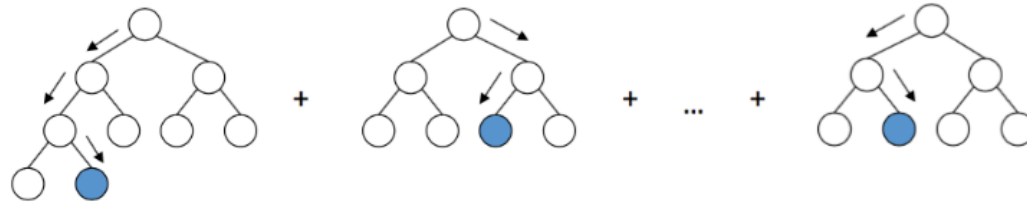

Решение задачи многомерной регрессии методом Gradient Boosting (CatBoost)



Матвеев И. Ю.

Необходимо выполнить следующие пункты:

- Провести предварительный анализ данных;
- Построить регрессионную модель. Обосновать выбор модели;
- Сравнить ошибки на тестовых и тренировочных данных. Использовать метрику R^2 ;
- Построить зависимости ошибок на тестовых и тренировочных данных.

Dataset состоит из матриц признаков X и значений Y :

$$X = \{s_{mt}, s_{mq}, d, h_p\} \text{ и } Y = \{QW, DP\}$$

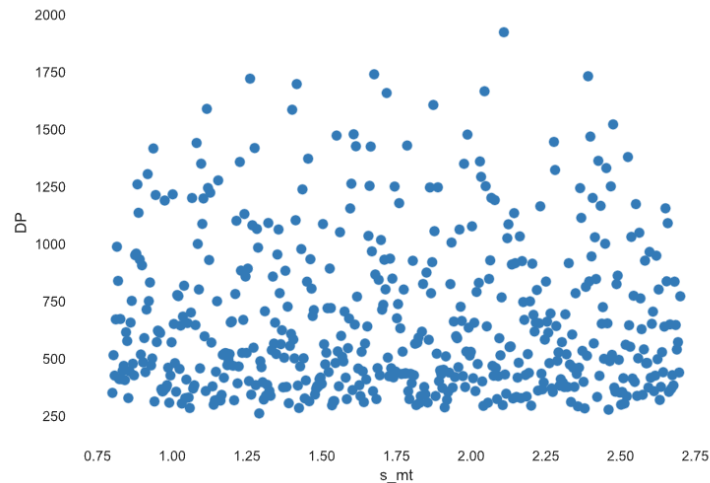
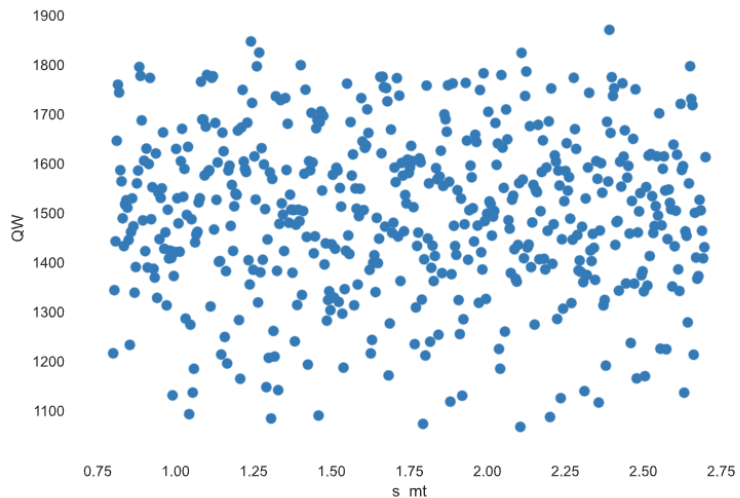
- где каждый объект матриц это вектор $\bar{v}_i, \in \mathbb{R}^{500 \times 1}$

