

## Overview

This project aims to predict the type of pitch thrown in baseball games using multiclass classification techniques. The workflow involves data exploration, feature engineering, model training, hyperparameter tuning, and model evaluation to develop a robust prediction model.

## Libraries and Imports

The following libraries are used for various tasks such as data manipulation, visualization, and machine learning:

- **Pandas:** Used for data manipulation and analysis. Provides DataFrame structures to handle large datasets efficiently.
- **NumPy:** Utilized for numerical operations, particularly array manipulations, which are integral in data preprocessing and model computations.
- **Matplotlib & Seaborn:** Libraries for creating static, animated, and interactive visualizations. Seaborn builds on Matplotlib and provides a high-level interface for drawing attractive and informative statistical graphics.
- **Scikit-learn:** A machine learning library in Python that features various classification, regression, and clustering algorithms, and tools for model evaluation and preprocessing.
- **Warnings:** Handles and suppresses warnings during code execution to maintain clean output.

## Data Loading and Exploration

- **Loading Data:** The dataset is read from a CSV file named 'pitches'. This data contains information on each pitch thrown in a series of baseball games.
- **Summary Statistics:** The `describe(include='all')` method provides a summary of the dataset, including count, mean, standard deviation, min, max, and percentile information for each feature. This helps in understanding the initial distribution and characteristics of the data.

## Pitch Type Distribution

- **Visualization:** A bar plot is generated to show the frequency of each pitch type in the dataset. This helps in identifying the most and least common pitch types, providing insight into the data's balance.

## Feature Engineering

- **Removing Rare Pitch Types:** Pitches labeled as 'AB' are removed due to their extremely low frequency, which could negatively affect model performance and cross-validation.
- **Creating Binary Features:**
  - **same\_stance\_throw:** This binary feature indicates whether the batter and pitcher have the same handedness. It is created based on the assumption that same-handedness influences the pitch type.
  - **on\_1b, on\_2b, on\_3b:** These features indicate the presence of a player on first, second, or third base, respectively. They are converted from their original format to binary, reflecting whether a base is occupied.

- **Pitch Tendencies:** Fastball tendency (**FF\_tendency**) for each pitcher is calculated as a running average. This feature reflects a pitcher's likelihood to throw a fastball based on historical data.
- **Times Through Rotation:** This feature counts how many times a pitcher has faced the batting order. It is based on the idea that pitchers might change their pitch selection strategy as they face the lineup multiple times.

## Data Preparation

- **Scaling Features:** StandardScaler is used to standardize the features by removing the mean and scaling to unit variance. This step is crucial as it ensures that all features contribute equally to the model, improving the convergence speed and performance of many machine learning algorithms.
- **Train-Test Split:** The dataset is split into training and testing sets using an 80-20 split. This allows for the evaluation of the model's performance on unseen data, ensuring that the model generalizes well.

## Model Training

Several classifiers are trained and evaluated to determine the best model for pitch type prediction:

- **KNeighborsClassifier:** A simple, instance-based learning algorithm where predictions are based on the majority class among the nearest neighbors.
- **GaussianNB:** A probabilistic classifier based on Bayes' theorem, assuming Gaussian distribution of features.
- **RandomForestClassifier:** An ensemble learning method that constructs multiple decision trees and outputs the mode of the classes for classification tasks.
- **LogisticRegression:** A linear model for classification that predicts the probability of categorical outcomes using a logistic function.

## Primary Metric

Brier Score Loss is used to evaluate the quality of predicted probabilities. This metric is preferable for multiclass classification as it measures the accuracy of probabilistic predictions rather than just the final classification.

## Hyperparameter Tuning

- **GridSearchCV:** Used to perform an exhaustive search over specified parameter values for an estimator. It helps in finding the optimal hyperparameters for the Logistic Regression model within a pipeline that includes scaling and calibration.
- **StratifiedKFold Cross-Validation:** This ensures that the proportion of each class is approximately the same in each fold, which is particularly important for datasets with imbalanced classes.

## Model Evaluation

- **Confusion Matrix:** A table used to describe the performance of a classification model by showing the actual versus predicted classifications. It is visualized using a heatmap for clarity.

- **Calibration Curve:** Plots the relationship between predicted probabilities and actual outcomes, helping to assess the accuracy of the predicted probabilities for each class.

## Key Functions and Variables

### Variables

- **data:** DataFrame containing the pitch data.
- **X, y:** Feature matrix and target vector.
- **X\_train, X\_test, y\_train, y\_test:** Train-test split data.
- **scaler:** Instance of StandardScaler for feature scaling.
- **classifiers:** List of classifiers to evaluate.
- **best\_classifier, best\_brier\_score:** Track the best-performing classifier and its score.
- **pipeline, grid\_search:** Used for hyperparameter tuning and model calibration.

### Functions

- **preprocess\_data:** Performs feature engineering steps to prepare the data for modeling.
- **train\_and\_evaluate\_models:** Trains multiple classifiers and evaluates them using Brier Score Loss.
- **evaluate\_model:** Generates and visualizes the confusion matrix and calibration curves for the final model.

### Illustrative Images

- **Pitch Type Distribution:** A bar plot showing the frequency of each pitch type, aiding in understanding the class distribution.
- **Confusion Matrix Heatmap:** Visual representation of model performance, indicating how often predictions match actual labels.
- **Calibration Curves:** Graphs showing the calibration of predicted probabilities, assessing the accuracy of probabilistic predictions for each class.