

Forecasting and Analyzing Baseball

Double Play Analytics



Overview

This project analyzes the factors associated with forecasting baseball game outcomes and provides key recommendations for success to industry leaders.

Business Understanding

In the highly competitive and financially driven landscape of professional baseball, accurately predicting game outcomes can significantly impact team performance, fan engagement, and revenue generation. By leveraging historical game data our predictive modeling approach aims to provide teams, broadcasters, and stakeholders with actionable insights into future game outcomes. This predictive analytics tool will empower teams to optimize their strategies, such as player selection, pitching rotations, and in-game tactics, leading to improved win rates and overall performance. Ultimately, a reliable predictive model for baseball game outcomes has the potential to drive fan excitement, increase ticket sales, and attract valuable sponsorships, thus contributing to the long-term success and profitability of baseball organizations.

Data and Analysis Limitations

This analysis spanned the COVID-19 pandemic. The pandemic introduced unique challenges for analyzing baseball data, requiring careful consideration of context, data quality, and potential biases when interpreting and modeling data from 2020 and beyond. In 2020 there were 60 regular season games per team, less than half of the standard 162. The shortened and interrupted season impacted player salaries in 2020 and may have affected player performance metrics such as batting averages, earned run averages (ERAs), and fielding

percentages. Some players may have performed better or worse than expected due to factors like altered training routines, health concerns, or changes in game dynamics. COVID-19 outbreaks and health protocols may have impacted player availability due to positive cases, contact tracing, or precautionary measures. This could affect team rosters, playing time, and lineup strategies, leading to shifts in player statistics and team performance. Many sports events in 2020 were held without spectators or with limited attendance to comply with public health guidelines. The absence of fans in stadiums and arenas could have affected player morale, home field advantage, and game dynamics, potentially influencing game outcomes and performance metrics.

Baseball has one of the richest data caches in sports, offering a plethora of statistics and metrics for analysis. Nonetheless, navigating this vast sea of information poses challenges, demanding a focused approach, organization, and an understanding of big data modeling techniques to extract meaningful insights from the wealth of available data.

Data and Analysis Preparation

```
In [1]: # Bringing in packages for EDA, pre-processing, modeling, and visualizations
import pandas as pd
import numpy as np

import matplotlib.pyplot as plt
import seaborn as sns
%matplotlib inline

from sklearn.preprocessing import OneHotEncoder, StandardScaler
from sklearn.impute import MissingIndicator, SimpleImputer
from sklearn.dummy import DummyClassifier

from sklearn.model_selection import train_test_split, cross_val_score
from sklearn.ensemble import RandomForestClassifier
from sklearn.tree import DecisionTreeClassifier, plot_tree
from sklearn import tree
from sklearn.linear_model import LogisticRegression
from sklearn.model_selection import GridSearchCV
from sklearn.metrics import classification_report
from sklearn.metrics import accuracy_score, precision_score, confusion_matrix,
                           recall_score, f1_score

import torch
import torch.nn as nn
import torch.optim as optim
from torch.utils.data import DataLoader, TensorDataset
from torchvision import transforms
import torch.nn.functional as F
```

Create a single dataset

```
In [2]: # Each year is saved in a separate .txt file on retrosheets.org
file_paths = ["data/gl2010.txt", "data/gl2011.txt", "data/gl2012.txt", "data/g
            "data/gl2015.txt", "data/gl2016.txt", "data/gl2017.txt", "data/g
            "data/gl2020.txt", "data/gl2021.txt", "data/gl2022.txt", "data/g
```

```
In [3]: dfs=[]
```

```
In [4]: # Pulling together the individual year files
for file_path in file_paths:
    with open(file_path, "r") as f:
        data = f.readlines()
        data_split = [line.strip().split(",") for line in data]
        df_initial = pd.DataFrame(data_split)
    dfs.append(df_initial)
```

```
In [5]: # Concatenate all DataFrames in the list
df = pd.concat(dfs, ignore_index=True)
```

General Data Understanding

This section does a high level over view of the data. It looks at the initial head, tail, shape, description, missing values, data types, and counts.

```
In [6]: # Setting up display
# Set the display width to accommodate more characters per row
pd.set_option('display.width', 1000) # Adjust as needed

# Set display options
pd.set_option('display.max_rows', None)
pd.set_option('display.max_columns', None)

# Display the head of the DataFrame
df.info()
```

```
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 32484 entries, 0 to 32483
Columns: 179 entries, 0 to 178
dtypes: object(179)
memory usage: 44.4+ MB
```

In [7]: df.head()

Out[7]:

	0	1	2	3	4	5	6	7	8	9	10	11	12	13	14	15	16
0	"20100404"	"0"	"Sun"	"NYA"	"AL"	1	"BOS"	"AL"	1	7	9	51	"N"	""	""	""	"BOS07"
1	"20100405"	"0"	"Mon"	"MIN"	"AL"	1	"ANA"	"AL"	1	3	6	51	"N"	""	""	""	"ANA01"
2	"20100405"	"0"	"Mon"	"CLE"	"AL"	1	"CHA"	"AL"	1	0	6	51	"D"	""	""	""	"CHI12"
3	"20100405"	"0"	"Mon"	"DET"	"AL"	1	"KCA"	"AL"	1	8	4	54	"D"	""	""	""	"KAN06"
4	"20100405"	"0"	"Mon"	"SEA"	"AL"	1	"OAK"	"AL"	1	5	3	54	"N"	""	""	""	"OAK01"



In [8]: df.tail()

Out[8]:

	0	1	2	3	4	5	6	7	8	9	10	11	12	13	14	15
32479	"20231001"	"0"	"Sun"	"SDN"	"NL"	162	"CHA"	"AL"	162	2	1	66	"D"	""	""	""
32480	"20231001"	"0"	"Sun"	"CLE"	"AL"	162	"DET"	"AL"	162	2	5	51	"D"	""	""	""
32481	"20231001"	"0"	"Sun"	"NYA"	"AL"	162	"KCA"	"AL"	162	2	5	51	"D"	""	""	""
32482	"20231001"	"0"	"Sun"	"TEX"	"AL"	162	"SEA"	"AL"	162	0	1	51	"D"	""	""	""
32483	"20231001"	"0"	"Sun"	"TBA"	"AL"	162	"TOR"	"AL"	162	12	8	54	"D"	""	""	""



In [9]: df.describe()

Out[9]:

	0	1	2	3	4	5	6	7	8	9	10
count	32484	32484	32484	32484	32484	32484	32484	32484	32484	32484	32
unique	2417	3	7	31	2	163	31	2	163	27	26
top	"20200904"	"0"	"Sat"	"TOR"	"NL"	36	"PHI"	"NL"	19	3	3
freq	20	31518	5352	1091	16479	222	1088	16491	225	4625	4547



```
In [10]: df.shape
```

```
Out[10]: (32484, 179)
```

```
In [11]: nan_values = df.isna().any()  
nan_values
```

```
Out[11]: 0      False  
1      False  
2      False  
3      False  
4      False  
5      False  
6      False  
7      False  
8      False  
9      False  
10     False  
11     False  
12     False  
13     False  
14     False  
15     False  
16     False  
17     False  
18     False  
...    ...
```

```
In [12]: df.dtypes
```

```
Out[12]: 0      object  
1      object  
2      object  
3      object  
4      object  
5      object  
6      object  
7      object  
8      object  
9      object  
10     object  
11     object  
12     object  
13     object  
14     object  
15     object  
16     object  
17     object  
18     object  
...    ...
```

```
In [13]: df.count()
```

```
153    32484
154    32484
155    32484
156    32484
157    32484
158    32484
159    32484
160    32484
161      218
162      218
163      218
164      192
165      171
166      171
167       67
168       61
169       60
170       17
171       14
172       14
```

Feature preparation and data manipulation

This section renames columns, drops rows with missing column values of interest, drops columns, converts objects to numeric columns, and strips characters from object columns.

```
In [14]: # Dropping the columns with minimal information at the END of the dataset
# Dropping and individual player, manager, and umpire information; this analysis
df.drop(columns=df.columns[77:179], inplace=True)
```


In [15]: # Renaming the numeric columns to something more descriptive

```
new_column_names = {
    0: 'date',
    1: 'num_games',
    2: 'day_of_week',
    3: 'team_visiting',
    4: 'league_visiting',
    5: 'game_num_visiting',
    6: 'team_home',
    7: 'league_home',
    8: 'game_num_home',
    9: 'score_visiting',
    10: 'score_home',
    11: 'outs_in_game',
    12: 'time_of_day',
    13: 'game_completed',
    14: 'forfeit',
    15: 'protest',
    16: 'park_id',
    17: 'attendance',
    18: 'length_min',
    19: 'line_score_visiting',
    20: 'line_score_home',

    # Offense - Visiting
    21: 'at_bats_visiting',
    22: 'hits_visiting',
    23: 'double_visiting',
    24: 'triple_visiting',
    25: 'home_run_visiting',
    26: 'rbi_visiting',
    27: 'sacrifice_hit_visiting',
    28: 'sacrifice_fly_visiting',
    29: 'hit_by_pitch_visiting',
    30: 'walk_visiting',
    31: 'intent_walk_visiting',
    32: 'strikeout_visiting',
    33: 'stolen_base_visiting',
    34: 'caught_stealing_visiting',
    35: 'grounded_into_double_plays_visiting',
    36: 'first_catcher_interfere_visiting',
    37: 'left_on_base_visiting',

    # Pitching - Visiting
    38: 'pitchers_used_visiting',
    39: 'individual_earned_runs_visiting',
    40: 'team_earned_runs_visiting',
    41: 'wild_pitches_visiting',
    42: 'balks_visiting',

    # Defense - Visiting
    43: 'putouts_visiting',
    44: 'assists_visiting',
    45: 'errors_visiting',
    46: 'passed_balls_visiting',
    47: 'double_def_visiting',
    48: 'triple_def_visiting',}
```

```

# Offense - Home
49: 'at_bats_home',
50: 'hits_home',
51: 'double_home',
52: 'triple_home',
53: 'home_run_home',
54: 'rbi_home',
55: 'sacrifice_hit_home',
56: 'sacrifice_fly_home',
57: 'hit_by_pitch_home',
58: 'walk_home',
59: 'intent_walk_home',
60: 'strikeout_home',
61: 'stolen_base_home',
62: 'caught_stealing_home',
63: 'grounded_into_double_plays_home',
64: 'first_catcher_interfere_home',
65: 'left_on_base_home',

# Pitching - Home
66: 'pitchers_used_home',
67: 'individual_earned_runs_home',
68: 'team_earned_runs_home',
69: 'wild_pitches_home',
70: 'balks_home',

# Defense - Home
71: 'putouts_home',
72: 'assists_home',
73: 'errors_home',
74: 'passed_balls_home',
75: 'double_def_home',
76: 'triple_def_home',
}

# Rename columns
df.rename(columns=new_column_names, inplace=True)

```

In [16]: # Removing quotations around strings--

```

columns_to_convert_strings = [
    'num_games',
    'date',
    'day_of_week',
    'team_visiting',
    'league_visiting',
    'team_home',
    'league_home',
    'time_of_day',
    'park_id',
    'line_score_visiting',
    'line_score_home']