Диапазон значений входных данных представлен в таблице:

	s_{mt}	s_{mq}	d	h_p
Min	0.8	0.8	1.0	4.0
Max	2.7	2.1	3.0	10.0

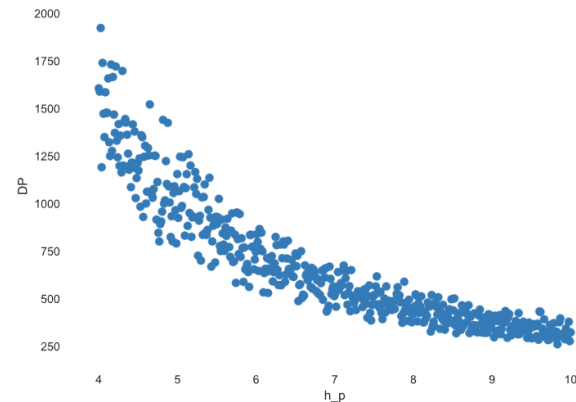
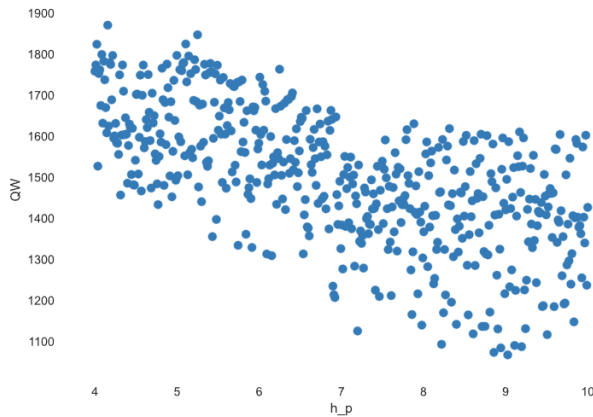
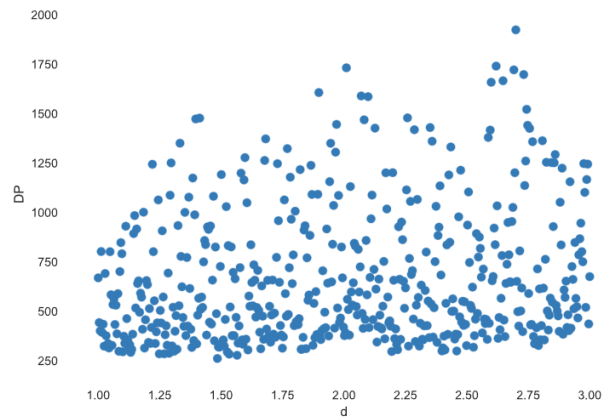
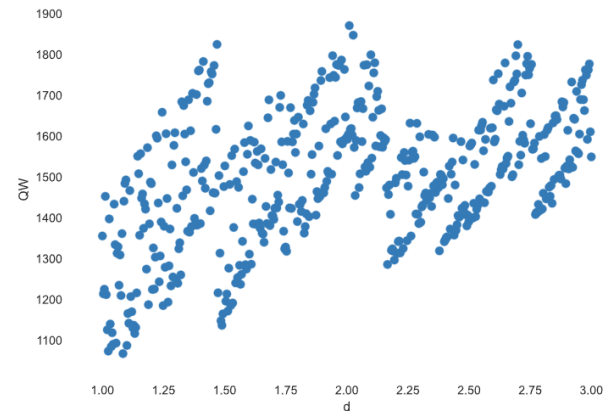
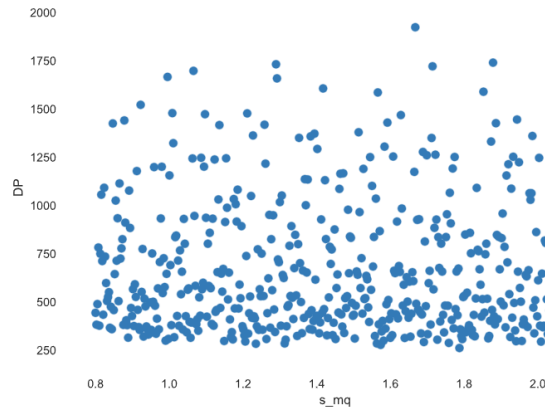
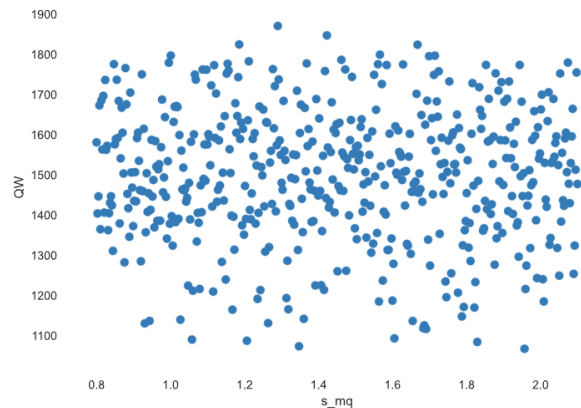
- Среди генеральной совокупности не обнаружены пустые ячейки значений переменных и ложные типы данных;
- Ниже продемонстрированы зависимости целевых переменных от признаков

