

Step 3: Building Serving ML Models with FastAPI

University: Tek-Up

Group: ING-4-SDIA

Tutor: Haythem Ghazouani

Contents

1	Overview	2
2	Step 1: Installation	2
3	Step 2: Understanding the Logic (app.py)	2
3.1	1. Model Loading	2
3.2	2. Input Validation (Pydantic)	2
3.3	3. Endpoints	2
4	Step 3: Running the Server	3
5	Step 4: Testing with Swagger UI	3
6	Step 5: Testing with Code (Client-Side)	3

1 Overview

FastAPI is a modern, high-performance framework for building APIs with Python. In this module, we use it to turn our trained Machine Learning pipeline into a usable service.

Implementation Note

Reference implementation: [code/app.py](#)

2 Step 1: Installation

To run the API, you need the framework and an ASGI server (Uvicorn). We also need `python-multipart` for file uploads.

```
1 pip install fastapi uvicorn python-multipart
```

Ensure these are in your `requirements.txt`.

3 Step 2: Understanding the Logic (app.py)

Our application follows a standard pattern:

3.1 1. Model Loading

We load the pickle file only **once** when the application starts, not for every request.

```
1 MODEL_PATH = "data/best_model_pipeline.pkl"
2 if os.path.exists(MODEL_PATH):
3     model = joblib.load(MODEL_PATH)
```

3.2 2. Input Validation (Pydantic)

We define a Class `CustomerData` that inherits from `BaseModel`. This acts as a strict contract. If a user sends "Male" as an integer, FastAPI will raise a `422 Validation Error`.

3.3 3. Endpoints

- **/health (GET)**: Returns status 200 if the model is loaded. Used by load balancers (like Kubernetes) to check if the pod is alive.
- **/predict (POST)**: Accepts JSON data for a single customer. Returns prediction and probability.
- **/predict_batch (POST)**: Accepts a `.csv` file using `UploadFile`, processes it with Pandas, and returns the result.

4 Step 3: Running the Server

Execute the following command from the root directory:

```
1 uvicorn code.app:app --reload
```

- `code.app`: Refers to the module `app.py` inside `code/` folder.
- `:app`: The instance of `FastAPI` created in the code.
- `--reload`: Automatically restarts the server when code changes (Dev mode only).

5 Step 4: Testing with Swagger UI

FastAPI automatically generates interactive documentation.

1. Open `http://localhost:8000/docs` in your browser.
2. Find the **POST /predict** endpoint.
3. Click **Try it out**.
4. Modify the JSON body and click **Execute**.

Common Pitfall

If you see "Internal Server Error", check your terminal. It usually means the input data format didn't match what the model pipeline expects (e.g., "Gender" vs "gender").

6 Step 5: Testing with Code (Client-Side)

You can verify your API programmatically:

```
1 import requests
2
3 url = "http://localhost:8000/predict"
4 payload = {
5     "Geography": "France",
6     "Gender": "Male",
7     "CreditScore": 600,
8     "Age": 40,
9     "Tenure": 3,
10    "Balance": 60000.0,
11    "NumOfProducts": 2,
12    "HasCrCard": 1,
13    "IsActiveMember": 1,
14    "EstimatedSalary": 50000.0
15 }
16 response = requests.post(url, json=payload)
17 print(response.json())
```