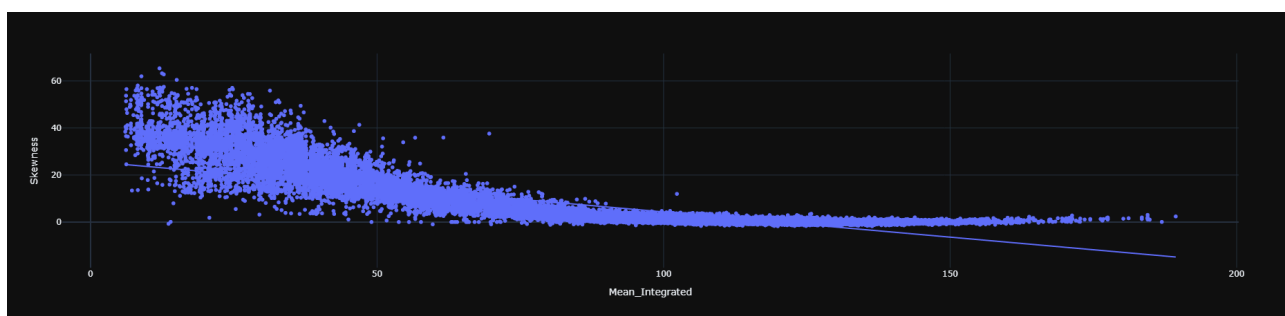


Figure 8: x - Mean Integrated, y - Skewness

Как мы видим наша нулевая гипотеза не подтвердилась, можно утверждать что с ростом Mean Integrated Skewness уменьшается.

2. Предположим альтернативную гипотезу: SD имеет негативное влияние на Skewness.



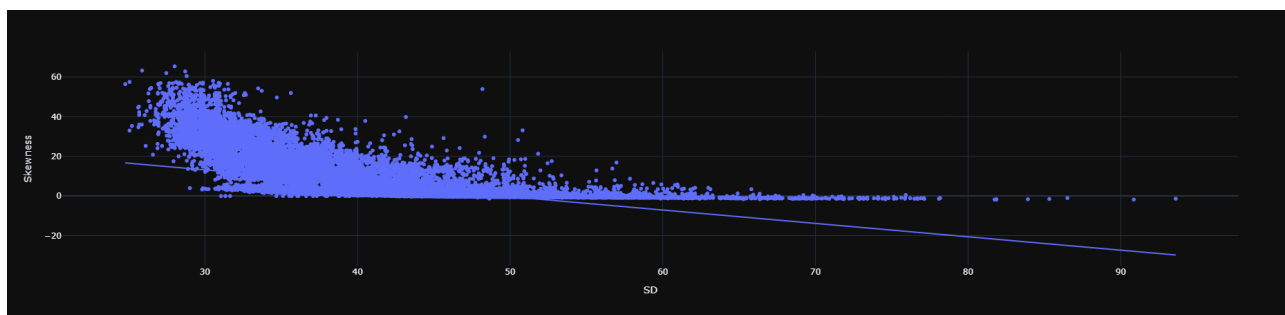


Figure 9: x - SD, y - Skewness

Как мы видим наша гипотеза подтвердилась.