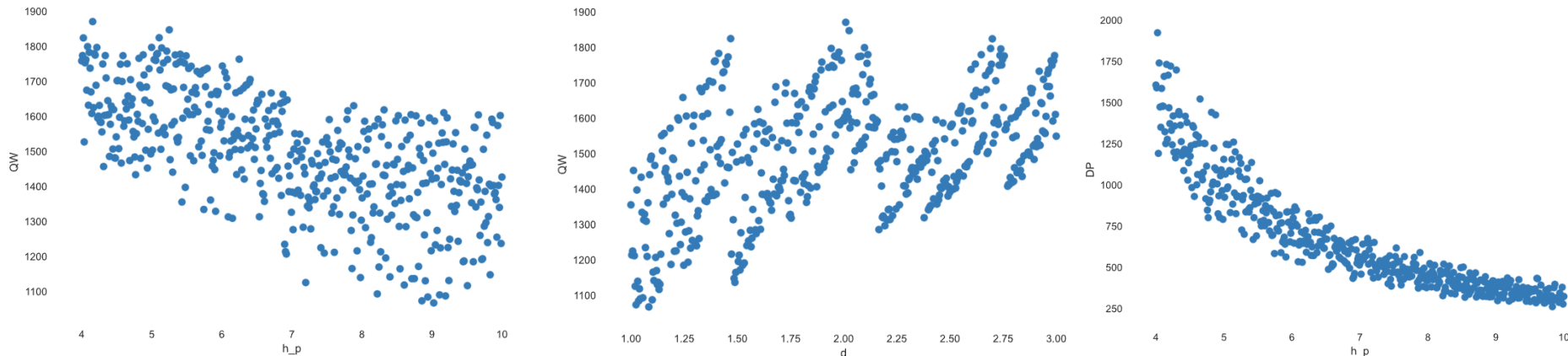
$$X = \{s_{mt}, s_{mq}, d, h_p\} \text{ и } Y = \{QW, DP\}:$$



Предварительный анализ данных

5/15





- Зависимости наблюдается у пар переменных: (h_p, QW) , (h_p, DP) , (d, QW) ;
- Далее необходимо построить матрицу корреляций переменных для вычисления коэффициентов Пирсона.



- Как видно из матрицы корреляций линейная зависимость наблюдается у группы признаков: $(h_p, QW), (h_p, DP), (QW, DP)$. При этом влияние признака h_p уменьшает значения целевых переменных QW, DP . Признак d имеет слабую линейную зависимость с целевой переменной QW ;
- Входные признаки между собой никак не коррелирует и следовательно они либо зависят нелинейно, либо связь полностью отсутствует.