```

```
In [17]: df[columns_to_convert_strings] = df[columns_to_convert_strings].applymap(lambda
```



```
In [18]: # Converting most numeric columns from strings to integers
```

```
columns_to_convert = [
    "num_games",
    "game_num_visiting",
    "game_num_home",
    "score_visiting",
    'game_num_home',
    'score_visiting',
    'score_home',
    "outs_in_game",
    "attendance",
    "length_min",
    "at_bats_visiting",
    "hits_visiting",
    "double_visiting",
    "triple_visiting",
    "home_run_visiting",
    "rbi_visiting",
    "sacrifice_hit_visiting",
    'sacrifice_fly_visiting',
    'hit_by_pitch_visiting',
    'walk_visiting',
    'intent_walk_visiting',
    'strikeout_visiting',
    'stolen_base_visiting',
    'caught_stealing_visiting',
    'grounded_into_double_plays_visiting',
    'first_catcher_interfere_visiting',
    'left_on_base_visiting',
    'pitchers_used_visiting',
    'individual_earned_runs_visiting',
    'team_earned_runs_visiting',
    'wild_pitches_visiting',
    'balks_visiting',
    'putouts_visiting',
    'assists_visiting',
    'errors_visiting',
    'passed_balls_visiting',
    'double_def_visiting',
    'triple_def_visiting',
    'at_bats_home',
    'hits_home',
    'double_home',
    'triple_home',
    'home_run_home',
    'rbi_home',
    'sacrifice_hit_home',
    'sacrifice_fly_home',
    'hit_by_pitch_home',
    'walk_home',
    'intent_walk_home',
    'strikeout_home',
    'stolen_base_home',
    'caught_stealing_home',
    'grounded_into_double_plays_home',
    'first_catcher_interfere_home',
    'left_on_base_home',
```

```
'pitchers_used_home',
'individual_earned_runs_home',
'team_earned_runs_home',
'wild_pitches_home',
'balks_home',
'putouts_home',
'assists_home',
'errors_home',
'passed_balls_home',
'double_def_home',
'triple_def_home']
```

```
In [19]: # Convert columns to numeric, coercing errors to NaN
for column in columns_to_convert:
    df[column] = pd.to_numeric(df[column], errors='coerce')
```

```
In [20]: df.head()
```

Out[20]:

	date	num_games	day_of_week	team_visiting	league_visiting	game_num_visiting	team
0	20100404	0	Sun	NYA	AL		1
1	20100405	0	Mon	MIN	AL		1
2	20100405	0	Mon	CLE	AL		1
3	20100405	0	Mon	DET	AL		1
4	20100405	0	Mon	SEA	AL		1



```
In [21]: df.describe()
```

Out[21]:

	num_games	game_num_visiting	game_num_home	score_visiting	score_home	outs_ir
count	32484.000000	32484.000000	32484.000000	32484.000000	32484.000000	32484.
mean	0.044607	80.078654	80.078408	4.356976	4.472663	53.
std	0.268992	46.946529	46.954208	3.157740	3.080320	4.
min	0.000000	1.000000	1.000000	0.000000	0.000000	27.
25%	0.000000	39.000000	39.000000	2.000000	2.000000	51.
50%	0.000000	79.000000	79.000000	4.000000	4.000000	54.
75%	0.000000	121.000000	121.000000	6.000000	6.000000	54.
max	2.000000	163.000000	163.000000	28.000000	29.000000	120.



In [22]: df.count()

```
Out[22]: date                                32484
num_games                               32484
day_of_week                             32484
team_visiting                            32484
league_visiting                           32484
game_num_visiting                        32484
team_home                                32484
league_home                               32484
game_num_home                            32484
score_visiting                            32484
score_home                               32484
outs_in_game                             32484
time_of_day                               32484
game_completed                            32484
forfeit                                   32484
protest                                    32484
park_id                                    32484
attendance                                 31549
length_min                                32446
line_score_visiting                      32484
line_score_home                            32484
at_bats_visiting                          32479
hits_visiting                            32484
double_visiting                           32446
triple_visiting                           32446
home_run_visiting                         32484
rbi_visiting                             32484
sacrifice_hit_visiting                   32484
sacrifice_fly_visiting                   32484
hit_by_pitch_visiting                    32484
walk_visiting                            32484
intent_walk_visiting                     32484
strikeout_visiting                       32484
stolen_base_visiting                     32484
caught_stealing_visiting                 32484
grounded_into_double_plays_visiting     32484
first_catcher_interfere_visiting        32484
left_on_base_visiting                    32484
pitchers_used_visiting                  32484
individual_earned_runs_visiting          32484
team_earned_runs_visiting                32484
wild_pitches_visiting                   32484
balks_visiting                            32484
putouts_visiting                           32484
assists_visiting                          32484
errors_visiting                           32484
passed_balls_visiting                   32484
double_def_visiting                      32484
triple_def_visiting                     32484
at_bats_home                             32484
hits_home                                32484
double_home                               32484
triple_home                               32484
home_run_home                            32484
rbi_home                                  32484
sacrifice_hit_home                       32484
sacrifice_fly_home                       32484
```

hit_by_pitch_home	32484
walk_home	32484
intent_walk_home	32484
strikeout_home	32484
stolen_base_home	32484
caught_stealing_home	32484
grounded_into_double_plays_home	32484
first_catcher_interfere_home	32484
left_on_base_home	32484
pitchers_used_home	32484
individual_earned_runs_home	32484
team_earned_runs_home	32484
wild_pitches_home	32484
balks_home	32484
putouts_home	32484
assists_home	32484
errors_home	32484
passed_balls_home	32484
double_def_home	32484
triple_def_home	32484
dtype: int64	

