

МИНИСТЕРСТВО НАУКИ И ВЫСШЕГО ОБРАЗОВАНИЯ РОССИЙСКОЙ  
ФЕДЕРАЦИИ

Федеральное государственное автономное  
образовательное учреждение высшего образования  
«Самарский национальный исследовательский университет  
имени академика С.П. Королева»  
(Самарский университет)

Институт информатики, математики и электроники  
Факультет информатики  
Кафедра технической кибернетики

**Отчет по лабораторной работе № 3  
по курсу «Инженерия данных»**

Тема: «Airflow и MLflow - логгирование экспериментов и версионирование  
моделей»

Выполнил Самигуллин Равиль  
Группа 6233 – 010402D  
Преподаватель Парингер Р.А

**Самара**

## Содержание

Содержание.....	2
Техническое задание.....	3
Ход выполнения .....	4
Задание 1: Пайплайн для инференса данных. ....	4
Задание 2: Пайплайн для хостинга лучшей модели .....	13
ЗАКЛЮЧЕНИЕ .....	20

## Техническое задание

### **Airflow и MLflow - логгирование экспериментов и версионирование моделей**

В рамках данной лабораторной работы предлагается построить два пайплайна:

1. Пайплайн, который обучает любой классификатор из sklearn по заданному набору параметров.
2. Пайплайн, который выбирает лучшую модель из обученных и производит её хостинг.

Для построения такого пайплайна воспользуемся следующими инструментами:

- Apache Airflow
- MLflow

#### Задание на лабораторную работу

##### Пайплайн для обучения классификаторов

Построенный пайплайн будет выполнять следующие действия поочередно:

1. Производить мониторинг целевой папки на предмет появления новых конфигурационных файлов классификаторов (в форматах .json или .yaml).
2. Обучать классификатор в соответствии с полученными параметрами.
3. Производить логгирование параметров модели в **MLflow**.
4. Производить логгирование процесса обучения **MLflow**.
5. Производить тестирование модели и сохранять его результаты в **MLflow**.
6. Сохранять обученный классификатор в model registry **MLflow**.

##### Пайплайн для хостинга лучшей модели

Построенный пайплайн будет выполнять следующие действия поочередно:

1. В соответствии с таймером производит валидацию новых моделей из model registry.
2. Модель с лучшим показателем метрики переводится на stage: Production
3. (Опционально) произвести хостинг лучшей модели

## Ход выполнения

### Задание 1: Пайплайн для инференса данных.

#### 1. Выбор данных и классификаторов для обучения:

Перед обучением и созданием выбрал несколько классификаторов из sklearn которые будут обрабатывать наши данные, для обработки данных так же нужно выбрать датасет на котором планирую обучать наши классификаторы, а так же я решил делать предобработку данных перед обучением, а так как я всё равно буду записывать аргументы и настройки классификаторов в json, то решил туда же добавлять всё то что хочу сделать с данными и то какие данные мы будем обучать, мне показалось это правильным так как, наш пайплайн сможет обрабатывать любые данные и обучать любые модели (из sklearn разумеется).

Я выбрал следующие модели:

- SGDClassifier
- RandomForestClassifier
- SVC
- KNeighborsClassifier
- DecisionTreeClassifier
- LogisticRegression
- Perceptron

В качестве датасета я выбрал titanic классический датасет который так же требует предобработки перед обучением, так как имеет множество не закодированных значений представленных в виде категорий, а так же значений которые не имеют смысла для классификации, как например номер билета или имя пассажира.

#### 2. Создание json файла с настройкой конфигурации.

```

data > {} conf.json > [ ] configs > { } 0 > { } args
1  {
2      "dataset_info": {
3          "name": "titanic",
4          "data_path": "/opt/airflow/data/lr3",
5          "drop_columns": ["PassengerId", "Name", "Ticket", "Cabin"],
6          "categorical_columns": ["Sex", "Embarked"]
7      },
8      "configs": [
9          {
10             "module": "sklearn.linear_model",
11             "classifier": "SGDClassifier",
12             "args": {
13                 "loss": "log_loss",
14                 "max_iter": 200
15             }
16         },
17         {
18             "module": "sklearn.ensemble",
19             "classifier": "RandomForestClassifier",
20             "args": {
21                 "n_estimators": 100,
22                 "max_depth": 10,
23                 "random_state": 42
24             }
25         },
26         {
27             "module": "sklearn.svm",
28             "classifier": "SVC",
29             "args": {
30                 "kernel": "rbf",
31                 "C": 1.0,
32                 "gamma": "scale"
33             }
34         },
35         {
36             "module": "sklearn.neighbors",
37             "classifier": "KNeighborsClassifier",
38             "args": {
39                 "n_neighbors": 5,
40                 "weights": "uniform"
41             }
42         }
43     ]
44 }

```

В json я записал данные по dataset расположение файлов csv, имена столбцов необходимые для кодировки, имена столбцов для удаления, после идут настройки классификаторов, модуль для импорта библиотеки, имя классификатора, и аргументы которые мы используем для обучения классификации.

