

eta_prediction

This is the **main project folder**. It contains all the components and files required for the ETA Provider Prediction project, including raw data, model artifacts, logs, src, and serving scripts. Below is a breakdown of the subfolders:

1. artifact

- This folder contains **pretrained models** and other important files needed for model inference. Each .pkl file represents a serialized object that can be loaded later for serving.
 - gradient_boosting_model.pkl: The gradient boosting model saved after training.
 - random_forest_model.pkl: The random forest model saved after training.
 - Kmeans_model.pkl: The kmeans used for clustering before classification.
-

2. data

- The folder responsible for **storing raw and processed datasets**.
 - **Subfolder:**
 - raw: This subfolder stores the initial raw data, in this case, the rides_data.pq file, which is a parquet file containing ride data used for training and evaluation.
-

3. exploratory_data_analysis

- This folder contains any scripts related to **exploratory data analysis (EDA)**. The EDA step is crucial to understanding the structure, distribution, and correlations within the data.
-

4. logs

- This folder stores logs and other tracking information, including **model weights for Neural Network Classifier**.
 - **Subfolder:**
 - model_weights/model1: This is where model weights are stored, specifically the model1.h5 file.
-

5. report

- This folder is intended for any **documentation or reporting** on the project. Reports include classification results of 3 different algorithm.
-

6. serving

- This folder is used for **model deployment** or serving the trained model. It contains all the necessary files to make the model accessible via API.
 - **Files:**
 - api.py: The script defining the API endpoint (using FastAPI) to serve the model for predictions.
 - model_service.py: A service or class that handles loading the model and running inference on incoming requests.
 - preprocessing.py: This file is used to preprocess the incoming data to ensure it's in the same format as the training data before passing it to the model.
 - **Subfolder:**
 - tests: Contains the test files to ensure that the API and preprocessing scripts are functioning correctly.
 - test_api.py: This is likely a test script for ensuring the API works as expected.
-

7. src

This folder represents the **source code** for the project. It houses all the code related to data processing, model training, and pipeline orchestration. This modular approach allows for clear separation between various components of the project, ensuring that each responsibility is isolated within its own file.

✓ data

- This folder contains all the **data preprocessing scripts**.
- **Files:**
 - __init__.py: This file is typically empty but makes the folder a Python package, allowing you to import from it.
 - data_processor.py: This file will contain all the logic related to **data preprocessing**. It includes functions to clean, process, and transform raw data into a format suitable for model training, including handling missing values, encoding, normalization, and feature engineering.

✓ models

- This folder contains all the **model training scripts**. Each model type is separated into its own file, adhering to the **Single Responsibility Principle** from SOLID design principles.
- **Files:**
 - `__init__.py`: Makes this folder a package.
 - `gradient_boosting_trainer.py`: This script handles **training of Gradient Boosting** models. It contains the model initialization, training process, and saving of the trained model.
 - `model_trainer.py`: This file is the **base class or abstract class** that provides a common interface for all model trainers. It defines shared methods and attributes that each model-specific trainer can inherit from.
 - `neural_network_trainer.py`: This file handles the **training of the neural network** (e.g., using TensorFlow). It contains the model architecture, hyperparameters, and training process for neural networks.
 - `random_forest_trainer.py`: This file is responsible for **training a Random Forest model**. It contains the model definition, training procedure, and any relevant hyperparameters.
 - `save_training_weights.py`: This file manages **saving model weights** during training. It handles checkpointing of models, useful in the case of neural networks.
 - `xgboost_trainer.py`: This script handles **training of XGBoost model**. It contains the initialization of the model, training routine, and saving process.

✓ pipeline

- This folder contains scripts responsible for **orchestrating the entire process** from data loading, model training, and evaluation. It centralizes all the steps required to complete the training process.
- **Files:**
 - `__init__.py`: A file to make this folder a Python package.
 - `train_pipeline.py`: This script is responsible for defining the **end-to-end pipeline** for training models. It likely combines steps like data preprocessing, model training, and evaluation. This ensures that you can run all steps sequentially without needing to call each part individually.
 - `main.py`: This is likely the **entry point** of the project. Running this file will execute the entire pipeline, kicking off all necessary steps for preprocessing, training, and saving models. This is the main file you run to trigger the workflow.

Additional Considerations:

- **Consistency:** The folder structure is modular and consistent, which makes it easier to maintain and navigate. Each component (e.g., data, models, logs, and serving) has its own folder, ensuring separation of concerns.
 - **Scalability:** This setup is scalable for future enhancements, like adding more models or expanding the API for additional functionalities.
-

This structured approach ensures a clean separation of different stages of the machine learning pipeline—data collection, model training, logging, evaluation, and serving—while adhering to clean code and design principles.