In [23]: df.dtypes

```
Out[23]: date                                     object
         num_games                                int64
         day_of_week                               object
         team_visiting                            object
         league_visiting                          object
         game_num_visiting                      int64
         team_home                                object
         league_home                               object
         game_num_home                           int64
         score_visiting                           int64
         score_home                               int64
         outs_in_game                            int64
         time_of_day                             object
         game_completed                          object
         forfeit                                 object
         protest                                 object
         park_id                                 object
         attendance                             float64
         length_min                            float64
         line_score_visiting                     object
         line_score_home                          object
         at_bats_visiting                       float64
         hits_visiting                           int64
         double_visiting                         float64
         triple_visiting                        float64
         home_run_visiting                      int64
         rbi_visiting                           int64
         sacrifice_hit_visiting                  int64
         sacrifice_fly_visiting                  int64
         hit_by_pitch_visiting                  int64
         walk_visiting                           int64
         intent_walk_visiting                   int64
         strikeout_visiting                     int64
         stolen_base_visiting                   int64
         caught_stealing_visiting                int64
         grounded_into_double_plays_visiting    int64
         first_catcher_interfere_visiting      int64
         left_on_base_visiting                  int64
         pitchers_used_visiting                 int64
         individual_earned_runs_visiting       int64
         team_earned_runs_visiting              int64
         wild_pitches_visiting                 int64
         balks_visiting                          int64
         putouts_visiting                        int64
         assists_visiting                        int64
         errors_visiting                         int64
         passed_balls_visiting                  int64
         double_def_visiting                    int64
         triple_def_visiting                    int64
         at_bats_home                           int64
         hits_home                              int64
         double_home                            int64
         triple_home                            int64
         home_run_home                          int64
         rbi_home                               int64
         sacrifice_hit_home                     int64
         sacrifice_fly_home                     int64
```

```
hit_by_pitch_home           int64
walk_home                   int64
intent_walk_home            int64
strikeout_home              int64
stolen_base_home            int64
caught_stealing_home        int64
grounded_into_double_plays_home int64
first_catcher_interfere_home int64
left_on_base_home           int64
pitchers_used_home          int64
individual_earned_runs_home int64
team_earned_runs_home       int64
wild_pitches_home           int64
balks_home                  int64
putouts_home                 int64
assists_home                 int64
errors_home                  int64
passed_balls_home            int64
double_def_home              int64
triple_def_home              int64
dtype: object
```

```
In [24]: # Dropping rows with missing values in the listed columns
# at_bats_visiting - 5 rows with missing data
# double_visiting - 38 rows with missing data
# triple_visiting - 38 rows with missing data
# A total of 38 rows will be deleted; there are 43 duplicates in the 81 rows i
df.dropna(subset=["at_bats_visiting", "double_visiting", "triple_visiting"], i
```

```
In [25]: # Drop the one row where the game ended in a tie; found this obs because when
df.drop(index=16957, inplace=True)
```

```
In [26]: # Drop rows where there was a protest
# A total of 11 rows will be deleted
df.dropna(subset=['protest'], inplace=True)
df = df[(df['protest'] != 'V') & (df['protest'] != 'H')]
```

```
In [27]: # This drops the partial 2020 season; not sure I want to drop it. It was an o
# but unless I'm looking at year as a factor what is the impact of keeping it
# df = df[df['year'] != 2020]
```

```
In [28]: df.drop(columns=["attendance", "num_games", "length_min", "game_completed", "f
```

```
In [29]: # FLO/MIA is the only team that switched three letter codes during this time p
df['team_visiting'].replace({"FLO": "MIA"}, inplace=True)
df['team_home'].replace({"FLO": "MIA"}, inplace=True)
```

In [30]: df.head()

Out[30]:

	date	team_visiting	team_home	score_visiting	score_home	outs_in_game	at_bats_visit
0	20100404	NYA	BOS	7	9	51	3
1	20100405	MIN	ANA	3	6	51	3
2	20100405	CLE	CHA	0	6	51	3
3	20100405	DET	KCA	8	4	54	3
4	20100405	SEA	OAK	5	3	54	3



In [31]: df.info()

```

<class 'pandas.core.frame.DataFrame'>
Int64Index: 32434 entries, 0 to 32483
Data columns (total 62 columns):
 #   Column           Non-Null Count Dtype
 ---  ----
 0   date             32434 non-null  object
 1   team_visiting    32434 non-null  object
 2   team_home         32434 non-null  object
 3   score_visiting   32434 non-null  int64
 4   score_home        32434 non-null  int64
 5   outs_in_game     32434 non-null  int64
 6   at_bats_visiting 32434 non-null  float64
 7   hits_visiting    32434 non-null  int64
 8   double_visiting   32434 non-null  float64
 9   triple_visiting   32434 non-null  float64
 10  home_run_visiting 32434 non-null  int64
 11  rbi_visiting     32434 non-null  int64
 12  sacrifice_hit_visiting 32434 non-null  int64
 13  sacrifice_fly_visiting 32434 non-null  int64
 14  hit_by_pitch_visiting 32434 non-null  int64
 15  walk_visiting     32434 non-null  int64
 16  intent_walk_visiting 32434 non-null  int64
 17  strikeout_visiting 32434 non-null  int64
 18  stolen_base_visiting 32434 non-null  int64
 19  caught_stealing_visiting 32434 non-null  int64
 20  grounded_into_double_plays_visiting 32434 non-null  int64
 21  first_catcher_interfere_visiting 32434 non-null  int64
 22  left_on_base_visiting 32434 non-null  int64
 23  pitchers_used_visiting 32434 non-null  int64
 24  individual_earned_runs_visiting 32434 non-null  int64
 25  team_earned_runs_visiting 32434 non-null  int64
 26  wild_pitches_visiting 32434 non-null  int64
 27  balks_visiting     32434 non-null  int64
 28  putouts_visiting   32434 non-null  int64
 29  assists_visiting   32434 non-null  int64
 30  errors_visiting   32434 non-null  int64
 31  passed_balls_visiting 32434 non-null  int64
 32  double_def_visiting 32434 non-null  int64
 33  triple_def_visiting 32434 non-null  int64
 34  at_bats_home       32434 non-null  int64
 35  hits_home          32434 non-null  int64
 36  double_home         32434 non-null  int64
 37  triple_home         32434 non-null  int64
 38  home_run_home      32434 non-null  int64
 39  rbi_home            32434 non-null  int64
 40  sacrifice_hit_home 32434 non-null  int64
 41  sacrifice_fly_home 32434 non-null  int64
 42  hit_by_pitch_home   32434 non-null  int64
 43  walk_home           32434 non-null  int64
 44  intent_walk_home    32434 non-null  int64
 45  strikeout_home      32434 non-null  int64
 46  stolen_base_home    32434 non-null  int64
 47  caught_stealing_home 32434 non-null  int64
 48  grounded_into_double_plays_home 32434 non-null  int64
 49  first_catcher_interfere_home 32434 non-null  int64
 50  left_on_base_home   32434 non-null  int64
 51  pitchers_used_home 32434 non-null  int64

