

Annot

задача:

Химическая промышленность: Веб-приложение, которое на основе ИИ будет использоваться для отслеживания химического состава природных водных ресурсов. Оно будет анализировать изменения в составе воды и предоставлять рекомендации для минимизации воздействия на окружающую среду и на здоровье человека.

```
Ввод [1]: 1 import os
2 import pandas as pd
3 import numpy as np
4 import pickle
5
6 from typing import List
7
8 import seaborn as sns
9 from matplotlib import pylab as plt
10 %matplotlib inline
11 plt.style.use('bmh')
12 plt.rcParams.update({"font.size": 14})
```

Данные + EDA

Данные найдены на Kaggle (<https://www.kaggle.com/datasets/adityakadiwal/water-potability>)
(<https://www.kaggle.com/datasets/adityakadiwal/water-potability>)

```
Ввод [2]: 1 data = pd.read_csv("../data/water_potability.csv")
2 print(f"Датасет содержит {data.shape[0]} строк и {data.shape[1]} столбцов")
3 init_columns = data.columns
```

Датасет содержит 3276 строк и 10 столбцов

```
Ввод [3]: 1 data.sample()
```

Out[3]:

	ph	Hardness	Solids	Chloramines	Sulfate	Conductivity	Organic_carbon	Trihalomethanes	Turbidity	Potability
2259	7.855906	209.751802	19850.182892	8.038872	300.292954	368.437561	18.525207	57.238874	4.15189	0

Определим пропущенные данные

```
Ввод [4]: 1 for col in data.columns:
2     print(f"{col:30s} :: {str(data[col].isna().sum()):4s} {data[col].isna().sum() / data.shape[0] * 100:.0f}%")
```

ph	:: 491	15%
Hardness	:: 0	0%
Solids	:: 0	0%
Chloramines	:: 0	0%
Sulfate	:: 781	24%
Conductivity	:: 0	0%
Organic_carbon	:: 0	0%
Trihalomethanes	:: 162	5%
Turbidity	:: 0	0%
Potability	:: 0	0%

Так как предполагается работа в продакшене, используем дополнительные фичи для маркировки пропусков в данных. Более того, подобный подход позволит обучить более устойчивую модель.

Ввод [5]:

```
1 for col in data.columns:
2     if col != "Potability":
3         data[f"{col}_nan"] = data[col].isna().astype(int)
4         data.loc[data[col].isna(), col] = 0.0
5 data.sample()
```

Out[5]:

	ph	Hardness	Solids	Chloramines	Sulfate	Conductivity	Organic_carbon	Trihalomethanes	Turbidity	Potability	ph_n
1020	5.21899	164.790992	27007.886906	6.592517	0.0	422.223802	14.070879	73.92032	4.383943	0	

Ввод [6]:

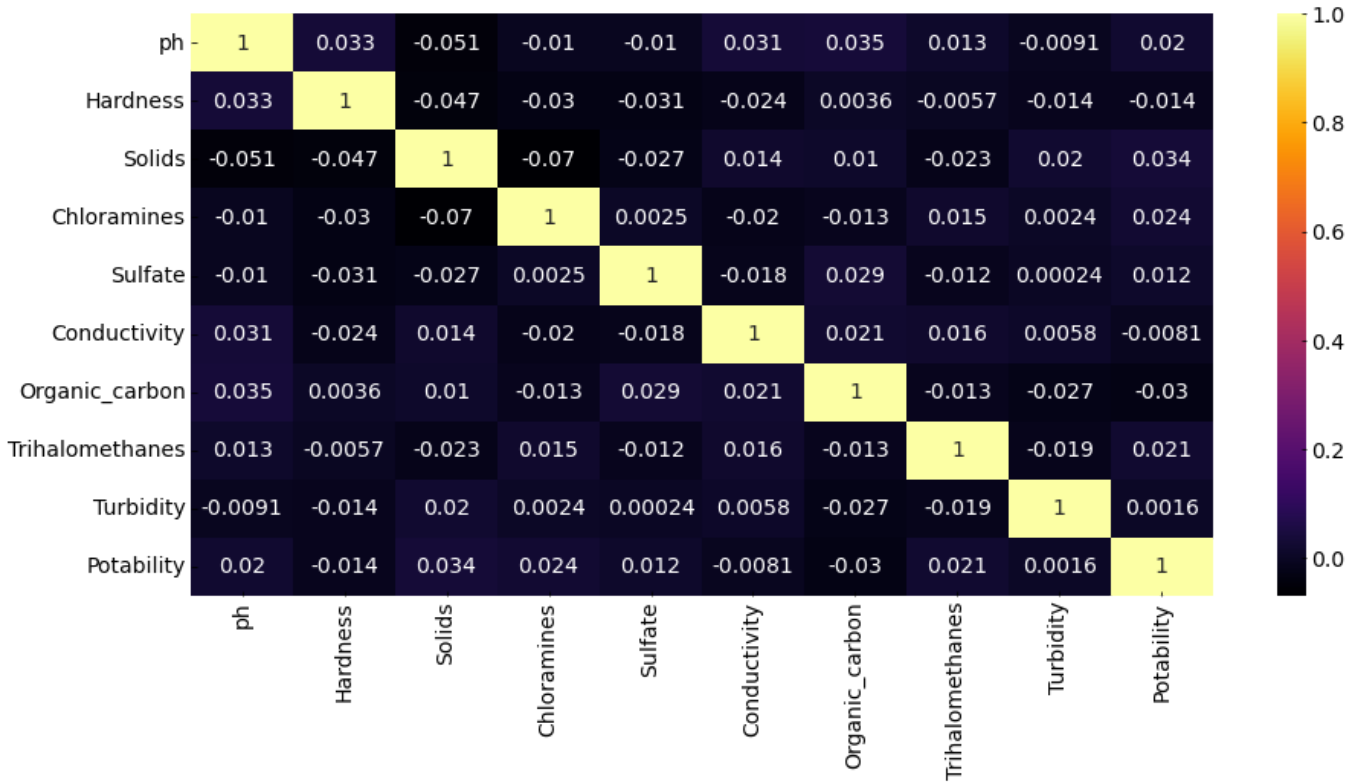
```
1 print(f"Пропусков в обогащенных данных: {data.isna().sum().sum()}")
```

Пропусков в обогащенных данных: 0

Посмотрим на корреляции между целевой переменной и исходными колонками

Ввод [7]:

```
1 corr = data[init_columns].corr()
2 plt.figure(figsize=(15, 7))
3 sns.heatmap(corr, cmap="inferno", annot=True)
4 plt.show()
```