Problem Definition

The project requirement was to select the most accurate ETA provider from a set of options. To address this, we formulated the task as a **classification problem**. For each ride, we calculated the error between the actual arrival time (ATA) and the ETA provided by each provider. A new column, `accurate_provider`, was introduced to indicate the provider with the least error, serving as the target class for the classification task.

Data Preprocessing

1. Handling Missing Values:

- Missing values (NaNs) in the ETA columns for providers were handled by replacing them with a large value, ensuring no rows were dropped and the model would treat these cases appropriately.

2. Feature Engineering:

- **Datetime Features:** From the `accept_event_timestamp`, we extracted the following features:
 - `accept_hour`
 - `accept_day_of_week`
 - `accept_month`
 - `is_weekend` (a binary feature indicating if the ride occurred on a weekend).
- **Error Calculation:** For each provider (A, B, C, D), we calculated the absolute error between the ETA and the ATA, transforming the problem into a classification task where the provider with the lowest error is the target.
- **Clustering:** To enhance feature representation, we applied **KMeans clustering** on the data, adding a cluster feature to indicate the predicted cluster for each ride.

- **Scaling:** We applied **StandardScaler** to transform the data, ensuring consistent feature scaling for better model performance.

Model Training

We implemented and trained three different algorithms:

1. Neural Network:

- Architecture:
 - Two dense hidden layers with 20 and 40 units, respectively, using **ReLU activation**.
 - Dense output layer with **softmax activation** for multi-class classification.
- Optimized for classification by minimizing **sparse_categorical_crossentropy** loss.

2. Gradient Boosting:

- Parameters:
 - `n_estimators = 200`
 - `max_depth = 6`
 - `random_state = 42`

3. Random Forest:

- Parameters:
 - `n_estimators = 200`
 - `max_depth = 6`
 - `random_state = 42`

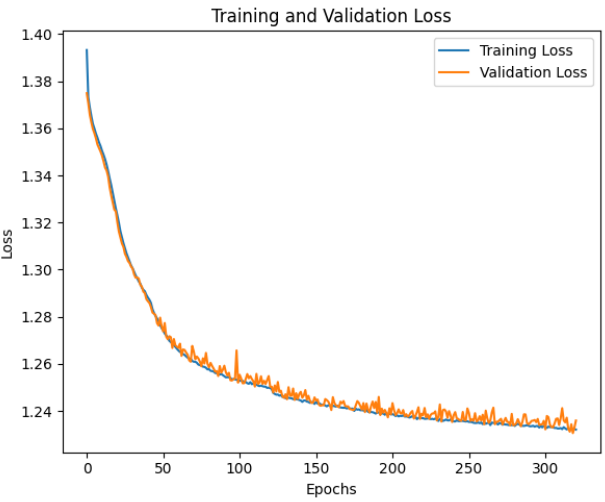
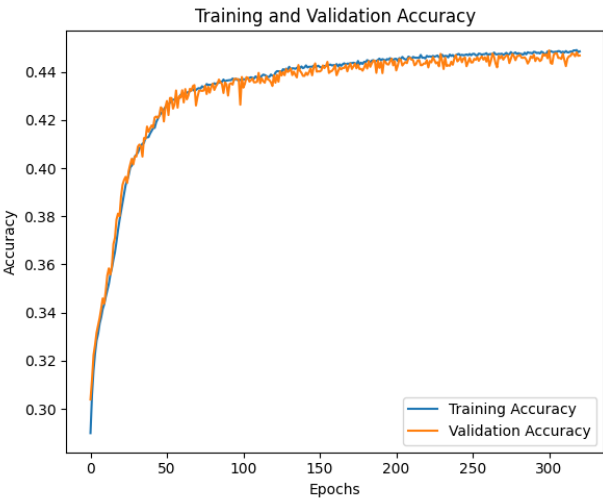
4. KNN:

- Parameters:
 - `n_neighbors = 10`

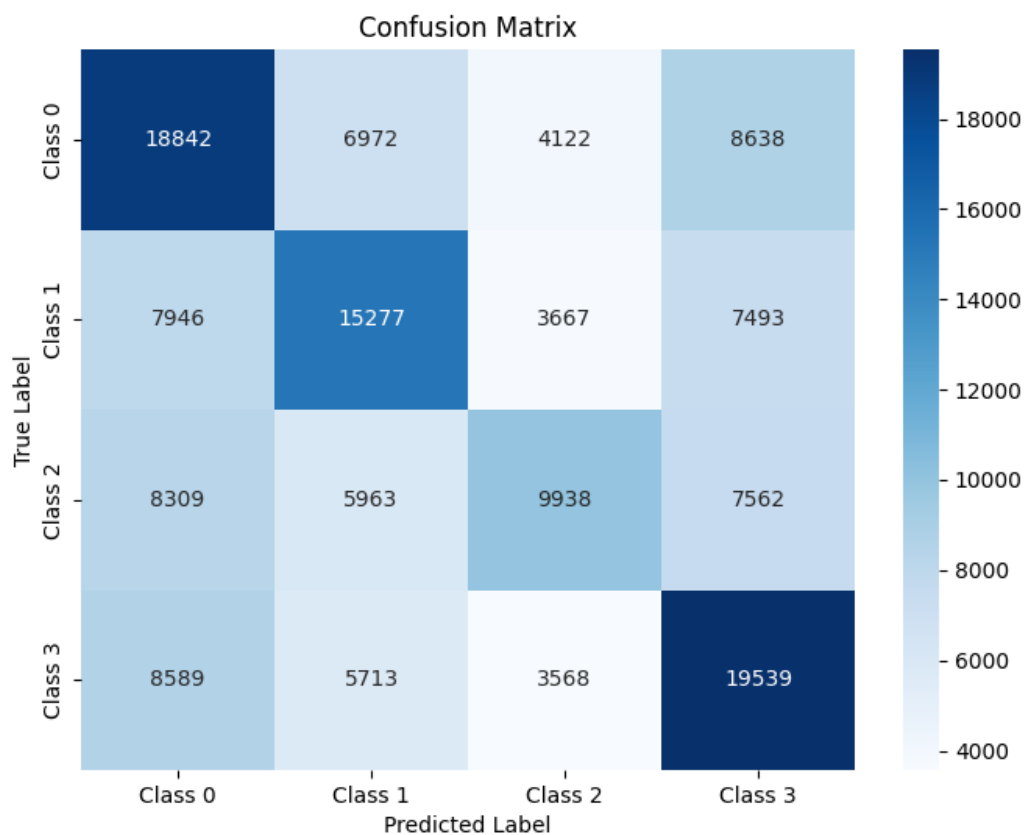
Results

The performance of each model was evaluated based on the precision, recall, and F1-score for each class (provider A, B, C, D). Below are the key results of the trained models:

✓ **Neural Network Results:**

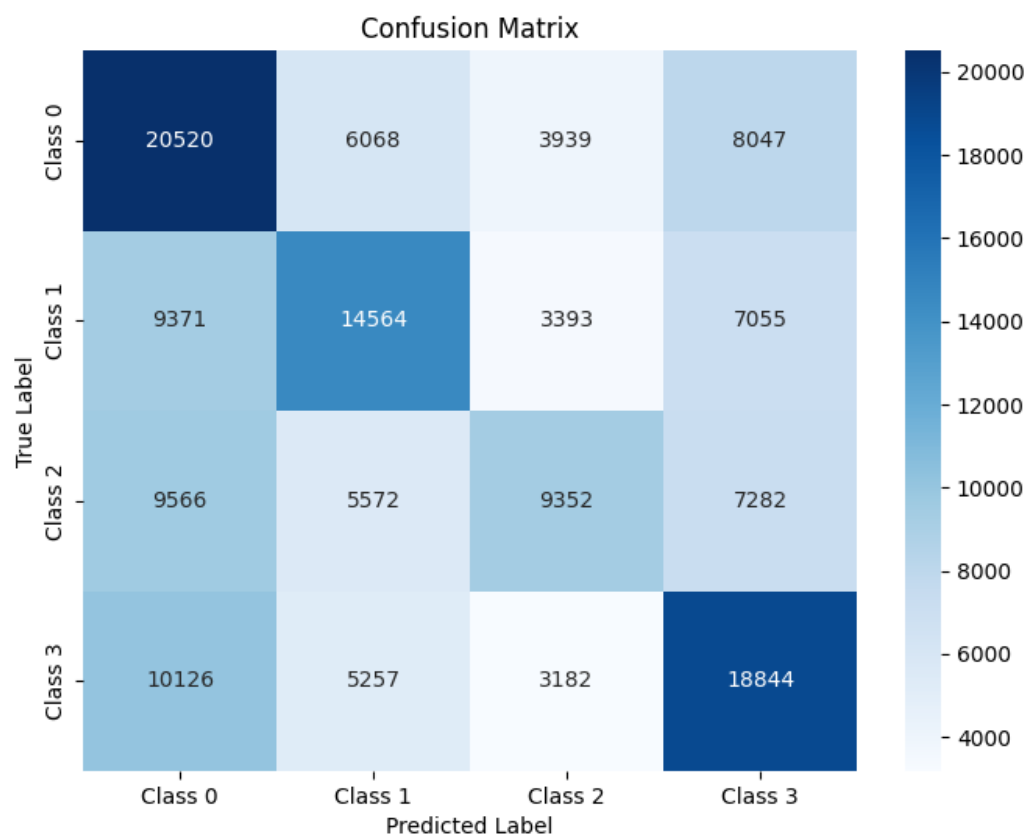


| Class | Precision | Recall | F1-score |
|-------------------------|-----------|--------|----------|
| Class 0 (provider_A) | 0.43 | 0.49 | 0.46 |
| Class 1 (provider_B) | 0.45 | 0.44 | 0.45 |
| Class 2 (provider_C) | 0.47 | 0.31 | 0.37 |
| Class 3 (provider_D) | 0.45 | 0.52 | 0.48 |



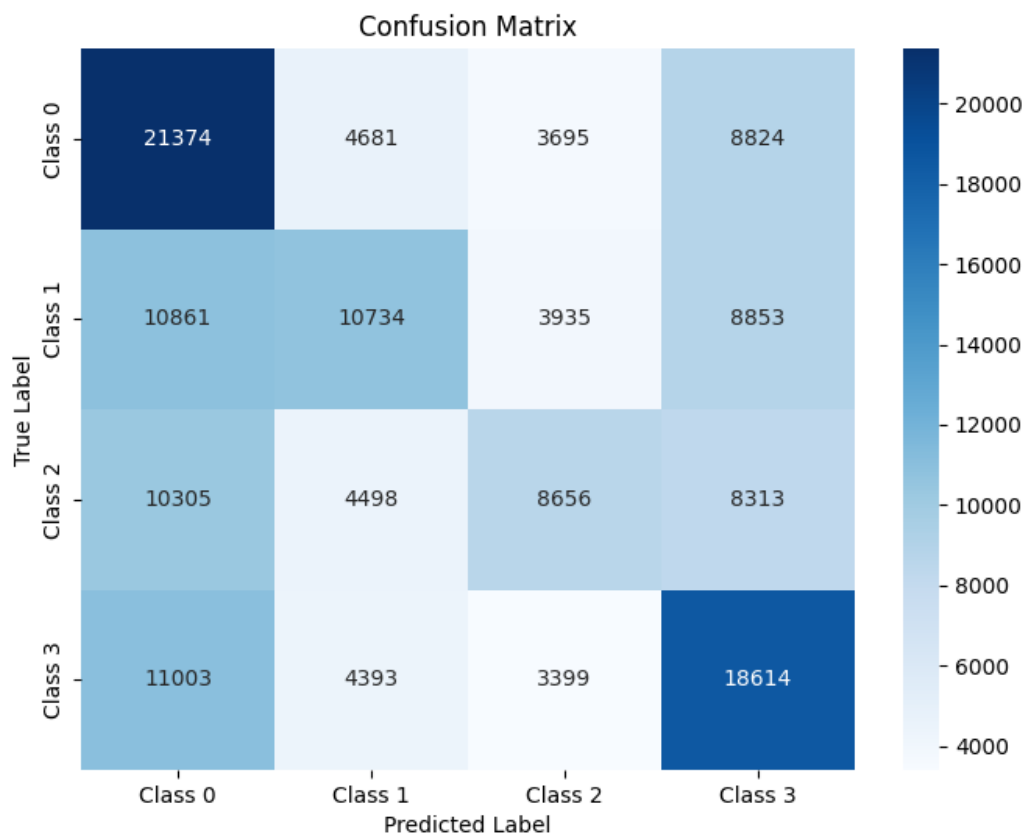
✓ **Random Forests:**

| Class | Precision | Recall | F1-score |
|-------------------------|-----------|--------|----------|
| Class 0 (provider_A) | 0.41 | 0.53 | 0.47 |
| Class 1 (provider_B) | 0.46 | 0.42 | 0.44 |
| Class 2 (provider_C) | 0.47 | 0.29 | 0.36 |
| Class 3 (provider_D) | 0.46 | 0.50 | 0.48 |