На основе предварительного анализа данных можно сделать вывод о применимости регрессионных моделей. Примеры моделей представлены ниже:

- Линейная регрессия
 - Регрессия LASSO
 - Гребневая регрессия
- Только для данных, которые сильно коррелируют*
- Нейросетевая модель
- Необходимо много данных для качественного обучения*
- Деревья решений/**Gradient Boosting**
- Отлично могут описать любую форму зависимости данных*

Перед обучением модели **Gradient Boosting** необходимо задать следующие гиперпараметры: **iterations = 100**, **learning rate = 0.1**, **Loss function = 'MultiRMSE'**



Оптимизация функционала ошибки на train, test

- В результате рассчитана метрика R^2 для тренировочного и тестового набора данных: $R^2_{train} = 0.975$, $R^2_{test} = 0.924$;
- Определены максимальные и минимальные относительные ошибки предсказаний \widehat{QW} , \widehat{DP} для каждой выходной переменной QW , DP :

Max Relative Error (%)	Train	Test
QW	6.798	11.428
DP	14.618	18.750

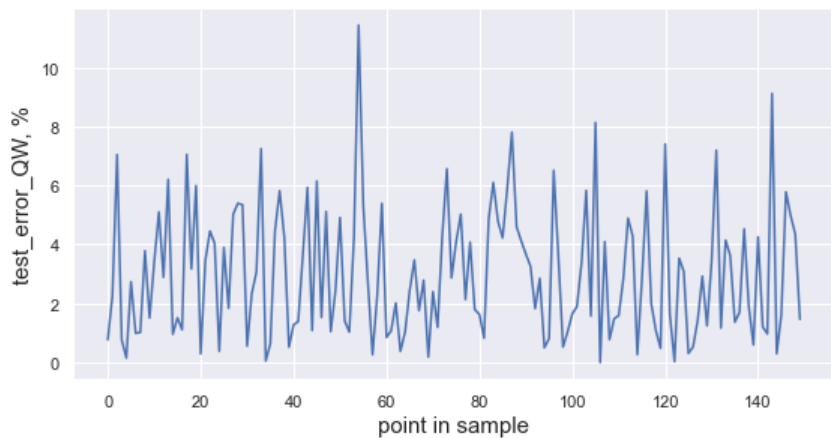
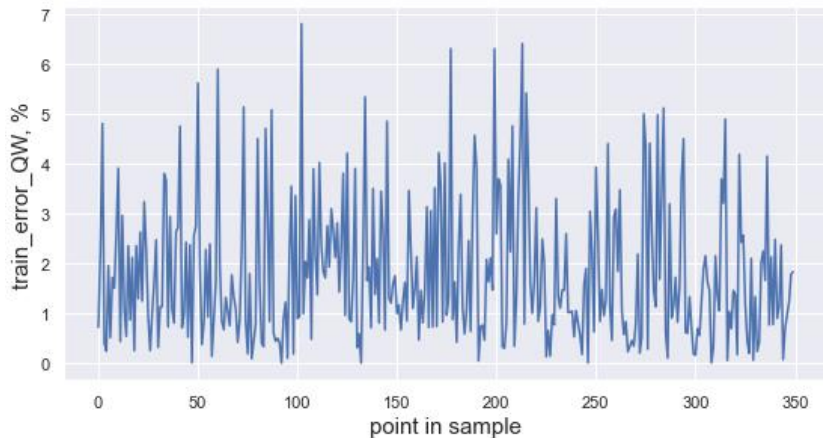
Min Relative Error (%)	Train	Test
QW	0.004	0.001
DP	0.011	0.009