### 3. Разработка DAG

```

dag_lr3_task_1.py > ...
1  from airflow import DAG
2  from datetime import datetime
3  from airflow.providers.docker.operators.docker import DockerOperator
4  from docker.types import Mount
5  from airflow.operators.bash_operator import BashOperator
6  from airflow.sensors.filesystem import FileSensor
7  import os
8
9  os.environ["AWS_ACCESS_KEY_ID"] = "minio"
10 os.environ["AWS_SECRET_ACCESS_KEY"] = "minio123"
11 os.environ["MLFLOW_S3_ENDPOINT_URL"] = "http://minio:9000"
12
13 default_args = {
14     'owner': 'darti',
15     'start_date': datetime(2024, 1, 1),
16     'retries': 1,
17 }
18
19 dag = DAG(
20     'lr3_learning_models',
21     default_args = default_args,
22     description='dag for learn sclearn models and data',
23     schedule_interval=None,
24 )
25
26 wait_for_new_file = FileSensor(
27     task_id='wait_for_new_file',
28     poke_interval=10, # Interval to check for new files (in seconds)
29     filepath= '/opt/airflow/data/lr3/data/*.json',
30     fs_conn_id='fs_connection_default', # Target folder to monitor
31     dag=dag,
32 )
33
34 train_model = BashOperator(
35     task_id='train_model',
36     bash_command='/opt/***/data/lr3/train.py {{ params.config_file }}',
37     params={'config_file': '/opt/airflow/data/lr3/data/conf.json'},
38     dag=dag,
39 )
40
41 wait_for_new_file >> train_model

```

В данном dag, мы сначала ждём файл wait for new file в папку мы ждём json.

После мы через bash operator выполняем скрипт train.py передавая ему в качестве аргумента config\_file.

## 4. Подготовка данных

```
predobr.py > ...
1  #!/usr/bin/env python3
2  import mlflow
3  import mlflow.sklearn
4  import pandas as pd
5  import numpy as np
6  from sklearn.model_selection import train_test_split
7  from sklearn.metrics import f1_score, accuracy_score, precision_score, recall_score
8  from sklearn.preprocessing import LabelEncoder
9
10 def get_data(data_path: str, labels: str):
11     """Загрузка данных из CSV файлов."""
12     df_features = pd.read_csv(f'{data_path}/{labels}_features.csv')
13     df_target = pd.read_csv(f'{data_path}/targets_{labels}.csv')
14     return df_features, df_target
15
16
17 def preprocess_data(df_features, df_labels, drop_columns, categorical_columns):
18     # Удаляем ненужные столбцы
19     df_features.drop(columns=drop_columns, inplace=True)
20
21     # Обрабатываем категориальные столбцы, преобразовывая их в числовые
22     for column in categorical_columns:
23         df_features[column] = pd.Categorical(df_features[column]).codes
24
25     # Заполняем пропуски в числовых столбцах медианой
26     df_features.fillna(df_features.median(numeric_only=True), inplace=True)
27
28     return df_features, df_labels
29
30
31
32 def split_data(x_data, y_data, test_size=0.2, random_state = 42):
33     """Разделение данных на тренировочные и тестовые выборки."""
34     x_train, x_test, y_train, y_test = train_test_split(x_data, y_data, test_size=test_size, random_state=random_state)
35     return x_train, x_test, y_train, y_test
36
37
38 def log_metrics(config, y_test, y_pred):
39     """Логирование метрик модели в MLflow."""
40     mlflow.log_params(config)
41     mlflow.log_metrics({
42         "f1_score": f1_score(y_test, y_pred, average='weighted'),
43         "accuracy": accuracy_score(y_test, y_pred),
44         "precision": precision_score(y_test, y_pred, average='weighted'),
45         "recall": recall_score(y_test, y_pred, average='weighted')
46     })
47
```

Используя функции написанные в файле predobr.py

Мы удаляем ненужные столбцы, обрабатываем категориальные переменные (преобразует их в числовые коды) и заполняем пропуски в числовых столбцах медианой.