### 4.2.3. Корреляции

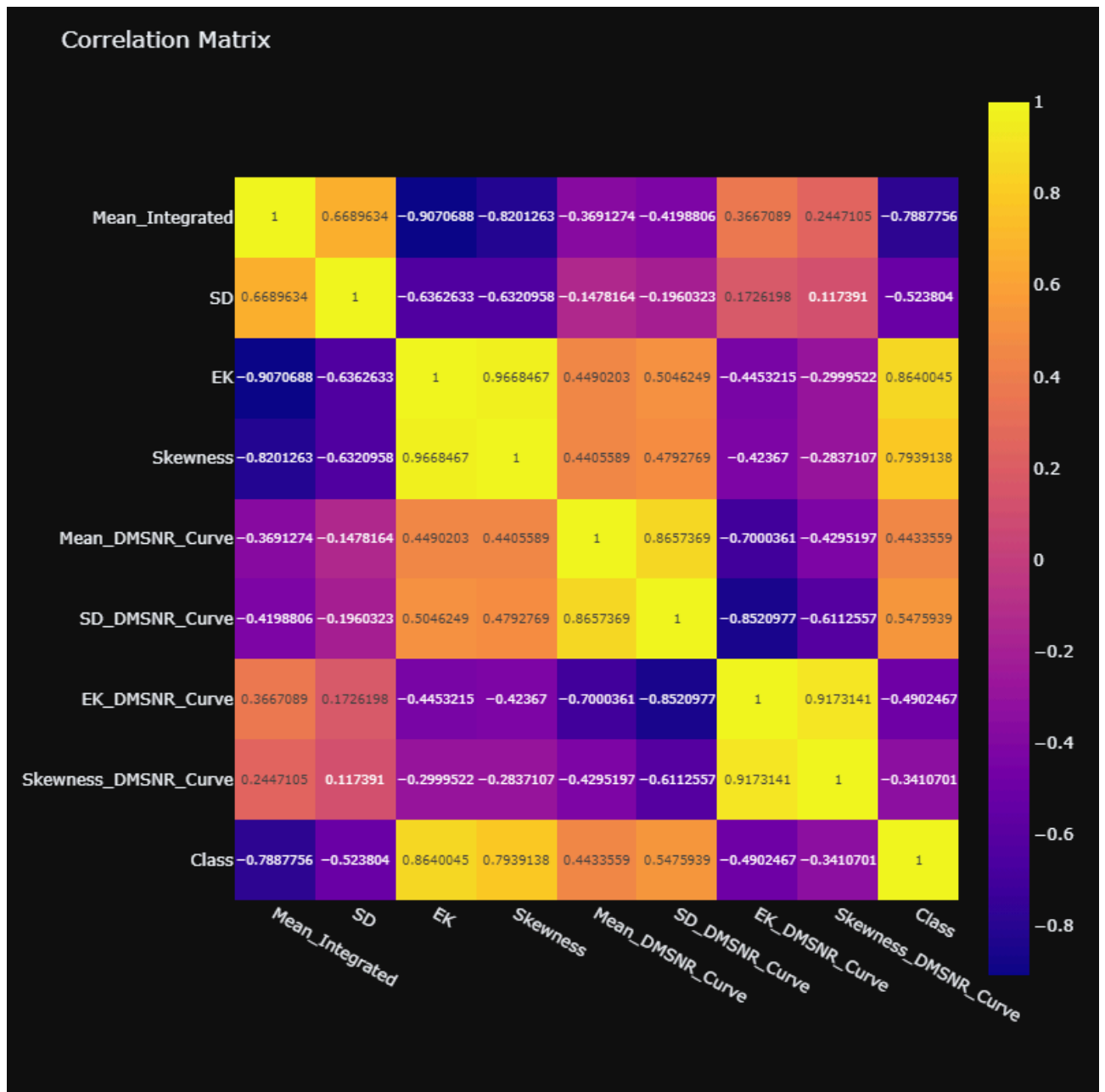


Figure 10: Корреляционная матрица

### 4.2.4. Стохастический и градиентный спуски

Проведем на основании у - Skewness и х - SD

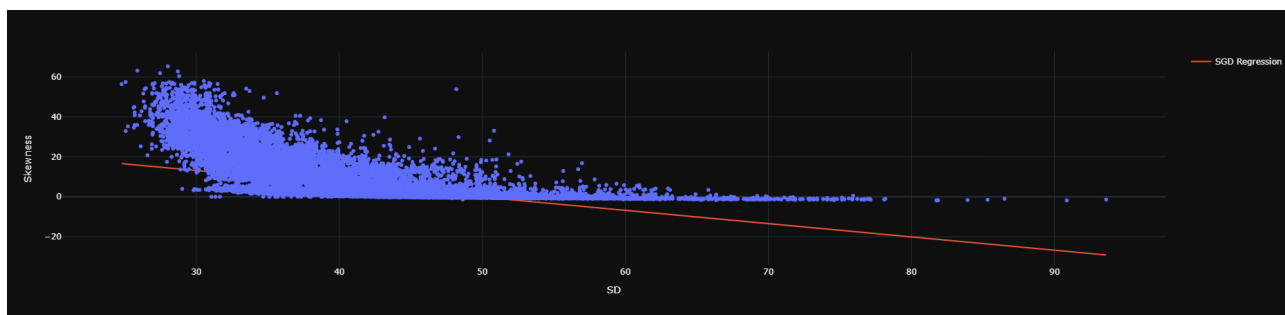


Figure 11: Стохастический градиентный спуск

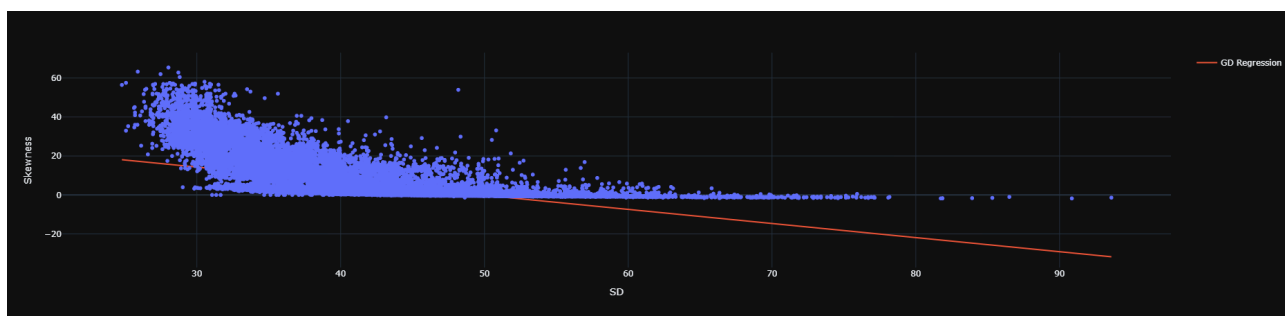


Figure 12: Обычный градиентный спуск

#### 4.2.5. Вывод

Мы проанализировали датасет Pulsar. Как мы увидели все параметры так или иначе связаны с Skewness, кроме Skewness\_DMSNR\_Curve, тк на нём довольно низкий уровень корреляции.

### 5. Вывод

Мы использовали средства python и IDE pycharm для анализа датасетов seaborn.mpg и Pulsar.

## Список Литературы

- [1] “Binary Classification with a Tabular Pulsar Dataset.” Accessed: Jan. 27, 2025.  
[Online]. Available: <https://www.kaggle.com/competitions/playground-series-s3e10/code>