

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 и кибернетики  
Кафедра суперкомпьютеров и общей информатики

**Отчет по лабораторной работе №3**

Дисциплина: «Инженерия данных»

Тема: «**Airflow и MLflow – логгирование экспериментов и  
версионирование моделей**»

Выполнил: Неженский М.С.

Группа: 6233-010402D

Самара 2023

## 1. Пайплайн для обучения классификаторов

В рамках данной лабораторной работы предлагается построить два пайплайна. Первый – для обучения классификаторов. Для начала стоит создать конфигурационный файл. Я назвал его `config.json`, он написан на JavaScript. В качестве классификатор были выбраны:

- `AdaBoostClassifier`;
- `RandomForestClassifier`;
- `GaussianProcessClassifier`;
- `KNeighborsClassifier`;
- `MLPClassifier`.

Пример кода:

```
{
  "module": "sklearn.ensemble",
  "classifier": "AdaBoostClassifier",
  "args": {
    "n_estimators": 50,
    "learning_rate": 2.0
  }
}
```

В `module` записывается библиотека, из которой берем классификатор. В `classifier` записывается название самого классификатора. Далее в `args` записываются некоторые параметры классификатора, которые выбираем не по умолчанию.

Например, в данном случае выбраны самостоятельно максимальное количество оценок и размер шага на каждой итерации.

Таким образом создаем файл `config.json` с пятью выбранными классификаторами.

Код DAG для первого задания с комментариями:

```

airflow > dags > dag_train_data.py > ...
1 from datetime import datetime
2 from airflow import DAG
3 from airflow.sensors.filesystem import FileSensor
4 from airflow.operators.bash_operator import BashOperator
5 import os
6 os.environ["AWS_ACCESS_KEY_ID"] = "minio" # переменные среды для сохранения модели
7 os.environ["AWS_SECRET_ACCESS_KEY"] = "minio123"
8 os.environ["MLFLOW_S3_ENDPOINT_URL"] = "http://minio:9000"
9
10 default_args = { # аргументы по умолчанию
11     'owner': 'airflow', # владелец (можно сортировать в airflow)
12     'start_date': datetime(2023, 1, 1), # начало работы
13     'retries': 1, # количество повторных попыток, которое должно быть выполнено, прежде чем задача будет провалена
14 }
15
16 dag = DAG(
17     'train_models', # название DAG
18     default_args = default_args, # определение аргументов по умолчанию
19     description = 'DAG, which is needed to train models from sklearn', # описание DAG
20     schedule_interval = None, # автоматический запуск не планируется, запуск "извне"
21 )
22
23 monitoring = FileSensor( # используется для обнаружения файлов в локальной директории( в данном случае для обнаружения файла config.json)
24     task_id = 'monitoring', # название задания
25     poke_interval = 20, # время, в течение которого задание должно ждать между каждой попыткой
26     filepath = '/opt/airflow/data/config.json', # путь по которому происходит поиск файла и название файла
27     fs_conn_id = 'connection_inference', # название подключения (использовано, что и в прошлой ЛР)
28     dag = dag, # использование DAG
29 )
30
31 train_data = BashOperator(
32     task_id = "train_data", # название задания
33     bash_command = "python /opt/airflow/data/train_data.py", # запуск файла python из данной директории
34     dag = dag # использование DAG
35 )
36
37 monitoring >> train_data # последовательность заданий для выполнения, выполняются поочередно