```

```
52 individual_earned_runs_home           32434 non-null int64
53 team_earned_runs_home                32434 non-null int64
54 wild_pitches_home                  32434 non-null int64
55 balks_home                         32434 non-null int64
56 putouts_home                        32434 non-null int64
57 assists_home                        32434 non-null int64
58 errors_home                         32434 non-null int64
59 passed_balls_home                  32434 non-null int64
60 double_def_home                    32434 non-null int64
61 triple_def_home                   32434 non-null int64
dtypes: float64(3), int64(56), object(3)
memory usage: 15.6+ MB
```

Create new features

This section creates several new features: year, winning team, losing team, run differential, total runs, and on base percentage (OBP).

```
In [32]: # Creating an integer four-digit year feature
# Extract the first four characters
year_digits = df['date'].str[0:4]

# Convert the extracted substring to numeric
df['year'] = pd.to_numeric(year_digits)
df['year'] = df['year'].astype(int)
```

```
In [33]: # Creating a feature for the LOSING team for each game/row
def compare_and_get_value(row):
    if row['score_visiting'] < row['score_home']:
        return row['team_visiting']
    else:
        return row['team_home']

# Create a new feature based on the comparison
df['losing_team'] = df.apply(compare_and_get_value, axis=1)
```

```
In [34]: df['losing_team'].value_counts()
```

```
Out[34]: MIA    1204  
COL    1186  
KCA    1182  
PIT    1154  
CHA    1151  
BAL    1150  
SDN    1138  
ARI    1138  
DET    1133  
CIN    1128  
MIN    1125  
SEA    1112  
OAK    1093  
ANA    1093  
PHI    1090  
NYN    1090  
CHN    1086  
TEX    1084  
TOR    1067  
WAS    1066  
HOU    1057  
SFN    1047  
MIL    1041  
CLE    1025  
BOS    1020  
ATL    992  
TBA    981  
SLN    979  
NYA    936  
LAN    886  
Name: losing_team, dtype: int64
```

```
In [35]: # Creating a feature for the WINNING team for each game/row  
def compare_and_get_value(row):  
    if row['score_visiting'] > row['score_home']:  
        return row['team_visiting']  
    else:  
        return row['team_home']  
  
# Create a new feature based on the comparison  
df['winning_team'] = df.apply(compare_and_get_value, axis=1)
```

```
In [36]: df['winning_team'].value_counts()
```

```
Out[36]: LAN    1277  
NYA    1226  
TBA    1183  
SLN    1182  
ATL    1170  
BOS    1142  
CLE    1136  
MIL    1126  
SFN    1115  
HOU    1106  
TOR    1097  
WAS    1090  
TEX    1083  
CHN    1077  
PHI    1073  
ANA    1070  
OAK    1070  
NYN    1068  
SEA    1052  
MIN    1039  
CIN    1031  
ARI    1027  
DET    1025  
SDN    1023  
BAL    1015  
CHA    1009  
PIT    1008  
COL     979  
KCA    978  
MIA    957  
Name: winning_team, dtype: int64
```

```
In [37]: # Creating a feature for run differential for each game/row  
def calculate_run_differential(row):  
    if row['winning_team'] == row['team_home']:  
        return row['score_home'] - row['score_visiting']  
    else:  
        return row['score_visiting'] - row['score_home']  
  
# Apply the function to create a new run differential column  
df['run_differential'] = df.apply(calculate_run_differential, axis=1)
```

```
In [38]: df['run_differential'].value_counts()
```

```
Out[38]: 1      9308  
2      5916  
3      4622  
4      3817  
5      2695  
6      1974  
7      1356  
8      919  
9      647  
10     432  
11     269  
12     162  
13     139  
14     72  
15     40  
16     28  
17     14  
18     9  
19     5  
20     4  
21     4  
23     1  
24     1  
Name: run_differential, dtype: int64
```

```
In [39]: # Total Runs per Game
df['total_runs'] = df['score_home'] + df['score_visiting']
df['total_runs'].value_counts()
```

```
Out[39]: 7      3551
9      3347
5      3175
8      2555
11     2389
6      2379
10     2133
3      2065
4      1655
13     1650
12     1554
15     1068
14     1060
2      662
1      628
17     608
16     606
18     353
19     289
20     200
21     174
22     94
23     87
25     47
24     47
26     19
28     9
27     9
29     9
30     4
31     3
32     3
38     1
33     1
Name: total_runs, dtype: int64
```

```
In [40]: # OBP = (Hits + Walks + Hit by Pitch) ÷ (At Bats + Walks + Hit by Pitch + Sac
df['obp'] = (df['hits_visiting'] + df['hits_home'] + df['walk_visiting'] + df[
df['obp']].describe()
```

```
Out[40]: count    32434.000000
mean        0.313536
std         0.067768
min         0.051724
25%        0.266667
50%        0.313433
75%        0.358974
max         0.600000
Name: obp, dtype: float64
```

```
In [41]: df.head()
```

Out[41]:

	date	team_visiting	team_home	score_visiting	score_home	outs_in_game	at_bats_visit
0	20100404	NYA	BOS	7	9	51	3
1	20100405	MIN	ANA	3	6	51	3
2	20100405	CLE	CHA	0	6	51	3
3	20100405	DET	KCA	8	4	54	3
4	20100405	SEA	OAK	5	3	54	3