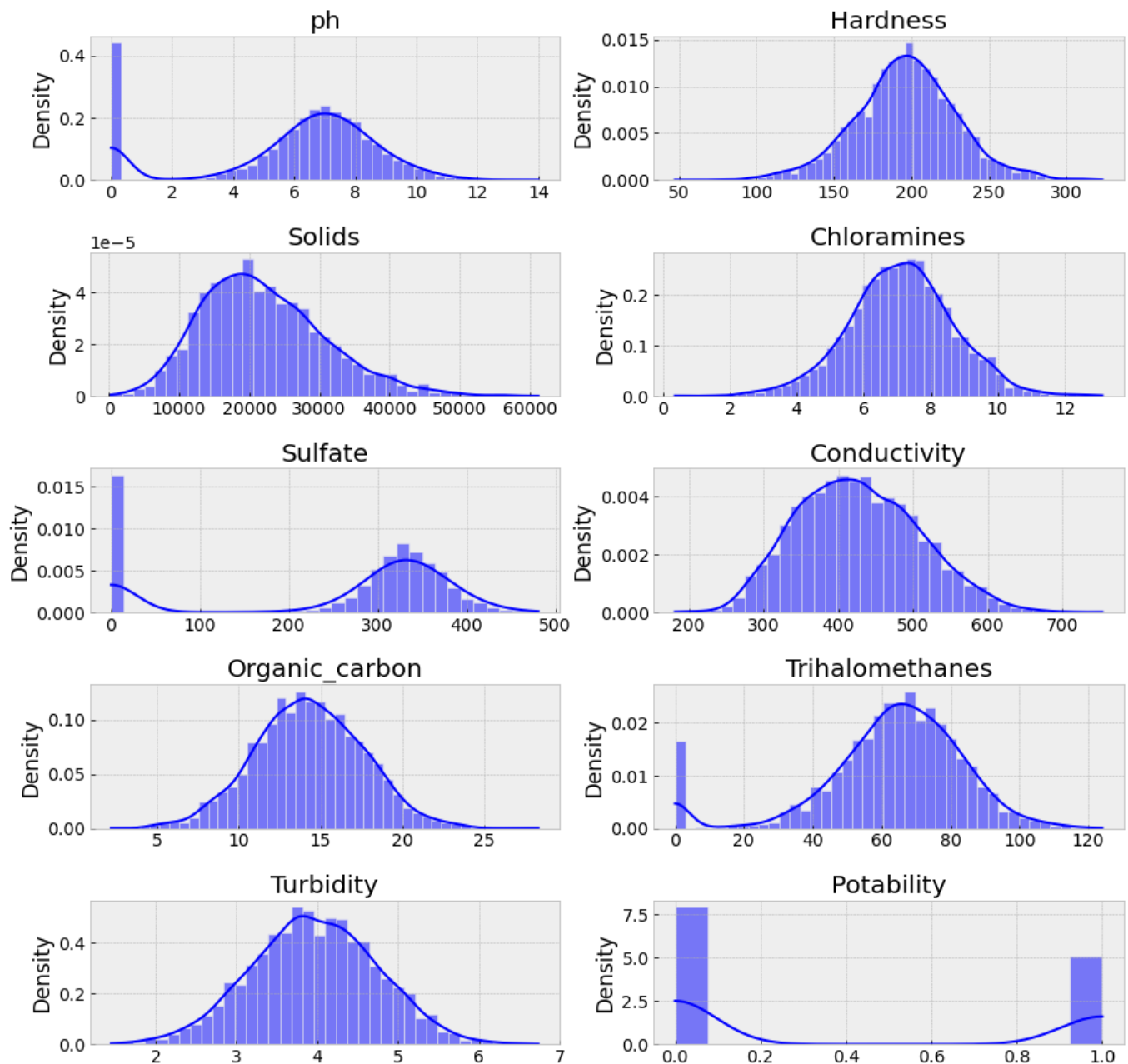


Данные достаточно чистые, в плане внутренних корреляций. Нет сильно коррелированных фичей.

Визуализируем распределения значений

Ввод [8]:

```
1 fig, ax = plt.subplots(5, 2, figsize=(15, 15))
2
3 for idx, col in enumerate(init_columns):
4     sns.histplot(data[col], ax=ax[idx // 2, idx % 2], kde=True, stat='density')
5     ax[idx // 2, idx % 2].set_title(col)
6     ax[idx // 2, idx % 2].set_xlabel("")
7 plt.subplots_adjust(hspace=0.5)
```