```

Код на python, содержащий загрузку данных, деление выборки, запуск эксперимента для первого задания с комментариями:

```

airflow > data > train_data.py > ...
1 import json
2 import importlib
3 import pandas as pd
4 import numpy as np
5 from sklearn import datasets
6 from sklearn.model_selection import train_test_split
7 from sklearn.metrics import accuracy_score, precision_score, recall_score, f1_score
8 from sklearn.datasets import load_breast_cancer
9 import mlflow
10 import mlflow.sklearn
11 from mlflow.models import infer_signature
12
13 with open('/opt/airflow/data/config.json', 'r') as config_file: # открытие файла с конфигурацией классификаторов
14     config_data = json.load(config_file) # загрузка данных о классификаторах в переменную
15
16 breast_cancer = load_breast_cancer() # набор данных о раке молочной железы в штате Висконсин (классификация)
17
18 X_train, X_test, y_train, y_test = train_test_split( # разделение данных на тестовую и тренировочную выборку
19     breast_cancer.data, breast_cancer.target, test_size = 0.3, random_state = 33 # 70 к 30
20 )
21 X_validation, X_test, y_validation, y_test = train_test_split( # разделение данных на тестовую и валидационную выборку
22     X_test, y_test, test_size = 0.5, random_state = 33 # 50 к 50
23 )
24
25 pd.DataFrame(X_validation).to_csv('/opt/airflow/data/X_validation.csv', index=False) # запись в файл для второго задания
26 pd.DataFrame(y_validation).to_csv('/opt/airflow/data/y_validation.csv', index=False) # валидационной выборки
27
28 tracking_url = 'http://mlflow_server:5000' # задание ссылки для работы с MLflow
29 mlflow.set_tracking_uri(tracking_url) # подключение к серверу для работы с MLflow
30
31 experiment = mlflow.create_experiment("experiment02") # создание эксперимента
32 mlflow.set_experiment(experiment_id = experiment) # использование созданного эксперимента

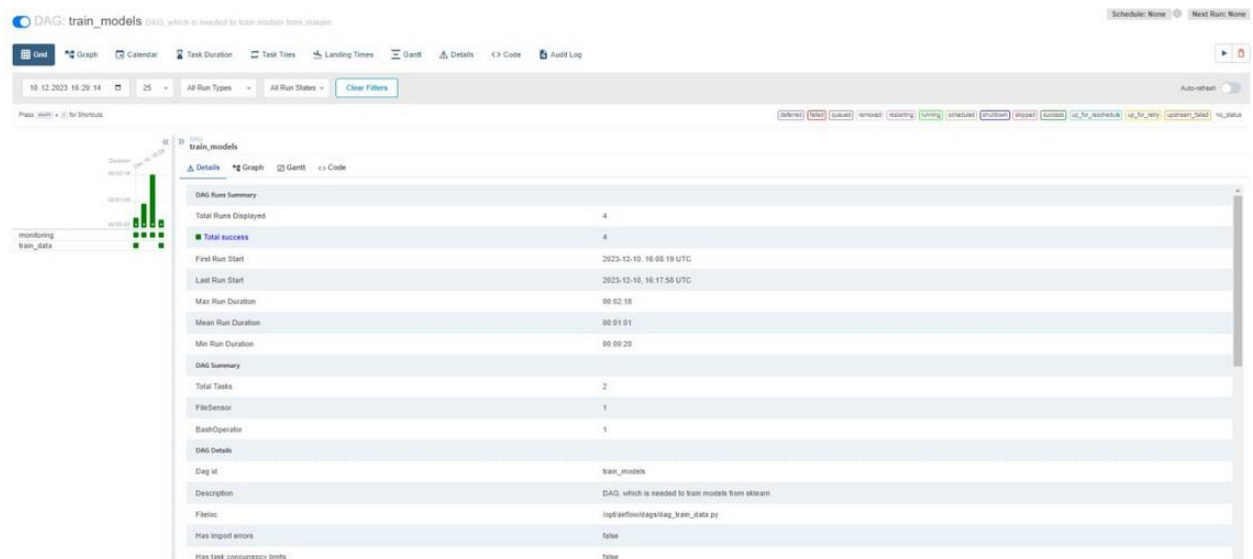
```

```

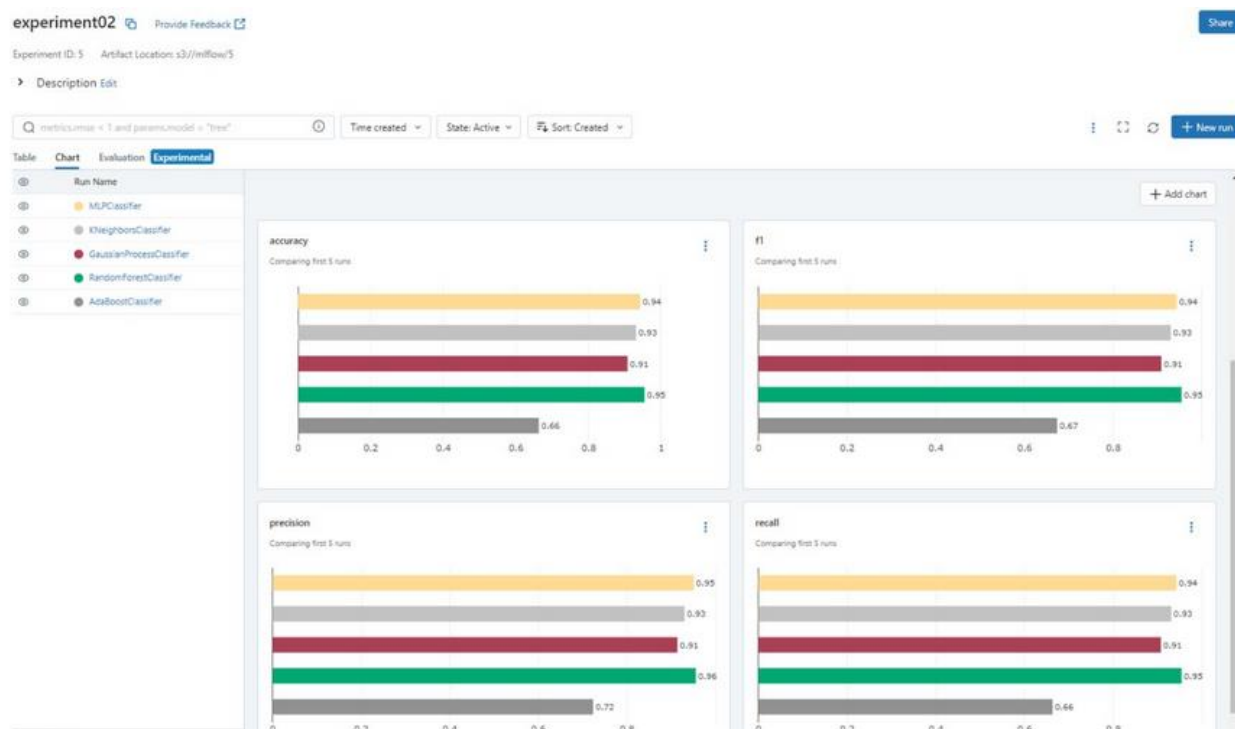
34 for i, config in enumerate(config_data['configs']): # для каждого классификатора из конфигурационного файла
35     mlflow.start_run(run_name = config['classifier'], experiment_id = experiment) # запуск в эксперименте
36     module = importlib.import_module(config['module']) # импортирование для обработки данных в конфигурационном файле
37     classifier = getattr(module, config['classifier']) # название классификатора из данных
38     model = classifier(**config['args']) # задание модели обучения через аргументы классификатора конфигурационного файла
39     model.fit(X_train, y_train) # обучение модели
40     prediction = model.predict(X_test) # предсказание модели
41     signature = infer_signature(X_test, prediction) # сигнатура
42
43     mlflow.log_params(config['args']) # логгирование параметров модели
44     mlflow.log_metrics({"accuracy": accuracy_score(y_test, prediction), # логгирование метрик
45                        "precision": precision_score(y_test, prediction, average = 'weighted'),
46                        "recall": recall_score(y_test, prediction, average = 'weighted'),
47                        "f1": f1_score(y_test, prediction, average = 'weighted')
48                        })
49
50     modelInfo = mlflow.sklearn.log_model( # логгирование модели
51         sk_model = model, # использованная модель для логгирования
52         artifact_path = config['module'], # модуль с помощью которого распознается модель
53         signature = signature, # сигнатура
54         registered_model_name = config['classifier']) # название модели
55
56     dataframe = pd.DataFrame({
57         "name": config['classifier'], # запись названия модели
58         "uri": modelInfo.model_uri # запись информации о модели
59     },
60     index=[i])
61     dataframe.to_csv('/opt/airflow/data/models.csv', mode = 'a', header = False) # запись моделей в файл
62     mlflow.end_run() # завершение работы