## 5. Обучение модели

Обучение модели а так же вызов всех вспомогательных функций реализован в файле train.py

```
train.py > ...
1  #!/usr/bin/env python3
2  import json
3  import importlib
4  import random as rnd
5  import pandas as pd
6  import numpy as np
7  import uuid
8  import mlflow
9  import mlflow.sklearn
10 import predobr as pred
11 from mlflow.models import infer_signature
12 import sys
13 from sklearn.preprocessing import LabelEncoder
14 from sklearn.metrics import f1_score, accuracy_score, precision_score, recall_score
15
16
17 # URL для подключения к MLflow Tracking Server
18 tracking_url = 'http://mlflow_server:5000'
19
20 # Получаем путь к конфигурационному файлу
21 config_file_path = sys.argv[1]
22
23 # Загрузка конфигурации из JSON файла
24 with open(config_file_path, 'r') as config_file:
25     config_data = json.load(config_file)
26
27 # Извлекаем информацию о наборе данных из JSON
28 dataset_name = config_data['dataset_info']['name']
29 data_path = config_data['dataset_info']['data_path']
30 drop_columns = config_data['dataset_info']['drop_columns']
31 categorical_columns = config_data['dataset_info']['categorical_columns']
32
33 # Загрузка данных
34 x_data, y_data = pred.get_data(data_path, dataset_name)
35
36 # Применяем предобработку к данным
37 # Применяем предобработку к данным
38 x_data, y_data, encoders = pred.preprocess_data(pd.DataFrame(x_data), pd.DataFrame(y_data))
39
40 # Разделение данных на тренировочные и тестовые выборки
41 x_train, x_test, y_train, y_test = pred.split_data(x_data, y_data)
42 x_val, x_temp, y_val, y_temp = pred.split_data(x_data, y_data, test_size=0.5, random_state
43
44 # Сохранение валидационных данных в CSV это для задания 2
45 x_val.to_csv(f'{data_path}/x_val.csv', index=False)
46 y_val.to_csv(f'{data_path}/y_val.csv', index=False)
47
48 # Удаление временных данных из памяти
```



```

train.py > ...
47
48 # Установка URI для подключения к MLflow Tracking Server
49 mlflow.set_tracking_uri(tracking_url)
50
51 # Генерация уникального ID эксперимента
52 id_current_experiment = str(uuid.uuid4())
53 exp_id = mlflow.create_experiment(id_current_experiment)
54 mlflow.set_experiment(exp_id)
55
56 # Сохранение ID эксперимента в файл так же надо для 2 задания
57 with open(f'{data_path}/experiment_id.txt', 'w') as f:
58     f.write(id_current_experiment)
59
60 # Проходим по каждой конфигурации в JSON
61 for i, config in enumerate(config_data['configs']):
62     mlflow.start_run(run_name=config['classifier'], experiment_id=exp_id)
63
64     # Загрузка модуля и классификатора
65     module = importlib.import_module(config['module'])
66     classifier_class = getattr(module, config['classifier'])
67     model = classifier_class(**config['args'])
68
69     # Проверка на наличие метода partial_fit
70     if hasattr(model, 'partial_fit'):
71         max_iter = config['args'].get('max_iter', 50)
72         for iter in range(max_iter):
73             model.partial_fit(x_train, y_train, classes=np.unique(y_train))
74             y_pred = model.predict(x_test)
75             mlflow.log_metrics({
76                 f"f1_score": f1_score(y_test, y_pred, average='weighted'),
77                 f"accuracy": accuracy_score(y_test, y_pred),
78                 f"precision": precision_score(y_test, y_pred, average='weighted'),
79                 f"recall": recall_score(y_test, y_pred, average='weighted')
80             })
81     else:
82         model.fit(x_train, y_train) # Стандартное обучение
83
84     # Тестирование модели на тестовых данных
85     y_pred_test = model.predict(x_test)
86     # Логирование метрик на тестовых данных
87     mlflow.log_metrics({
88         "f1_score_test": f1_score(y_test, y_pred_test, average='weighted'),
89         "accuracy_test": accuracy_score(y_test, y_pred_test),
90         "precision_test": precision_score(y_test, y_pred_test, average='weighted'),
91         "recall_test": recall_score(y_test, y_pred_test, average='weighted')
92     })
93
94     # Сохранение информации о модели в CSV
95     model_info = mlflow.sklearn.log_model(sk_model=model, artifact_path=f'{config["module"]}')
96
97     # Регистрация модели в Model Registry
98     model_uri = model_info.model_uri
99     registered_model = mlflow.register_model(model_uri=model_uri, name=config['classifier'])
100
101     df = pd.DataFrame({"name": config['classifier'], "uri": model_info.model_uri}, index=[i])
102     df.to_csv(f'{data_path}/models.csv', mode='a', header=False) # Добавляем заголовки только в первый раз
103
104     mlflow.end_run()

```

В основном скрипте мы:

