

**Model Deployment Process**

- In this task, I used Flask to create an API. I imported the ML classification model and wrote a function to predict and return a JSON response.
- Also, I used docker to containerized the API and serve the model easily.
- I wrote the unit tests to test the model for various scenarios like single/multiple prediction, missing valued data etc.

**Project structure**

root

- | - docker-compose.yml
- | - generate\_predictions.sh
- | - README.md
- | - .gitignore
- | - predictionapi
  - | - \_\_init\_\_.py
  - | - app.py
  - | - data\_preprocessing.py
  - | - requirements.txt
  - | - Dockerfile
  - | - unit\_test.py
  - | - endpoints
    - | - \_\_init\_\_.py
    - | - predict.py
- | - models
  - | - model.pkl
  - | - imputer.pkl
  - | - scaler.pkl
- | - test
  - | - \_\_init\_\_.py
  - | - test\_prediction\_api.py

**root**

- 'docker-compose.yml' is the compose file that documents and configures all of the application's service dependencies.
- generate\_predictions.sh' is the shell script that runs the 'docker-compose.yml' file.
- README.md is the project description file

## **predictionapi**

- 'app.py' is the main script that runs our app at port 1313.
- data\_preprocessing.py' is the script to preprocess the test data. It uses 'scaler.pkl' and 'imputer.pkl' to scale and impute the missing data, as well as create dummy variables and returns only those needed by the model.
- 'requirements.txt' is the file with all the packages needed to run the docker container.
- 'Dockerfile' is the dockerfile to create the docker image.
- 'unit\_test.py' is the script that runs all tests.

## **Endpoints**

### **Step 1: Create a prediction endpoint in predict.py**

- I initialized the API as a blueprint instead of a Flask app and called it "prediction\_api".

### **Step 2: Register the endpoint in app.py**

- Now that we have the prediction.py file to deal with the endpoint code, we can simply import the blueprint and use it in the Flask app.

## **Models**

I exported model and preprocessing variables as pickle file i.e 'model.pkl', 'imputer.pkl' and 'scaler.pkl' from the Notebook provided and used in data\_preprocessing.py' and predict.py

## **Test (unit testing)**

### **Step 1: Write tests(test\_prediction\_api.py)**

- I registered prediction\_api blueprint, defined in predict.py
- Later, I stored the app.test\_client() in a local variable called "tester". This will give us access to the API, as if we are hitting it with actual traffic.
- I defined the test functions, i.e. test\_predict\_single, test\_predict\_multiple, and test\_predict\_missing. Here, I send a POST request on the '/prediction' endpoint with the content\_type is JSON. The returned response is also in JSON format on which we can write assert statements..
- Here, I am only testing the API and not the model as the model has already been validated and it is assumed the predictions are accurate. Therefore, the assert statements are testing the status code returned(which should be 200) for the input JSON, and the response data is not null.

### **Step 2: Write a script to run all tests (unit\_test.py)**

I wrote a script called 'unit\_test.py' to run all the tests.