```

## Выполненный Dag:



Экран MLflow со значениями метрик для каждого классификатора:



## 2. Пайплайн для выбора лучшей модели

В данном задании выполняется все аналогично первому, DAG помещается в папку dags, а вспомогательный файл (validation.py) в папку data. Все это конечно в папке airflow. Код с комментариями приведен ниже. Код для DAG:

```

airflow > dags > dag_validation.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" # переменные среды для сохранения модели
7  os.environ["AWS_SECRET_ACCESS_KEY"] = "minio123"
8  os.environ["MLFLOW_S3_ENDPOINT_URL"] = "http://minio:9000"
9
10 default_args = { # аргументы по умолчанию
11     'owner': 'airflow', # владелец (можно сортировать в airflow)
12     'start_date': datetime(2023, 12, 10), # начало работы
13     'retries': 1, # количество повторных попыток, которое должно быть выполнено, прежде чем задача будет провалена
14 }
15
16 dag = DAG(
17     'validation', # название DAG
18     default_args=default_args, # определение аргументов по умолчанию
19     description='DAG for validate', # описание DAG
20     schedule_interval='@daily', # запуск один раз в день
21 )
22
23 validation = BashOperator(
24     task_id="validation", # название задания
25     bash_command="python /opt/airflow/data/validation.py", # запуск файла python из данной директории
26     dag=dag # использование DAG
27 )
28
29 validation # последовательность заданий для выполнения, выполняются поочередно