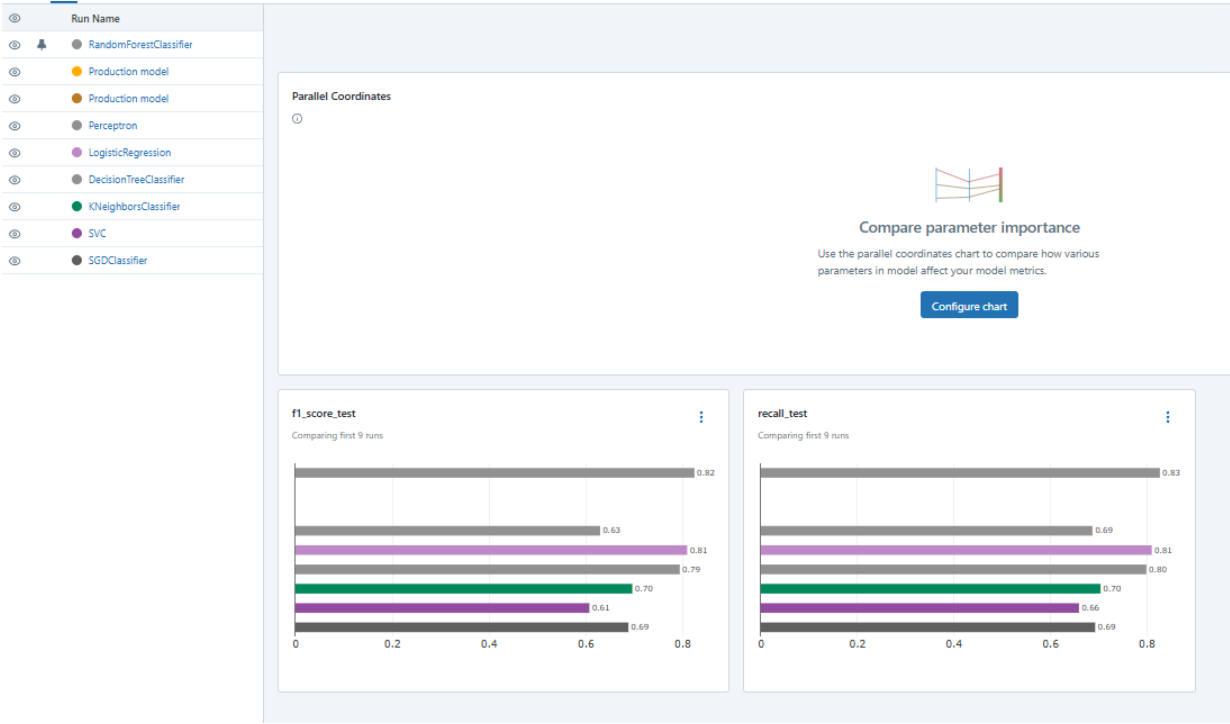
- Загружаем путь к конфигурационному файлу, читает его в формате JSON и извлекает необходимую информацию о наборе данных и настройках.
- Загружаем данные с помощью функции **get\_data**.

- Применяем предварительную обработку и разделяет данные на тренировочные и тестовые выборки.
- Сохраняем валидационные наборы данных в CSV. (Это так же для второго задания)
- Устанавливаем URI для подключения к MLflow Tracking Server.
- Создаем новый эксперимент в MLflow и логируем его ID. (Сейчас, я думаю, что лучше сохранять не в txt, а в каком-нибудь csv, чтобы хранить не только текущий эксперимент, но и все предыдущие)
- Проходим по каждой конфигурации в JSON и создаем классификатор, загружая соответствующий модуль и класс.
- Обучаем модель, если в ней реализован метод **partial\_fit** то обучаем по steps и сохраняем метрики каждого шага, я это сделал что бы можно было отследить процесс обучения у моделей у которых это возможно, или используем стандартное обучение.
- Логируем метрики на тестовых данных.
- Сохраняем информацию о модели в систему MLflow и регистрируем модель.
- Сохраняем информацию о модели в CSV файл. (Это нужно для второго шага)

## 6. Результат выполнения pipeline



Сначала сохранял метрики не правильно и каждая итерация была со своим графиком, после сохранил их с одним именем и получил хорошие графики.



Тут можно посмотреть точность моделей которые мы проверили на тестовых данных.