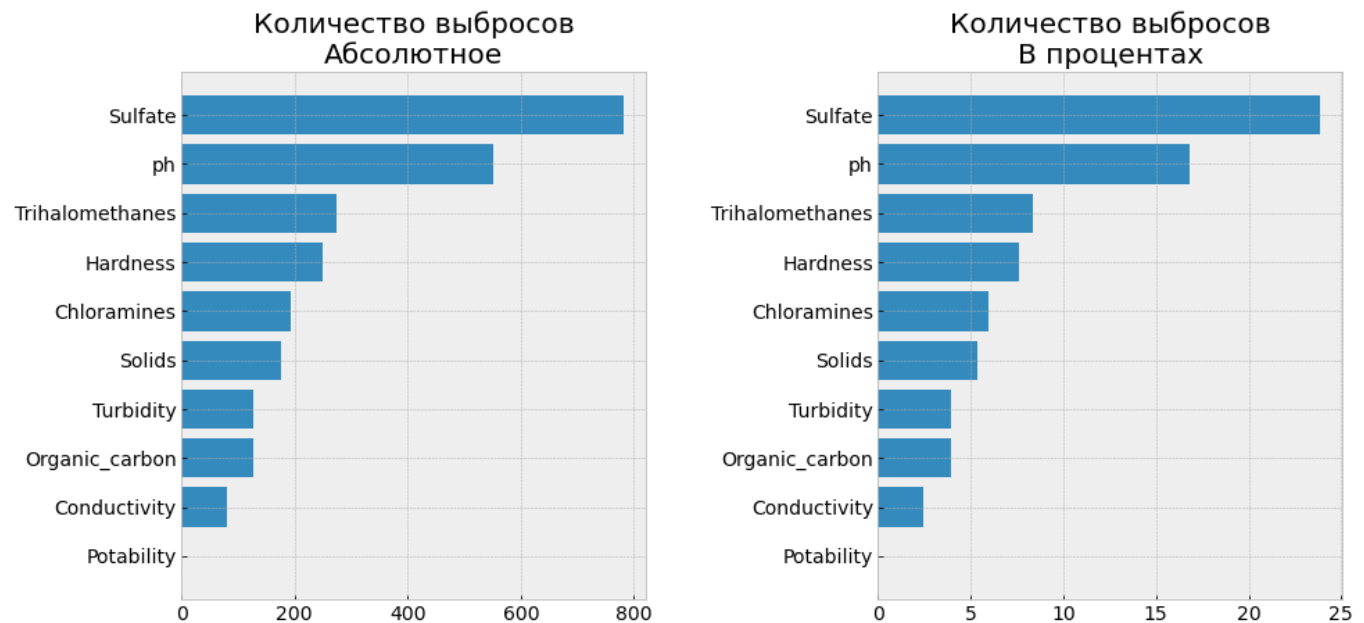
Данные распределены близко к нормальному распределению. Только Solids немного смещен. Проверим на выбросы. Нули - замещенные NaN

Ввод [9]:

```
1 def get_outliners(x: np.ndarray) -> int:
2     """Реализация выявления выбросов, которые не укладываются в 1.5 * IQR
3
4     Args:
5         x (np.ndarray): входной набор данных (вектор)
6
7     Returns:
8         (int): количество выбросов
9     """
10    iqr = 1.5 * (np.quantile(x, 0.75) - np.quantile(x, 0.25))
11    low_bond = x.median() - iqr
12    upper_bound = x.median() + iqr
13    return np.where(x > upper_bound)[0].size + np.where(x < low_bond)[0].size
```

Ввод [10]:

```
1 _temp_series = pd.Series({col: get_outliners(data[col]) for col in init_columns})
2 _temp_series = _temp_series.sort_values()
3
4 fig, ax = plt.subplots(1, 2, figsize=(15, 7))
5 ax[0].barh(_temp_series.index, _temp_series.values)
6 ax[0].set_title("Количество выбросов\nАбсолютное")
7 ax[1].barh(_temp_series.index, _temp_series.values / len(data) * 100)
8 ax[1].set_title("Количество выбросов\nВ процентах")
9 plt.subplots_adjust(wspace=0.5)
10 plt.show()
```



Так как выбросов либо нет, либо они не единичные (от 2.5 % значений фиичи) можно заключить, что это не выбросы, а **особенности данных**

Train

Будем учить XGBoost [на GPU, чтоб побыстрее, в инференсе будет CPU]. Параметры найдем перебором.

План обучения:

- baseline на чистом датасете
- повышение устойчивости путем генерации пропусков в данных

Ввод [11]:

```
1 import time
2 import pickle
3
4 from xgboost import XGBClassifier
5
6 from sklearn.model_selection import GridSearchCV
7 from sklearn.model_selection import train_test_split
8 from sklearn.metrics import f1_score, recall_score, roc_auc_score, accuracy_score, mutual_info_score
```

```

Ввод [12]: 1 def get_metrics(classifier, x: pd.DataFrame, y: pd.Series, name: str, training_time: str, num_classes: int)
2           """Функция для оценки качества обучения
3           Args:
4               classifier объект классификатора с методами predict и predict_proba
5               x (pd.DataFrame) признаки
6               y (pd.Series) лейблы
7               name (str) имя классификатора
8               training_time (str) время обучения (включая тюнинг) в формате чч:мм:сс
9               num_classes (int) число классов
10          Return:
11              (dict) структуру строки для датасета с результатами
12          """
13          preds = classifier.predict(x)
14          preds_proba = classifier.predict_proba(x)
15          return {
16              "model": [name],
17              "time": [training_time],
18              "accuracy": [accuracy_score(y, preds)],
19              "recall": [recall_score(y, preds)],
20              "roc_auc": [roc_auc_score(np.eye(num_classes)[y], preds_proba)],
21              "f1_score": [f1_score(y, preds)],
22              "mutal_score": [mutual_info_score(y, preds)],
23          }