```

Код validation.py с комментариями:

```

airflow > data > validation.py > ...
1  import os
2  import operator
3  import pandas as pd
4  import numpy as np
5  from sklearn import datasets
6  # from sklearn.model_selection import train_test_split
7  from sklearn.metrics import f1_score
8
9  import mlflow
10 import mlflow.sklearn
11 from mlflow import MlflowClient
12
13 tracking_uri = 'http://mlflow_server:5000' # задание ссылки для работы с Mlflow
14 mlflow.set_tracking_uri(tracking_uri) # подключение к серверу для работы с Mlflow
15 mlflow.set_experiment("experiment02") # использование созданного эксперимента
16
17 X_validation = np.asarray(pd.read_csv(f"/opt/airflow/data/X_validation.csv"), dtype = np.float32) # считывание данных из файла валидационной выборки
18 y_validation = pd.read_csv(f"/opt/airflow/data/y_validation.csv")
19
20 list_models = {} # задание пустого объекта(модели)
21 mlflow.start_run(run_name = "Production model") # запуск в эксперименте Production model
22 models = pd.read_csv("/opt/airflow/data/models.csv", header = None) # считывание моделей
23 for model_info in models.iterrows(): # для каждой модели считывается записанное название и информация о модели
24     name_model = model_info[1][1]
25     uri_model = model_info[1][2]
26     list_models[name_model + " " + uri_model] = mlflow.sklearn.load_model(uri_model) # запись моделей
27

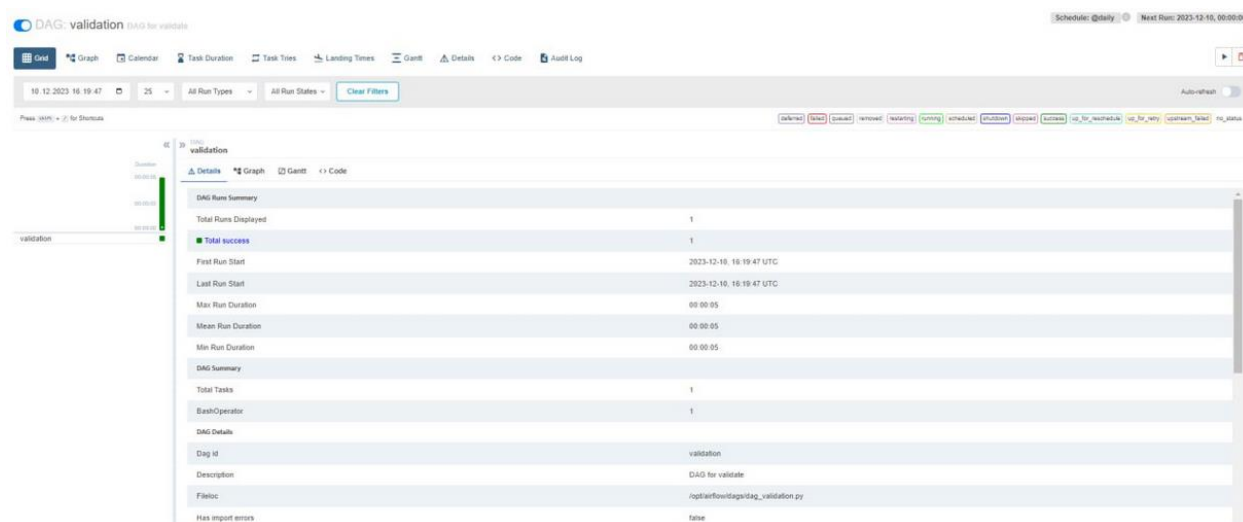
```

```

28 results = {} # задание пустого объекта(результат)
29 for i, j in list_models.items(): # посчет метрики f1 для каждой модели
30     pred = j.predict(X_validation)
31     results[i] = f1_score(
32         y_validation,
33         pred,
34         average="weighted"
35     )
36
37 best_model = max(results, key=results.get) # из всех моделей выбирается та, в которой метрика f1 максимальна
38 # запись лучшей версии и переопределение stage у лучшей версии с None на Production
39 version = MlflowClient().search_model_versions(f"name = '{best_model.split(' ')[0]}' and run_id = '{best_model.split(' ')[1].split('/')[1]}'")[0].version
40 MlflowClient().transition_model_version_stage(name = best_model.split(' ')[0], version = version, stage = "Production")
41 mlflow.end_run() # завершение работы

```

Выполненный DAG:



Зарегистрированные модели, как видим лучше всего в классификацией справился RandomForestClassifier:

Registered Models Create Model

Filter registered models by name or tags

Name	Latest version	Staging	Production	Created by	Last modified	Tags
AdaBoostClassifier	Version 1	—	—		2023-12-10 20:18:08	—
GaussianProcessClassifier	Version 1	—	—		2023-12-10 20:18:13	—
KNeighborsClassifier	Version 1	—	—		2023-12-10 20:18:15	—
QuadraticDiscriminantAnalysis	Version 1	—	—		2023-12-10 20:18:18	—
RandomForestClassifier	Version 1	—	Version 1		2023-12-10 20:19:51	—

## ЗАКЛЮЧЕНИЕ

В результате выполнения лабораторной работы были получены навыки по работе с Docker, Apache Airflow, MLflow. В данной работе были построены два пайплайна – для обучения классификаторов и выбора лучшей модели.