### SGDClassifier

Run ID: 7ef538f854ea4840b7425ea3d0c97ffd

Status: FINISHED

> Description [Edit](#)

> Datasets

> Parameters

▼ Metrics (804)

Name	Value
accuracy_iter_0 <a href="#">📄</a>	0.732
accuracy_iter_1 <a href="#">📄</a>	0.654
accuracy_iter_10 <a href="#">📄</a>	0.732
accuracy_iter_100 <a href="#">📄</a>	0.754
accuracy_iter_101 <a href="#">📄</a>	0.76
accuracy_iter_102 <a href="#">📄</a>	0.771
accuracy_iter_103 <a href="#">📄</a>	0.642
accuracy_iter_104 <a href="#">📄</a>	0.76
accuracy_iter_105 <a href="#">📄</a>	0.76
accuracy_iter_106 <a href="#">📄</a>	0.821
accuracy_iter_107 <a href="#">📄</a>	0.788
accuracy_iter_108 <a href="#">📄</a>	0.821

## f1\_score

Completed Runs <sup>?</sup>

1/1

Points: ☐Line Smoothness <sup>?</sup> 78

X-axis:

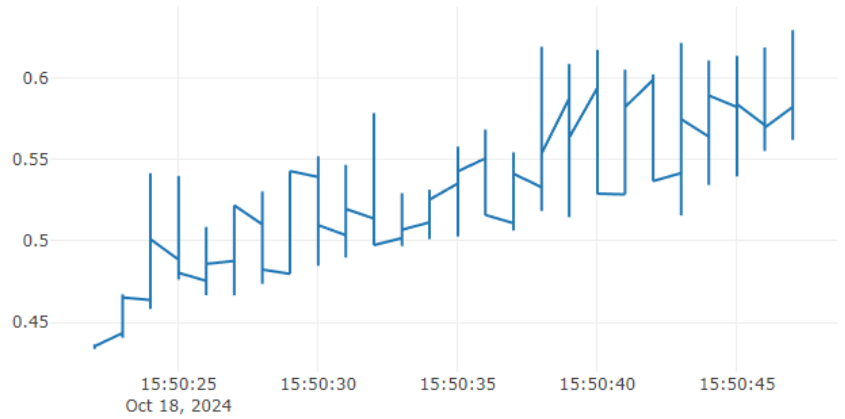
☐ Step☒ Time (Wall)☐ Time (Relative)

Y-axis:

f1\_score X

Y-axis Log Scale: ☐

Download CSV



Metric	Latest	Min	Max
f1_score	0.568 (step=0)	0.434 (step=0)	0.711 (step=0)

## f1\_score

Completed Runs <sup>?</sup>

1/1

Points: ☐Line Smoothness <sup>?</sup> 78

X-axis:

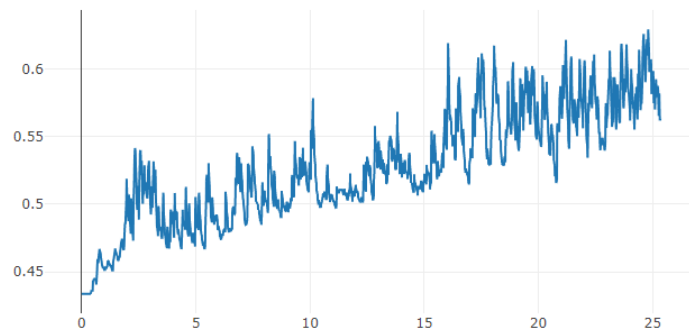
☐ Step☐ Time (Wall)☒ Time (Relative)

Y-axis:

f1\_score X

Y-axis Log Scale: ☐

Download CSV



Metric	Latest	Min	Max
f1_score	0.568 (step=0)	0.434 (step=0)	0.711 (step=0)

И теперь можно посмотреть, как выглядит график обучения процептона или логистической регрессии

## Задание 2: Пайплайн для хостинга лучшей модели

Хостинг лучшей модели.

### 1. Разработка dag

```
dag_lr3_task_2.py > ...
1  from datetime import datetime
2  from airflow import DAG
3  from airflow.operators.bash_operator import BashOperator
4  import os
5
6  os.environ["AWS_ACCESS_KEY_ID"] = "minio"
7  os.environ["AWS_SECRET_ACCESS_KEY"] = "minio123"
8  os.environ["MLFLOW_S3_ENDPOINT_URL"] = "http://minio:9000"
9
10 default_args = {
11     'owner': 'darti',
12     'start_date': datetime(2024, 10, 20),
13     'retries': 1,
14 }
15
16 dag = DAG(
17     'validate_and_promote_model',
18     default_args=default_args,
19     description='DAG for validating and promoting the best model from MLflow',
20     schedule_interval='@daily',
21 )
22
23 # Задача для валидации модели
24 validate_model = BashOperator(
25     task_id="validate_model",
26     bash_command="python /opt/airflow/data/lr3/validate.py",
27     dag=dag
28 )
29
30 # Задача для продвижения модели
31 promote_best_model_and_host = BashOperator(
32     task_id="promote_best_model_and_host",
33     bash_command="python /opt/airflow/data/lr3/promote_model.py",
34     dag=dag
35 )
36
37 validate_model >> promote_best_model_and_host # Зависимость между задачами
38
39
```

Сначала производим валидацию моделей после, переводим лучшую модель в promote и пытаемся захостить её.

