

Python Script Documentation: XGBoost on Tom Brady's Superbowl Win Prediction

1 Introduction

This document provides detailed documentation for a Python script that uses the XGBoost classifier to predict Superbowl wins based on NFL data. The script loads, processes data, trains, and evaluates models for different year ranges. This documentation covers:

- A breakdown of the code, explaining each section.
- A layered explanation of loops and nested operations.
- A flowchart visualizing the structure of the script.
- Libraries used, main functionality, and error handling.

2 Libraries Used

The script relies on the following libraries:

- `os`: Provides functionality for handling directories and file paths.
- `pandas`: Used for loading, manipulating, and preprocessing the dataset.
- `sklearn.preprocessing.StandardScaler`: For standardizing feature data.
- `sklearn.metrics`: Provides tools for calculating performance metrics such as accuracy, precision, recall, F1 score, and confusion matrices.

- `xgboost.XGBClassifier`: Implements the XGBoost classifier, a powerful gradient boosting model for classification tasks.
- `matplotlib.pyplot`: Used for plotting the confusion matrix and setting visual styles for plots.

3 Main Functionality

The script performs the following tasks:

- Load NFL data from a CSV file.
- Preprocess the data (standardize features and handle missing values).
- Train and evaluate XGBoost models for predicting Superbowl wins over specific year ranges.
- Calculate and save performance metrics (accuracy, precision, recall, F1 score) and confusion matrices.
- Visualize the results and save plots and metrics to a specified directory.

4 Code Breakdown

4.1 Function: `xgboost_classification`

This function trains and evaluates an XGBoost model for a specified range of years.

- **Input Parameters:**
 - `train_year_start`: The start year of the training period.
 - `train_year_end`: The end year of the training period.
 - `features`: A pandas DataFrame containing the feature variables.
 - `target`: A pandas Series representing the target variable (Superbowl Win).
 - `save_directory`: The path where model outputs (plots and metrics) will be saved.

- **Steps:**
 - Split the data into training and testing sets based on the specified year range.
 - Standardize the features using **StandardScaler**.
 - Create and fit an XGBoost model with the training data.
 - Generate predictions for the test set.
 - Plot the confusion matrix and save the plot as a PNG file.
 - Calculate evaluation metrics: accuracy, precision, recall, and F1 score.
 - Save these metrics to a text file.

4.2 Main Function

The **main** function orchestrates the loading, processing, and model training workflow.

- Define paths for loading data and saving results.
- Load NFL data from a CSV file using pandas.
- Set the DataFrame index to **Season** for easier year-based operations.
- Split the dataset into features and the target variable (Superbowl Win).
- Remove columns with missing values from the feature set.
- Call the **xgboost_classification** function for two different year ranges (2002–2009 and 2002–2019) to train and evaluate models.

5 Explanation of Loops and Nested Operations

Although there are no explicit loops in the script, nested operations are present in the following key steps:

- **StandardScaler:** The **fit_transform** method is used to standardize the training data, while **transform** is applied to the test data.

- **XGBoost Classifier:** The `fit` method trains the model on the provided training data. Internally, this involves iterative optimization across multiple boosting rounds.

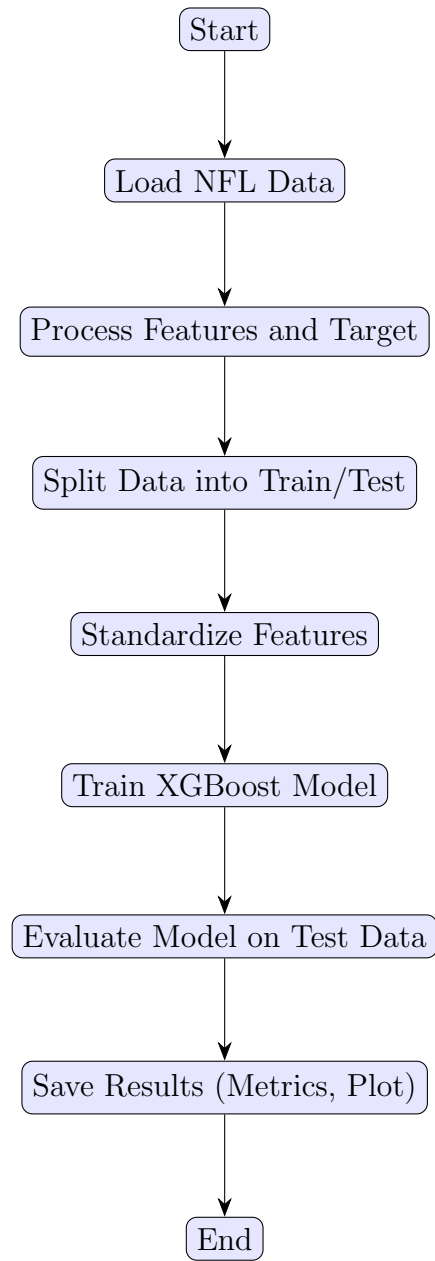
6 Error Handling

The script currently does not implement explicit error handling, but some improvements could be made:

- Ensure that the dataset exists before attempting to load it.
- Check if the save directories exist, and create them if necessary.
- Handle any issues that may arise during model training or metric calculations with `try-except` blocks.

7 Flowchart

The following flowchart illustrates the structure of the script, showing the main functions and the flow of data from loading to evaluation:



8 Conclusion

This Python script trains and evaluates an XGBoost model to predict Superbowl wins based on NFL data. It preprocesses data by standardizing features, trains the model on a specified year range, and evaluates its performance using accuracy, precision, recall, and F1 score. Confusion matrices and performance metrics are saved for further analysis. This LaTeX documentation outlines the structure and function of the script, including its main operations and possible areas for error handling.