✓ **Gradient Boosting:**

| Class | Precision | Recall | F1-score |
|-------------------------|-----------|--------|----------|
| Class 0 (provider_A) | 0.40 | 0.55 | 0.46 |
| Class 1 (provider_B) | 0.44 | 0.31 | 0.37 |
| Class 2 (provider_C) | 0.44 | 0.27 | 0.34 |
| Class 3 (provider_D) | 0.42 | 0.50 | 0.45 |



Serving the Models

✓ Dockerization:

```
Terminal Local x Local (2) x + v
[+] Building 175.5s (10/11) docker:desktop-linux
=> CACHED [3/7] COPY artifact /app/artifact 0.0s
=> CACHED [4/7] COPY logs /app/logs 0.0s
=> CACHED [5/7] COPY serving /app/serving 0.0s
=> CACHED [6/7] COPY requirements.txt /app/requirements.txt 0.0s
=> [7/7] RUN pip install --no-cache-dir --upgrade pip && pip install --no-cache-dir -r requirements.txt 166.7s
=> # Downloading zipp-3.20.1-py3-none-any.whl.metadata (3.7 kB)
=> # Collecting pyasn1<0.7.0,>=0.4.6 (from pyasn1-modules==0.2.1->google-auth<3,>=1.6.3->tensorboard<2.14,>=2.13->tensorflow->-r requirements.txt (line 1))
=> # Downloading pyasn1-0.6.1-py3-none-any.whl.metadata (8.4 kB)
=> # Collecting oauthlib>=3.0.0 (from requests-oauthlib==0.7.0->google-auth-oauthlib<1.1,>=0.5->tensorboard<2.14,>=2.13->tensorflow->-r requirements.txt (line 1))
=> # Downloading oauthlib-3.2.2-py3-none-any.whl.metadata (7.5 kB)
=> # Downloading tensorflow-2.13.1-cp38-cp38-manylinux_2_17_x86_64_manylinux2014_x86_64.whl (479.6 MB)
```

docker desktop PERSONAL Search for images, containers, volumes, extensions and ... Ctrl+K Sign In

Containers

Images

Volumes

Builds

Docker Scout

Extensions

beautiful_liskov
ata_pred_app:latest
ba561429f981
8080:5000

STATUS
Running (0 seconds ago)

Logs Inspect Bind mounts Exec Files Stats

2024-09-13 17:19:21 Collecting usage statistics. To deactivate, set browser.gatherUsageStats to false.
2024-09-13 17:19:21
2024-09-13 17:19:21 You can now view your Streamlit app in your browser.
2024-09-13 17:19:21 URL: http://0.0.0.0:5000
2024-09-13 17:21:02 2024-09-13 13:51:02.774497: I tensorflow/tsl/cuda/cudart_stub.cc:28] Could not find cuda drivers on your machine, GPU will not be used.
2024-09-13 17:21:02 2024-09-13 13:51:02.816246: I tensorflow/tsl/cuda/cudart_stub.cc:28] Could not find cuda drivers on your machine, GPU will not be used.
2024-09-13 17:21:02 2024-09-13 13:51:02.816672: I tensorflow/core/platform/cpu_feature_guard.cc:182] This TensorFlow binary is optimized to use available CPU instructions in performance-critical operations.
2024-09-13 17:21:02 To enable the following instructions: AVX2 FMA, in other operations, rebuild TensorFlow with the appropriate compiler flags.
2024-09-13 17:21:03 2024-09-13 13:51:03.598548: W tensorflow/compiler/tf2tensorrt/utils/py_utils.cc:38] TF-TRT Warning: Could not find TensorRT
2024-09-13 17:21:16
2001-01-01 00:00:00 [=====] - ETA: 0s
2001-01-01 00:00:00 [=====] - 0s 66ms/step

Engine running RAM 5.11 GB CPU ---% Disk --- GB avail. of --- GB BETA Terminal New version available 5

✓ **Streamlit UI:**



Deploy

ETA Prediction

City ID

C

Accept Event Timestamp

2024-09-10T12:34:56Z

Origin Latitude

35.69

- +

Origin Longitude

51.39

- +

Destination Latitude

36.29

- +

Destination Longitude

52.39

- +

EDD

12000

- +

Provider A

3600

- +

Provider B

3700

- +

3600

- +

Provider B

3700

- +

Provider C

3400

- +

Provider D

3300

- +

Predict

Prediction: {'neural_network': [0], 'random_forest': [2]}