### 2. Создания скрипта validate

```

1  #!/usr/bin/env python3
2  import pandas as pd
3  import numpy as np
4  from sklearn.metrics import accuracy_score
5  import mlflow
6  import mlflow.sklearn
7
8  tracking_url = 'http://mlflow_server:5000'
9  data_path = '/opt/airflow/data/lr3'
10
11 # Чтение ID текущего эксперимента
12 with open(f'{data_path}/experiment_id.txt', 'r') as f:
13     id_current_experiment = f.read()
14
15 mlflow.set_tracking_uri(tracking_url)
16 mlflow.set_experiment(id_current_experiment)
17
18 # Загрузка валидационных данных
19 x_val = pd.read_csv(f'{data_path}/x_val.csv')
20 y_val = pd.read_csv(f'{data_path}/y_val.csv')
21
22 list_models_for_validate = {}
23
24 with mlflow.start_run(run_name="Production model") as start_run:
25     models_file = pd.read_csv(f'{data_path}/models.csv", header=0)
26     for model_info in models_file.iterrows():
27         name = model_info[1][1]
28         uri = model_info[1][2]
29         print(f"Loading model: {name} from URI: {uri}") # Логируем информацию о моделях
30         list_models_for_validate[name + " " + uri] = mlflow.sklearn.load_model(uri)
31
32     current_results = {}
33     for name, model in list_models_for_validate.items():
34         prediction = model.predict(x_val)
35         current_results[name] = accuracy_score(y_val, prediction)
36
37     # Находим лучшую модель
38     best_model_in_list_validate_model = max(current_results, key=current_results.get)
39     model_name, model_uri = best_model_in_list_validate_model.split(" ")
40
41     # Сохранение лучшей модели в CSV
42     best_model_data = pd.DataFrame({
43         "name": [model_name],
44         "uri": [model_uri],
45         "accuracy": [current_results[best_model_in_list_validate_model]]
46     })
47
48     # Проверяем, существует ли уже файл с лучшими моделями

```

```

# Проверяем, существует ли уже файл с лучшими моделями
try:
    best_models_csv = pd.read_csv(f'{data_path}/best_models.csv")
    # Добавляем новую строку с лучшей моделью
    best_model_data.to_csv(f'{data_path}/best_models.csv", mode='a', header=False, index=False)
except FileNotFoundError:
    # Если файл не существует, создаем его с заголовками
    best_model_data.to_csv(f'{data_path}/best_models.csv", index=False)

print(f"Best model saved: {model_name} with URI: {model_uri}")

```

Этот скрипт делает следующее.

#### Инициализация и конфигурация:

- Устанавливает URL-адрес MLflow для отслеживания (`tracking_url`), путь к данным (`data_path`) и идентификатор текущего эксперимента, который читается из файла `experiment_id.txt`.
- Устанавливает текущий эксперимент для MLflow.

#### Загрузка данных:

- Загружает валидационные данные `x_val` и `y_val` из CSV-файлов, чтобы использовать их для проверки моделей.

#### Загрузка моделей для валидации:

- Загружает файл `models.csv`, содержащий список моделей с их именами и URI.
- Для каждой модели из файла загружает её из MLflow и добавляет в словарь `list_models_for_validate`, используя комбинацию имени и URI в качестве ключа для каждой модели.

#### Оценка моделей:

- Для каждой модели в `list_models_for_validate` делает предсказания на данных `x_val` и вычисляет точность предсказаний с помощью метрики `accuracy_score`.
- Сохраняет результаты точности каждой модели в словаре `current_results`.

#### Поиск лучшей модели:

- Определяет модель с наивысшей точностью на валидационных данных.
- Извлекает имя и URI этой модели.

#### Сохранение лучшей модели:

- Создаёт DataFrame `best_model_data` с информацией о лучшей модели (имя, URI и точность).
- Логирует сообщение о сохранённой лучшей модели.

