

MINISTRY OF EDUCATION OF REPUBLIC OF MOLDOVA
Technical University of Moldova
Faculty of Computers, Informatics and Microelectronics
Department of Software Engineering and Automatics

Report
on Artificial Intelligence Fundamentals
Laboratory Work nr. 4

Performed by:

Vasile Boaghi, Ion Dodon, FAF-172

Verified by:

Mihail Gavrilă, asist. univ.

Chişinău, 2021

Contents

Model Deployment	2
The Task	2
Solution Description	2
Code and Mentions	4
Conclusions	8
References	9

Model Deployment

Software deployment is all of the activities that make a software system available for use.

The general deployment process consists of several interrelated activities with possible transitions between them. These activities can occur at the producer side or at the consumer side or both. Because every software system is unique, the precise processes or procedures within each activity can hardly be defined. Therefore, "deployment" should be interpreted as a general process that has to be customized according to specific requirements or characteristics.

In software engineering, CI/CD or CICD generally refers to the combined practices of continuous integration and either continuous delivery or continuous deployment.

CI/CD bridges the gaps between development and operation activities and teams by enforcing automation in building, testing and deployment of applications. Modern day DevOps practices involve continuous development, continuous testing, continuous integration, continuous deployment and continuous monitoring of software applications throughout its development life cycle. The CI/CD practice or CI/CD pipeline forms the backbone of modern day DevOps operations.

The Task

The task is to make the linear regression model from laboratory work number 3 available for anyone and from anywhere. Also, is it needed to setup CI/CD workflow for the project. In GitHub action should be prepared a workflow that is triggered each time there is a push on the main branch and this workflow should test and build the application. A deployment workflow should be prepared as well and it should be triggered each time a new release has been made.

Solution Description

In laboratory work number 3 we had the model as a simple python file. That project only show that there is a model that does something. But, now we have to deliver the model model to the client to be used.

First of all we have to decide how the client will communicate with our model. One of the simplest solution is to use the HTTP protocol. This approach is probably not the best, because there can be a big amount of data that comes to the model to process, for example when a client what to make more predictions at once. But, we will use this protocol for simplicity.

If we want to communicate over HTTP, we have to prepare a HTTP server. Since our model is made and runs on python, the HTTP will also be made with python. Since in the case of this application, the most important part is the model and the server is used just to open an endpoint for communication, we will use Flask, because it is lightweight. An endpoint *GET /predict* will be available and it will take *application/json* content in the request's body.

Now, what is that data sanded by the client to the server/model? It is the data needed to make predictions. And we let it to be able to take more data sets in order to make more

predictions. The following is an example when it takes just one set of data and the model will make one prediction:

Listing 1: Input data example

```
1 {
2   "1":{
3     "1":-117.130000,
4     "2":-118.040000
5   },
6   "2":{
7     "1":32.700000,
8     "2":33.760000
9   },
10  "complexAge":{
11    "1":42.000000,
12    "2":25.000000
13  },
14  "totalRooms":{
15    "1":1210.000000,
16    "2":4061.000000
17  },
18  "totalBedrooms":{
19    "1":292.000000,
20    "2":545.000000
21  },
22  "complexInhabitants":{
23    "1":945.000000,
24    "2":1623.000000
25  },
26  "apartmentsNr":{
27    "1":258.000000,
28    "2":527.000000
29  },
30  "8":{
31    "1":0.899100,
32    "2":7.157200
33  },
34  "medianComplexValue":{
35    "1":78900.000000,
36    "2":294900.000000
37  }
38 }
```

The output for this input will be an array with the prediction for each data set from the input

Listing 2: Output data example

```
1 [
2   114116.93713370664,
3   332237.6029283451
4 ]
```

We made a simple service that will wrap the model and tested this service with two tests - a positive test and a negative test. We decided that if the client send non-valid data the client will receive a message in the form `{ "errorMessage": "Error. Check test data format!" }`. For simplicity, this message will be returned for any type of error in the server.

After having the server prepared and tested locally, we made a GitHub workflow that is triggered each time a push or pull request on the main branch has been made. This workflow runs the tests and if there are any non-passed tests the developer will be notified that the new code does not assure the correctness of the application. Then we create a workflow that build the application, creates a docker image, publishes the image on Docker Hub, spins a DigitalOcean^[4] droplet and on that droplet run a container with the new image.

Code and Mentions

The server is a simple Flask application with just one endpoint.