Dataframe manipulation

In this section we change the dataset from one row per game to two rows per game, one for each team for each game. Each paired row receives the same game ID to ensure pairs are kept together in the train_test_split to prevent data leakage.

```
In [42]: # Add unique numeric IDs
df['ID'] = range(1, len(df) + 1)
df['ID'] = df['ID'].astype(str) + '_w' # Adding suffix '_w' to original IDs

# Duplicate rows and modify IDs for duplicates
duplicates = df.copy()
duplicates['ID'] = duplicates['ID'].apply(lambda x: x.replace('_w', '_l')) #

# Concatenate original DataFrame and duplicates
result = pd.concat([df, duplicates], ignore_index=True)
```

```
In [43]: result['win'] = result['ID'].apply(lambda x: 1 if x.endswith('_w') else 0)
```

In [44]: `result.info()`

```

<class 'pandas.core.frame.DataFrame'>
RangeIndex: 64868 entries, 0 to 64867
Data columns (total 70 columns):
 #   Column           Non-Null Count Dtype
 ---  ----
 0   date             64868 non-null  object
 1   team_visiting    64868 non-null  object
 2   team_home         64868 non-null  object
 3   score_visiting   64868 non-null  int64
 4   score_home        64868 non-null  int64
 5   outs_in_game      64868 non-null  int64
 6   at_bats_visiting 64868 non-null  float64
 7   hits_visiting     64868 non-null  int64
 8   double_visiting   64868 non-null  float64
 9   triple_visiting   64868 non-null  float64
 10  home_run_visiting 64868 non-null  int64
 11  rbi_visiting      64868 non-null  int64
 12  sacrifice_hit_visiting 64868 non-null  int64
 13  sacrifice_fly_visiting 64868 non-null  int64
 14  hit_by_pitch_visiting 64868 non-null  int64
 15  walk_visiting     64868 non-null  int64
 16  intent_walk_visiting 64868 non-null  int64
 17  strikeout_visiting 64868 non-null  int64
 18  stolen_base_visiting 64868 non-null  int64
 19  caught_stealing_visiting 64868 non-null  int64
 20  grounded_into_double_plays_visiting 64868 non-null  int64
 21  first_catcher_interfere_visiting 64868 non-null  int64
 22  left_on_base_visiting 64868 non-null  int64
 23  pitchers_used_visiting 64868 non-null  int64
 24  individual_earned_runs_visiting 64868 non-null  int64
 25  team_earned_runs_visiting 64868 non-null  int64
 26  wild_pitches_visiting 64868 non-null  int64
 27  balks_visiting     64868 non-null  int64
 28  putouts_visiting   64868 non-null  int64
 29  assists_visiting   64868 non-null  int64
 30  errors_visiting    64868 non-null  int64
 31  passed_balls_visiting 64868 non-null  int64
 32  double_def_visiting 64868 non-null  int64
 33  triple_def_visiting 64868 non-null  int64
 34  at_bats_home        64868 non-null  int64
 35  hits_home          64868 non-null  int64
 36  double_home         64868 non-null  int64
 37  triple_home         64868 non-null  int64
 38  home_run_home       64868 non-null  int64
 39  rbi_home            64868 non-null  int64
 40  sacrifice_hit_home 64868 non-null  int64
 41  sacrifice_fly_home 64868 non-null  int64
 42  hit_by_pitch_home   64868 non-null  int64
 43  walk_home           64868 non-null  int64
 44  intent_walk_home    64868 non-null  int64
 45  strikeout_home       64868 non-null  int64
 46  stolen_base_home    64868 non-null  int64
 47  caught_stealing_home 64868 non-null  int64
 48  grounded_into_double_plays_home 64868 non-null  int64
 49  first_catcher_interfere_home 64868 non-null  int64
 50  left_on_base_home   64868 non-null  int64
 51  pitchers_used_home 64868 non-null  int64

```

```

52 individual_earned_runs_home      64868 non-null int64
53 team_earned_runs_home          64868 non-null int64
54 wild_pitches_home             64868 non-null int64
55 balks_home                   64868 non-null int64
56 putouts_home                  64868 non-null int64
57 assists_home                  64868 non-null int64
58 errors_home                   64868 non-null int64
59 passed_balls_home            64868 non-null int64
60 double_def_home               64868 non-null int64
61 triple_def_home              64868 non-null int64
62 year                          64868 non-null int32
63 losing_team                   64868 non-null object
64 winning_team                  64868 non-null object
65 run_differential              64868 non-null int64
66 total_runs                    64868 non-null int64
67 obp                           64868 non-null float64
68 ID                            64868 non-null object
69 win                           64868 non-null int64
dtypes: float64(4), int32(1), int64(59), object(6)
memory usage: 34.4+ MB

```

In [45]: `result_sorted = result.sort_values(by='ID')`
`result_sorted.head(6)`

Out[45]:

		date	team_visiting	team_home	score_visiting	score_home	outs_in_game	at_bats_
42433	20140422		NYA	BOS	9	3	54	
9999	20140422		NYA	BOS	9	3	54	
42434	20140422		KCA	CLE	8	2	54	
10000	20140422		KCA	CLE	8	2	54	
42435	20140422		CHA	DET	6	8	51	
10001	20140422		CHA	DET	6	8	51	



In [46]: `# Since I divided the winning team and losing team from each row and stacked them, I will now create four new DataFrames:`
`# (1) Home team wins`
`# (2) Visiting team wins`
`# (3) Home team loses`
`# (4) Visiting team loses`

```

home_win_df = result[(result['win'] == 1) & (result['winning_team'] == result['team_home'])]
visiting_win_df = result[(result['win'] == 1) & (result['winning_team'] == result['team_visiting'])]
home_lose_df = result[(result['win'] == 0) & (result['losing_team'] == result['team_home'])]
visiting_lose_df = result[(result['win'] == 0) & (result['losing_team'] == result['team_visiting'])]