### 3. Создание скрипта promote\_model

```
promote_model.py > ...
1  #!/usr/bin/env python3
2  import pandas as pd
3  import mlflow
4  from mlflow import MlflowClient
5  import subprocess
6  import requests
7  import time
8
9  tracking_url = 'http://mlflow_server:5000'
10 data_path = '/opt/airflow/data/lr3'
11
12 mlflow.set_tracking_uri(tracking_url)
13 client = MlflowClient()
14
15 models = client.search_registered_models()
16 for model in models:
17     print(model.name)
18
19 # Чтение последней строки из файла CSV с лучшими моделями
20 best_model_df = pd.read_csv(f"{data_path}/best_models.csv")
21 best_model_info = best_model_df.iloc[-1] # Берем последнюю строку
22
23 model_name = best_model_info['name']
24 model_uri = best_model_info['uri']
25
26 # Извлечение run_id из model_uri
27 run_id = model_uri.split('/')[1]
28
29 # Получение версии модели
30 print(f"Поиск модели {model_name} и run_id: {run_id}")
31 results = client.search_model_versions(f"name='{model_name}' and run_id='{run_id}'")
32 if not results:
33     print(f"No model versions found for name: {model_name} and run_id: {run_id}")
34 else:
35     version = results[0].version
36     client.transition_model_version_stage(name=model_name, version=version, stage="Production")
37     print(f"Model {model_name} version {version} transitioned to Production.")
38     model_version = f"models://{model_name}/{version}"
39     # Команда для запуска хостинга модели через mlflow serve
40     host_command = f"mlflow models serve -m {model_version} --host 0.0.0.0 --port 5001"
41     print(f"Запуск хостинга модели командой: {host_command}")
42
43
44
45 # Запуск команды для хостинга модели
46 return_code = subprocess.call(host_command, shell=True)
47 print(f"Return code: {return_code}")
48
```

```
time.sleep(5)

try:
    # Отправка POST-запроса для проверки
    response = requests.post(f'http://localhost/invocations',
                             json={
                                 "columns": ["Pclass", "Sex", "Age", "SibSp", "Parc
                                 "data": [
                                     [1, 0, 58.0, 0, 0, 26.55, 2],
                                     [2, 0, 23.0, 0, 0, 13.7917, 0],
                                     [3, 1, 23.0, 0, 0, 7.8542, 2]
                                 ]
                             })
    print(f"Response from model: {response.json()}")
except requests.exceptions.RequestException as e:
    print(f"Error during request: {e}")
```

Инициализация и подключение:

- Устанавливается URL-адрес для MLflow (tracking\_url), путь к данным (data\_path), а также подключение к клиенту MlflowClient, что позволяет управлять моделями в MLflow.



Поиск зарегистрированных моделей:

- Выводит на экран имена всех зарегистрированных моделей, доступных в MLflow, используя `search_registered_models()`.

Определение лучшей модели:

- Читает последний (последнюю строку) записанный результат из файла `best_models.csv`, где содержится информация о лучшей модели.
- Извлекает имя модели (`model_name`) и URI (`model_uri`), а также `run_id`, который используется для идентификации версии модели.

Поиск и перевод модели в "Production":

- Поиск версии модели с `model_name` и `run_id`.
- Если версия модели найдена, переводит её на "Production" с помощью `transition_model_version_stage`. Это обновляет стадию модели в MLflow.

Запуск сервера модели:

- Создает команду для запуска MLflow сервера на порту 5001 с использованием версии модели (`model_version`).
- Запускает команду через `subprocess.call`, что активирует сервер, который начнёт принимать запросы на предсказание.
- После запуска сервера делает паузу на 5 секунд, чтобы сервер успел инициализироваться.

Проверка работоспособности модели:

- Отправляет тестовый POST-запрос к запущенному серверу (по адресу `/invocations`), передавая несколько данных для проверки предсказания.
- Полученный ответ выводится на экран. Если возникает ошибка при выполнении запроса, выводит сообщение об ошибке.

Так должен был работать скрипт в теории, и да мне удалось переместить модель в `promote` состояние, но возникли проблемы с хостингом.

Сначала я подумал, что всё прошло отлично и команда для хостинга исполнилась так как пришла со статусом 0, и я подумал, что мне удалось замостить модель, но пост запрос не сработал и при проверке портов, можно было понять, что там пусто, я проверял порты в контейнере, для этого даже пришлось команду поставить в контейнер, но порт был пуст.

```

File /usr/local/lib/python3.10/site-packages/mlflow/store/artifact/s3_artifact_repo.py, in
lient
import boto3
ModuleNotFoundError: No module named 'boto3'
PS C:\Users\darti> docker exec -it ca7b68006b3d /bin/sh
#
# curl http://localhost:5003/invocations
/bin/sh: 2: curl: not found
# wget -qO- http://localhost:5003/invocations
/bin/sh: 3: wget: not found
# apt-get update && apt-get install -y curl
Get:1 http://deb.debian.org/debian bullseye InRelease [116 kB]
Get:2 http://deb.debian.org/debian-security bullseye-security InRelease [27.2 kB]
Get:3 http://deb.debian.org/debian bullseye-updates InRelease [44.1 kB]
Get:4 http://deb.debian.org/debian bullseye/main amd64 Packages [8066 kB]
Get:5 http://deb.debian.org/debian-security bullseye-security/main amd64 Packages [304 kB]
Get:6 http://deb.debian.org/debian bullseye-updates/main amd64 Packages [18.8 kB]
Fetched 8575 kB in 3s (2564 kB/s)

```

Но всё же запускать эту команду в airflow наверно не совсем корректно.