Анализ результатов алгоритма

11/15

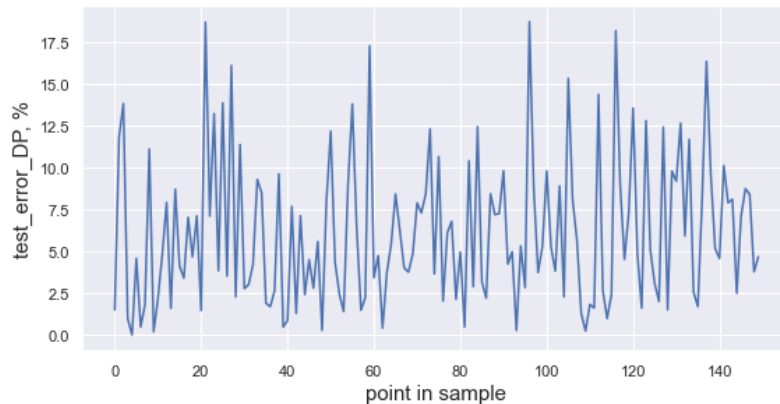
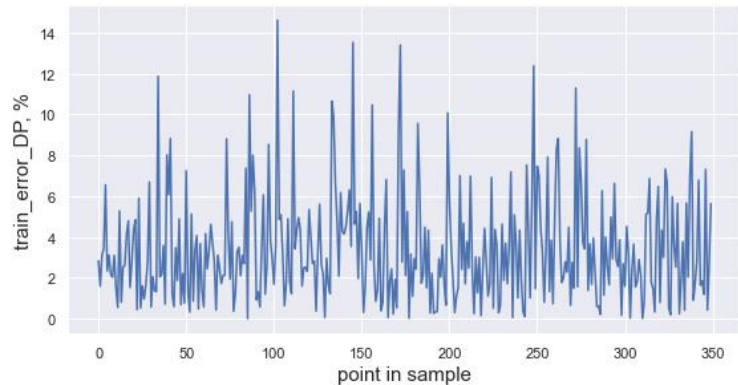
Характер изменения относительных ошибок для величин QW, DP на тренировочных и тестовых данных.

- По оси абсцисс указаны номера объектов в train/test data;
- По оси ординат распределена относительная ошибка (%) для каждого объекта в указанной выборке.