```



In [47]: `home_win_df.info()`

```

<class 'pandas.core.frame.DataFrame'>
Int64Index: 17359 entries, 0 to 32432
Data columns (total 70 columns):
 #   Column           Non-Null Count Dtype  
 ---  -- 
 0   date             17359 non-null  object  
 1   team_visiting    17359 non-null  object  
 2   team_home         17359 non-null  object  
 3   score_visiting   17359 non-null  int64   
 4   score_home        17359 non-null  int64   
 5   outs_in_game     17359 non-null  int64   
 6   at_bats_visiting 17359 non-null  float64 
 7   hits_visiting    17359 non-null  int64   
 8   double_visiting   17359 non-null  float64 
 9   triple_visiting   17359 non-null  float64 
 10  home_run_visiting 17359 non-null  int64   
 11  rbi_visiting     17359 non-null  int64   
 12  sacrifice_hit_visiting 17359 non-null  int64   
 13  sacrifice_fly_visiting 17359 non-null  int64   
 14  hit_by_pitch_visiting 17359 non-null  int64   
 15  walk_visiting     17359 non-null  int64   
 16  intent_walk_visiting 17359 non-null  int64   
 17  strikeout_visiting 17359 non-null  int64   
 18  stolen_base_visiting 17359 non-null  int64   
 19  caught_stealing_visiting 17359 non-null  int64   
 20  grounded_into_double_plays_visiting 17359 non-null  int64   
 21  first_catcher_interfere_visiting 17359 non-null  int64   
 22  left_on_base_visiting 17359 non-null  int64   
 23  pitchers_used_visiting 17359 non-null  int64   
 24  individual_earned_runs_visiting 17359 non-null  int64   
 25  team_earned_runs_visiting 17359 non-null  int64   
 26  wild_pitches_visiting 17359 non-null  int64   
 27  balks_visiting     17359 non-null  int64   
 28  putouts_visiting   17359 non-null  int64   
 29  assists_visiting   17359 non-null  int64   
 30  errors_visiting    17359 non-null  int64   
 31  passed_balls_visiting 17359 non-null  int64   
 32  double_def_visiting 17359 non-null  int64   
 33  triple_def_visiting 17359 non-null  int64   
 34  at_bats_home       17359 non-null  int64   
 35  hits_home          17359 non-null  int64   
 36  double_home         17359 non-null  int64   
 37  triple_home         17359 non-null  int64   
 38  home_run_home      17359 non-null  int64   
 39  rbi_home            17359 non-null  int64   
 40  sacrifice_hit_home 17359 non-null  int64   
 41  sacrifice_fly_home 17359 non-null  int64   
 42  hit_by_pitch_home   17359 non-null  int64   
 43  walk_home           17359 non-null  int64   
 44  intent_walk_home    17359 non-null  int64   
 45  strikeout_home      17359 non-null  int64   
 46  stolen_base_home    17359 non-null  int64   
 47  caught_stealing_home 17359 non-null  int64   
 48  grounded_into_double_plays_home 17359 non-null  int64   
 49  first_catcher_interfere_home 17359 non-null  int64   
 50  left_on_base_home   17359 non-null  int64   
 51  pitchers_used_home 17359 non-null  int64

```

```
52 individual_earned_runs_home      17359 non-null int64
53 team_earned_runs_home          17359 non-null int64
54 wild_pitches_home            17359 non-null int64
55 balks_home                   17359 non-null int64
56 putouts_home                  17359 non-null int64
57 assists_home                  17359 non-null int64
58 errors_home                  17359 non-null int64
59 passed_balls_home            17359 non-null int64
60 double_def_home              17359 non-null int64
61 triple_def_home              17359 non-null int64
62 year                          17359 non-null int32
63 losing_team                  17359 non-null object
64 winning_team                 17359 non-null object
65 run_differential             17359 non-null int64
66 total_runs                   17359 non-null int64
67 obp                           17359 non-null float64
68 ID                            17359 non-null object
69 win                           17359 non-null int64
dtypes: float64(4), int32(1), int64(59), object(6)
memory usage: 9.3+ MB
```

In [48]: visiting_win_df.info()

```

<class 'pandas.core.frame.DataFrame'>
Int64Index: 15075 entries, 3 to 32433
Data columns (total 70 columns):
 #   Column           Non-Null Count Dtype
 ---  ----
 0   date             15075 non-null  object
 1   team_visiting    15075 non-null  object
 2   team_home         15075 non-null  object
 3   score_visiting   15075 non-null  int64
 4   score_home        15075 non-null  int64
 5   outs_in_game     15075 non-null  int64
 6   at_bats_visiting 15075 non-null  float64
 7   hits_visiting    15075 non-null  int64
 8   double_visiting   15075 non-null  float64
 9   triple_visiting   15075 non-null  float64
 10  home_run_visiting 15075 non-null  int64
 11  rbi_visiting     15075 non-null  int64
 12  sacrifice_hit_visiting 15075 non-null  int64
 13  sacrifice_fly_visiting 15075 non-null  int64
 14  hit_by_pitch_visiting 15075 non-null  int64
 15  walk_visiting     15075 non-null  int64
 16  intent_walk_visiting 15075 non-null  int64
 17  strikeout_visiting 15075 non-null  int64
 18  stolen_base_visiting 15075 non-null  int64
 19  caught_stealing_visiting 15075 non-null  int64
 20  grounded_into_double_plays_visiting 15075 non-null  int64
 21  first_catcher_interfere_visiting 15075 non-null  int64
 22  left_on_base_visiting 15075 non-null  int64
 23  pitchers_used_visiting 15075 non-null  int64
 24  individual_earned_runs_visiting 15075 non-null  int64
 25  team_earned_runs_visiting 15075 non-null  int64
 26  wild_pitches_visiting 15075 non-null  int64
 27  balks_visiting     15075 non-null  int64
 28  putouts_visiting   15075 non-null  int64
 29  assists_visiting   15075 non-null  int64
 30  errors_visiting   15075 non-null  int64
 31  passed_balls_visiting 15075 non-null  int64
 32  double_def_visiting 15075 non-null  int64
 33  triple_def_visiting 15075 non-null  int64
 34  at_bats_home       15075 non-null  int64
 35  hits_home          15075 non-null  int64
 36  double_home         15075 non-null  int64
 37  triple_home         15075 non-null  int64
 38  home_run_home      15075 non-null  int64
 39  rbi_home            15075 non-null  int64
 40  sacrifice_hit_home 15075 non-null  int64
 41  sacrifice_fly_home 15075 non-null  int64
 42  hit_by_pitch_home   15075 non-null  int64
 43  walk_home           15075 non-null  int64
 44  intent_walk_home    15075 non-null  int64
 45  strikeout_home      15075 non-null  int64
 46  stolen_base_home    15075 non-null  int64
 47  caught_stealing_home 15075 non-null  int64
 48  grounded_into_double_plays_home 15075 non-null  int64
 49  first_catcher_interfere_home 15075 non-null  int64
 50  left_on_base_home   15075 non-null  int64
 51  pitchers_used_home 15075 non-null  int64