Listing 3: The endpoint and error handler

```
1 app = Flask(__name__)
2
3
4 @app.route('/predict')
5 def predict():
6     test_data = request.data.decode('utf-8')
7     prediction = predict_apartment_value(test_data)
8     return json.dumps(prediction)
9
10
11 @app.errorhandler(werkzeug.exceptions.HTTPException)
12 def handle_error(e):
13     response = e.get_response()
14     response.data = json.dumps({"errorMessage": "Error. Check test data format!"})
15     response.content_type = "application/json"
16     return response
```

We created a class that will represent the linear regression model. When the class is instantiated, the model will be trained.

The server is a simple Flask application with just one endpoint.

We cannot check the predicted values if they are correct or not, because the model will return different value each time it is re-started, because before it being trained, the model will get random weights and finally each time the model will make different predictions. Instead of checking the predicted values, we check is the number of predictions is correct.

Listing 4: Test the number of returned predictions

```
1 def test_predict_apartment_value_valid_json():
2     test_data: str = """{
3         "1":{
4             "1":-117.130000
5         },
6         "2":{
7             "1":32.700000
8         },
9         "complexAge":{
10            "1":42.000000
11        },
12        "totalRooms":{
13            "1":1210.000000
14        },
15        "totalBedrooms":{
16            "1":292.000000
17        },
18        "complexInhabitants":{
19            "1":945.000000
20        },
21        "apartmentsNr":{
22            "1":258.000000
23        },
24        "8":{
25            "1":0.899100
26        },
27        "medianCompexValue":{
28            "1":78900.000000
29        }
30    }"""
31     prediction = predict_apartment_value(test_data)
32     assert len(prediction) == 1
```

The following is the workflow that is triggered on each new release. It will create a Docker image and then run it on a DigitalOcean droplet and so the server will be available for any client.

Listing 5: *build_deploy.yml*

```
1 name: Build and deploy
2
3 on:
4   push:
5     tags:
6       - R*
7
8 jobs:
9   build:
10    runs-on: ubuntu-latest
11    steps:
12      - name: Checkout files
13        uses: actions/checkout@v2
14
15      - name: Set up Python 3.8.5
16        uses: actions/setup-python@v2
17        with:
18          python-version: 3.8.5
19
20      - name: Install dependencies
21        run: |
22          python -m pip install --upgrade pip
23          pip install flake8 pytest
24          if [ -f requirements.txt ]; then pip install -r requirements.txt; fi
25
26      - name: Lint with flake8
27        run: |
28          # stop the build if there are Python syntax errors or undefined names
29          flake8 . --count --select=E9,F63,F7,F82 --show-source --statistics
30          # exit-zero treats all errors as warnings. The GitHub editor is 127
31          chars wide
32          flake8 . --count --exit-zero --max-complexity=10 --max-line-length=127
33          --statistics
34
35      - name: Test with pytest
36        run: |
37          pytest -v
38
39      - name: Build and push Docker images
40        uses: docker/build-push-action@v1
41        with:
42          username: ${ secrets.DOCKERHUB.USERNAME }
43          password: ${ secrets.DOCKERHUB.PASSWORD }
44          repository: iondodon/lrmodel
```

```

43         tag_with_ref: true
44 #         tag_with_sha: true
45
46
47 deploy:
48     needs: build
49
50
51     runs-on: ubuntu-latest
52
53     steps:
54     - name: Checkout files
55       uses: actions/checkout@v2
56
57     - name: Get tag name
58       uses: olegtarasov/get-tag@v2.1
59
60     - name: Install doctl
61       uses: digitalocean/action-doctl@v2
62       with:
63         token: ${ secrets.DIGITALOCEAN_ACCESS_TOKEN }
64
65     - name: Provision Droplet and deploy container
66       run: doctl compute droplet create "$GIT_TAG_NAME" --image docker-18-04
        --size s-1vcpu-1gb --region nyc1 --user-data-file deploy.sh --wait

```


Conclusions

Working on this laboratory work we have learned how to continuously integrate out linear regression model by create a build workflow with GitHub actions and have also learned how to prepare a workflow for continuously delivery. We have understood how much the CI/CD principle help us to easier deliver any service after each new release and also how important CI is to assure that at any time we have a good working system. This was our first time when we prepared pipeline for CI/CD and now we are familiar with GitHub actions. We practiced how to make a docker image and publish it and how to work with Digital Ocean to host out service.

The application can now be tested by making a HTTP GET request to <http://157.245.253.166:5000/predict> and in the body should be a JSON as in the example from above.

References

- [1] Vasile Boaghi, Ion Dodon. Source code for the laboratory work. Accessed February 28, 2021. https://github.com/iondodon/FIA_Lab4.
- [2] Software Deployment https://en.wikipedia.org/wiki/Software_deployment.
- [3] Deploying to DigitalOcean With GitHub Actions https://www.digitalocean.com/community/tech_talks/deploying-to-digitalocean-with-github-actions
- [4] DigitalOcean <https://www.digitalocean.com/>