Анализ результатов алгоритма

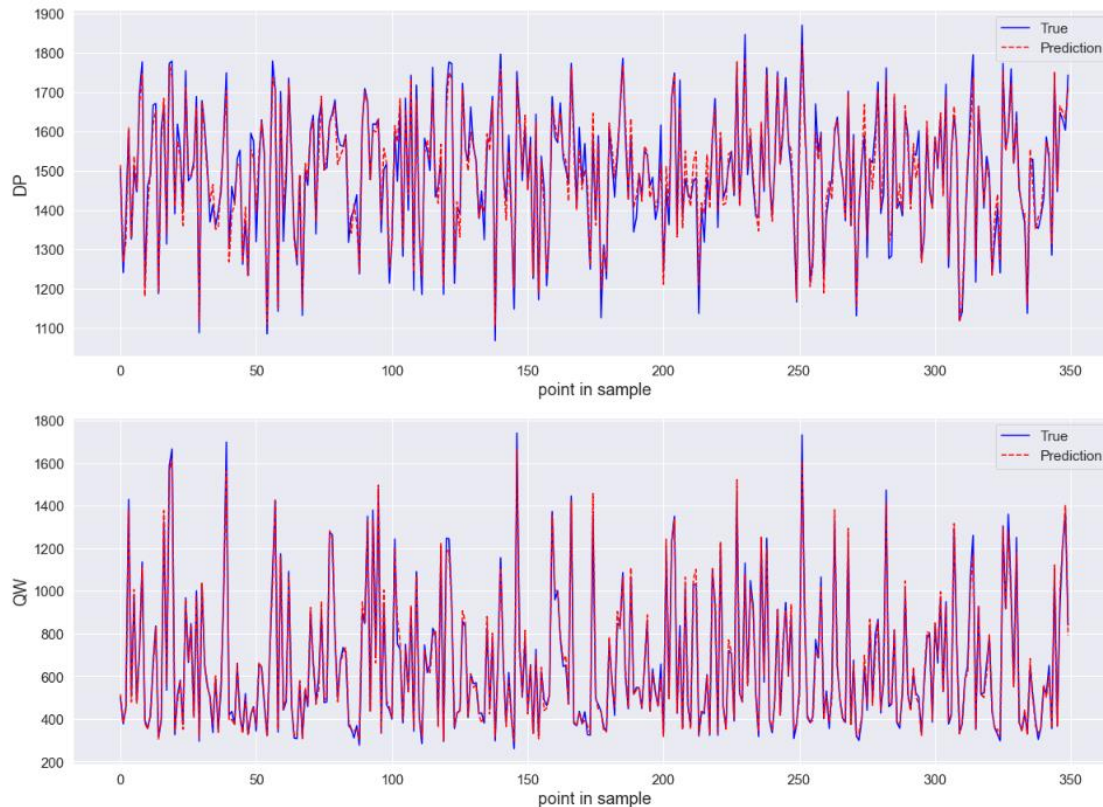
12/15



Анализ результатов алгоритма

13/15

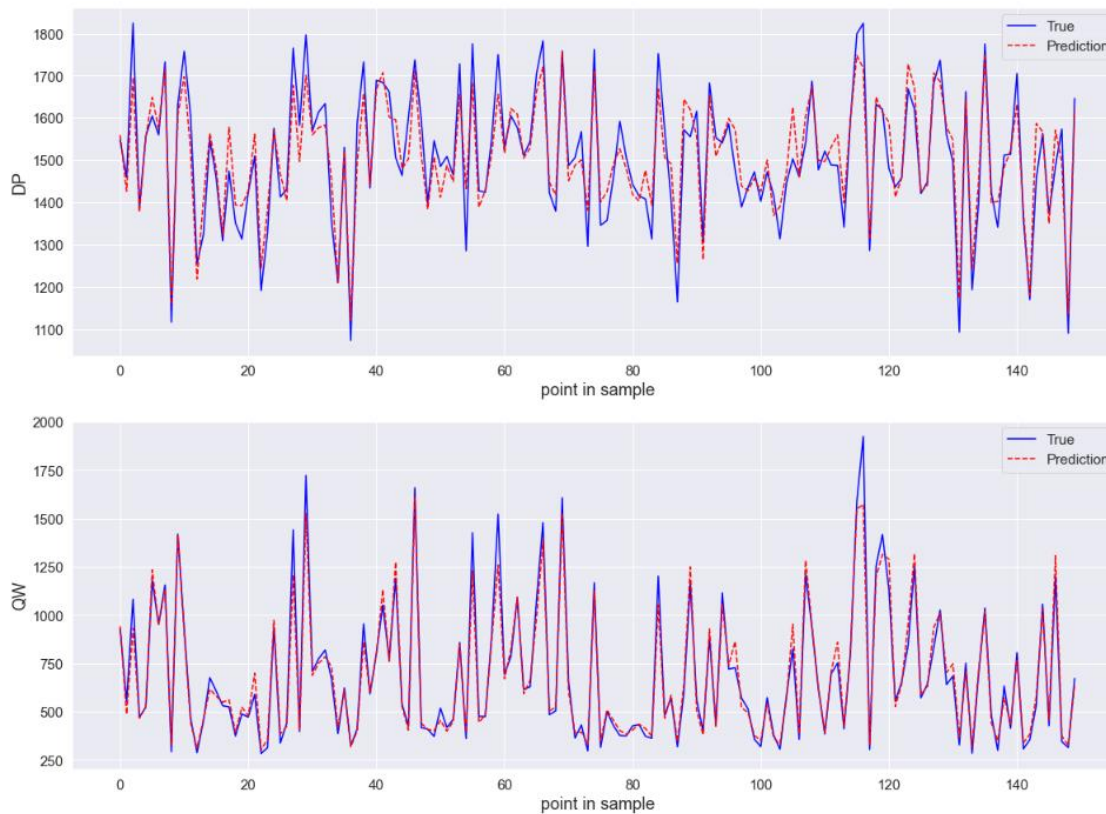
Аппроксимация на тренировочных данных для величин QW , DP



Анализ результатов алгоритма

14/15

Аппроксимация на тестовых данных для величин QW , DP



- Выполнен предварительный анализ данных;
- Построена модель регрессии **Gradient Boosting**;
- Определены и построены относительные ошибки на тестовых и тренировочных данных;
- Рассчитаны метрики $R^2_{train} = 0.975$, $R^2_{test} = 0.924$;
- Построены графики аппроксимации данных на базе **Gradient Boosting**.