```

Обучение baseline

```

Ввод [13]: 1 X_train, X_test, y_train, y_test = train_test_split(
2           data.drop(columns=['Potability']), data['Potability'], random_state=6875, train_size=0.8
3           )

```

```

Ввод [14]: 1 model = XGBClassifier()
2           parameters = {
3               "n_estimators": [10, 25, 50, 100, 150, 200, 250],
4               "max_depth": [5, 10, 15, 20],
5               "learning_rate": [1e-2, 1e-3, 1e-5],
6               "objective": ["binary:logistic"],
7               "eval_metric": ["auc"],
8               "tree_method": ["gpu_hist"],
9           }

```

```

Ввод [15]: 1 clf = GridSearchCV(model, parameters, cv=2, verbose=False)
2           start = time.time()
3           clf.fit(X_train, y_train)
4           train_time = time.time() - start

```

```

Ввод [16]: 1 metrics_frame = pd.DataFrame(get_metrics(clf.best_estimator_, X_test, y_test, "search", train_time, 2))
2           metrics_frame

```

```

Out[16]:
   model  time  accuracy  recall  roc_auc  f1_score  mutal_score
0  search  85.754553  0.682927  0.369919  0.682144  0.466667  0.038577

```

```

Ввод [20]: 1 pickle.dump(clf.best_estimator_, open("../models/baseline.bin", "wb"))

```

На тестовых данных ROC AUC в районе 0.68, что в целом, соответствует результатам из открытых источников (расшаренный код)

Обогащение данных

Для обогащения данных в случайных местах пропускаем данные и ставим колонку в соответствующее состояние. Увеличим датасет в 10 раз (по количеству фичей + исходный)

Ввод [21]:

```
1 X_train_large = X_train.copy()
2 y_train_large = y_train.copy()
3 for col in init_columns:
4     if col != 'Potability':
5         _temp_frame = X_train.copy()
6         _temp_frame[col] = 0.0
7         _temp_frame[f'{col}_nan'] = 1
8     X_train_large = pd.concat([X_train_large, _temp_frame], ignore_index=True)
9     y_train_large = pd.concat([y_train_large, y_train.copy()], ignore_index=True)
10 print(f"Размерность нового датасета: {len(X_train_large)}")
```

Размерность нового датасета: 28820

Ввод [22]:

```
1 clf = GridSearchCV(model, parameters, cv=2, verbose=False)
2 start = time.time()
3 clf.fit(X_train_large, y_train_large)
4 train_time = time.time() - start
```

Ввод [23]:

```
1 metrics_frame = pd.concat([
2     metrics_frame, pd.DataFrame(get_metrics(clf.best_estimator_, X_test, y_test, "large", train_time, 2))
3 ], ignore_index=True)
4 pickle.dump(clf.best_estimator_, open("../models/robust.bin", "wb"))
```

Ввод [24]:

```
1 metrics_frame
```

Out[24]:

	model	time	accuracy	recall	roc_auc	f1_score	mutal_score
0	search	85.754553	0.682927	0.369919	0.682144	0.466667	0.038577
1	large	237.992343	0.660061	0.398374	0.672606	0.467780	0.027405

Как видим качество не сильно упало, однако повысилась устойчивость - теперь модель может работать с неполными данными



Сформируем базу данных для Web-аpi

Ввод [11]:

```
1 import sqlite3
2 con = sqlite3.connect("../main_base.db")
```

Ввод [12]:

```
1 data.to_sql('train_data', con)
```

Out[12]: 3276

Ввод []:

```
1
```