После я решил попробовать ввести команду напрямую и зашел в контейнер mlflow, но и тут ничего не получилось, контейнер mlflow, просто не воспринимал команду, которую я использую говорил, что mlflow вообще нет тут.

В конечном итоге я пытался перестроить контейнер и добавить новую службу для хостинга модели.

```

0
1 model_server:
2   restart: unless-stopped
3   image: python:3.8 # Укажите базовый образ с Python и MLflow
4   container_name: model_server
5   depends_on:
6     - web
7   ports:
8     - "5002:5000" # Порт для доступа к модели
9   environment:
10    - MLFLOW_S3_ENDPOINT_URL=http://minio:9000
11    - AWS_ACCESS_KEY_ID=${AWS_ACCESS_KEY_ID}
12    - AWS_SECRET_ACCESS_KEY=${AWS_SECRET_ACCESS_KEY}
13   command: >
14     mlflow models serve -m runs/6c5de147109f491fb88538d1ddbcfe06/sklearn.ensemble/RandomForestClassifier
15     --host 0.0.0.0 --port 5000

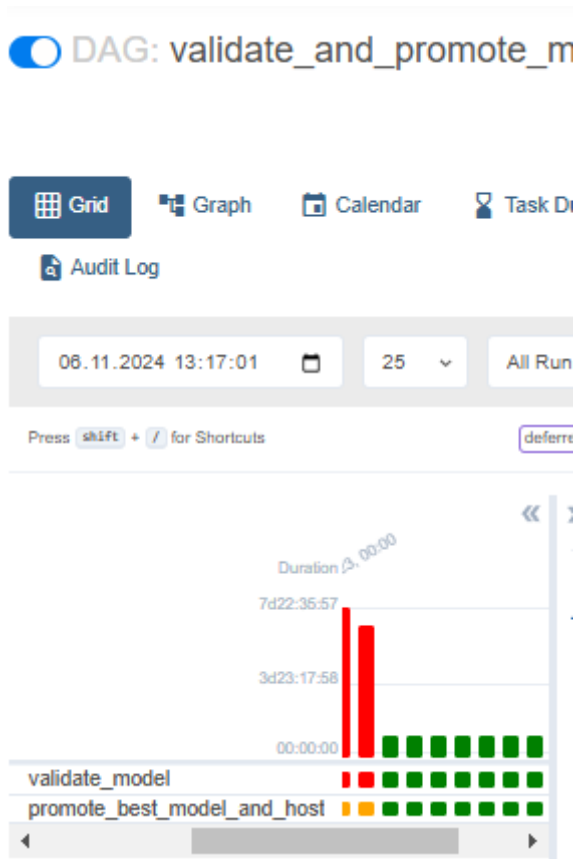
```

В общем я не уверен в том сделал ли я всё правильно.

Это тоже не совсем сработало, я так и не смог нормально захостить модель.

Решил посмотреть примеры в сети, и почему-то все люди ставили либо себе на пк либо на удалённый сервер, в итоге мне не удалось это реализовать для контейнера.

4. Выполнение dag



Так что мне пришлось довольствоваться переводом модели в promote

Но тем не менее лучшая модель оказалась в promote и основные задания были выполнены.

registered models

Filter registered models... ? Q

Name <span>⇅</span>	Latest version	Staging	Production	Created by	Last modified	Tags
DecisionTreeClassifier	Version 3	—	—		2024-10-20 23:...	—
KNeighborsClassifier	Version 3	—	—		2024-10-20 23:...	—
LogisticRegression	Version 2	—	—		2024-10-20 23:...	—
Perceptron	Version 1	—	—		2024-10-20 23:...	—
RandomForestClassifier	Version 3	—	Version 3		2024-10-20 23:...	—
SGDClassifier	Version 3	—	—		2024-10-20 23:...	—
SVC	Version 3	—	—		2024-10-20 23:...	—

## ЗАКЛЮЧЕНИЕ

В результате выполнения ЛР, мы создали пайплайн который принимает значения для предобработки данных и обучения моделей с заданными аргументами логирует все данные в mlflow. После мы создали пайплайн, который реализует валидацию моделей и выбирает лучшую модель и переводит её в promote, к сожалению хостинг не совсем получился, я понял что его можно реализовать через команду, но выполнить её в контейнере не получилось, возможно на настоящих серверах, с отдельной машиной с mlflow это бы получилось.