

File Formats Description

Flight Cancellation Prediction System

Data Science Project

January 9, 2026

1 Introduction

This document describes the file formats used for the “Prediction Objects” in the Flight Cancellation Prediction System. These files enable the server to load machine learning models and the data preparation pipeline at startup.

2 Manifest File (prediction_objects.json)

The entry point for the system configuration. It acts as a catalog, telling the server where to find the pipeline and model artifacts.

2.1 Format

JSON (JavaScript Object Notation)

2.2 Structure

```
1 {
2     "description": "String description of the file",
3     "pipeline_definition": "Path to the Python file (e.g., pipeline.py)"
4     ,
5     "models": {
6         "model_key_1": "Path to serialized model file",
7         "model_key_2": "Path to serialized model file"
8     },
9     "preprocessing": "Path to pipeline state file (e.g., models/pipeline
10 .joblib)",
11     "hyperparameters": {
12         "model_key": {
13             "algorithm": "Name of the algorithm class",
14             "params": {
15                 "param_name": "param_value"
16             }
17         }
18     },
19     "pipeline_description": {
20         "steps": ["List of strings describing steps"],
21         "scaling_method": "String name of scaling method",
22         "encoding_method": "String name of encoding method"
23     }
24 }
```

2.3 Key Fields

Field	Description
<code>description</code>	Human-readable description of the manifest
<code>pipeline_definition</code>	Path to the Python file containing the <code>PredictionPipeline</code> class
<code>models</code>	Dictionary mapping model keys to their <code>.joblib</code> file paths
<code>preprocessing</code>	Path to the consolidated pipeline state file
<code>hyperparameters</code>	Dictionary of hyperparameters used to train each model
<code>pipeline_description</code>	Steps and methods used in the data preparation pipeline

3 Model Files (*.joblib)

Each classification model is serialized into its own binary file using the `joblib` library. This allows for efficient storage and loading of large NumPy arrays often found in scikit-learn models.

3.1 Format

Joblib Binary Serialization (Python-specific)

3.2 Content

A single Python object corresponding to a trained scikit-learn estimator.

3.3 Available Models

File	Scikit-learn Class
<code>models/naive_bayes.joblib</code>	<code>sklearn.naive_bayes.GaussianNB</code>
<code>models/knn.joblib</code>	<code>sklearn.neighbors.KNeighborsClassifier</code>
<code>models/logistic_regression.joblib</code>	<code>sklearn.linear_model.LogisticRegression</code>
<code>models/decision_tree.joblib</code>	<code>sklearn.tree.DecisionTreeClassifier</code>
<code>models/mlp.joblib</code>	<code>sklearn.neural_network.MLPClassifier</code>
<code>models/random_forest.joblib</code>	<code>sklearn.ensemble.RandomForestClassifier</code>

4 Pipeline State File (models/pipeline.joblib)

The data preparation pipeline's internal state is consolidated into a single binary file to satisfy the “One single file” requirement from the project specification.

4.1 Format

Joblib Binary Serialization (Python-specific)

4.2 Content

A Python dictionary containing all necessary objects to transform raw data.

4.3 Dictionary Keys

Key	Content
"encoder"	Fitted <code>sklearn.preprocessing.OrdinalEncoder</code>
"scaler"	Fitted <code>sklearn.preprocessing.MinMaxScaler</code>
"cat_cols"	List of categorical column names
"shifts"	Dictionary mapping columns to shift values (to handle negative encoded integers)
"mvi_stats"	Dictionary containing mean/mode values for missing value imputation
"final_features"	List of column names expected by the trained models

5 Pipeline Definition (`pipeline.py`)

While not a data file, this Python source file defines the `PredictionPipeline` class logic. The server imports this class and populates it with the data loaded from `models/pipeline.joblib`.

5.1 Key Methods

- `__init__()` — Initializes the pipeline and loads all artifacts from disk.
- `transform(data)` — Applies the full transformation pipeline to input data.
- `predict_single(instance, model_name)` — Predicts a single instance using the specified model.
- `evaluate(df, model_name)` — Evaluates model performance on a validation dataset.