```

```
52 individual_earned_runs_home      15075 non-null int64
53 team_earned_runs_home          15075 non-null int64
54 wild_pitches_home            15075 non-null int64
55 balks_home                   15075 non-null int64
56 putouts_home                  15075 non-null int64
57 assists_home                  15075 non-null int64
58 errors_home                  15075 non-null int64
59 passed_balls_home            15075 non-null int64
60 double_def_home              15075 non-null int64
61 triple_def_home              15075 non-null int64
62 year                          15075 non-null int32
63 losing_team                  15075 non-null object
64 winning_team                 15075 non-null object
65 run_differential             15075 non-null int64
66 total_runs                   15075 non-null int64
67 obp                           15075 non-null float64
68 ID                            15075 non-null object
69 win                           15075 non-null int64
dtypes: float64(4), int32(1), int64(59), object(6)
memory usage: 8.1+ MB
```

In [49]: `home_lose_df.info()`

```

<class 'pandas.core.frame.DataFrame'>
Int64Index: 15075 entries, 32437 to 64867
Data columns (total 70 columns):
 #   Column           Non-Null Count Dtype  
 --- 
 0   date             15075 non-null  object  
 1   team_visiting    15075 non-null  object  
 2   team_home         15075 non-null  object  
 3   score_visiting   15075 non-null  int64   
 4   score_home        15075 non-null  int64   
 5   outs_in_game     15075 non-null  int64   
 6   at_bats_visiting 15075 non-null  float64 
 7   hits_visiting    15075 non-null  int64   
 8   double_visiting   15075 non-null  float64 
 9   triple_visiting   15075 non-null  float64 
 10  home_run_visiting 15075 non-null  int64   
 11  rbi_visiting     15075 non-null  int64   
 12  sacrifice_hit_visiting 15075 non-null  int64   
 13  sacrifice_fly_visiting 15075 non-null  int64   
 14  hit_by_pitch_visiting 15075 non-null  int64   
 15  walk_visiting     15075 non-null  int64   
 16  intent_walk_visiting 15075 non-null  int64   
 17  strikeout_visiting 15075 non-null  int64   
 18  stolen_base_visiting 15075 non-null  int64   
 19  caught_stealing_visiting 15075 non-null  int64   
 20  grounded_into_double_plays_visiting 15075 non-null  int64   
 21  first_catcher_interfere_visiting 15075 non-null  int64   
 22  left_on_base_visiting 15075 non-null  int64   
 23  pitchers_used_visiting 15075 non-null  int64   
 24  individual_earned_runs_visiting 15075 non-null  int64   
 25  team_earned_runs_visiting 15075 non-null  int64   
 26  wild_pitches_visiting 15075 non-null  int64   
 27  balks_visiting     15075 non-null  int64   
 28  putouts_visiting   15075 non-null  int64   
 29  assists_visiting   15075 non-null  int64   
 30  errors_visiting    15075 non-null  int64   
 31  passed_balls_visiting 15075 non-null  int64   
 32  double_def_visiting 15075 non-null  int64   
 33  triple_def_visiting 15075 non-null  int64   
 34  at_bats_home        15075 non-null  int64   
 35  hits_home          15075 non-null  int64   
 36  double_home         15075 non-null  int64   
 37  triple_home         15075 non-null  int64   
 38  home_run_home       15075 non-null  int64   
 39  rbi_home            15075 non-null  int64   
 40  sacrifice_hit_home 15075 non-null  int64   
 41  sacrifice_fly_home 15075 non-null  int64   
 42  hit_by_pitch_home   15075 non-null  int64   
 43  walk_home           15075 non-null  int64   
 44  intent_walk_home    15075 non-null  int64   
 45  strikeout_home      15075 non-null  int64   
 46  stolen_base_home    15075 non-null  int64   
 47  caught_stealing_home 15075 non-null  int64   
 48  grounded_into_double_plays_home 15075 non-null  int64   
 49  first_catcher_interfere_home 15075 non-null  int64   
 50  left_on_base_home   15075 non-null  int64   
 51  pitchers_used_home 15075 non-null  int64

```

```
52 individual_earned_runs_home      15075 non-null int64
53 team_earned_runs_home          15075 non-null int64
54 wild_pitches_home            15075 non-null int64
55 balks_home                   15075 non-null int64
56 putouts_home                  15075 non-null int64
57 assists_home                  15075 non-null int64
58 errors_home                  15075 non-null int64
59 passed_balls_home            15075 non-null int64
60 double_def_home              15075 non-null int64
61 triple_def_home              15075 non-null int64
62 year                          15075 non-null int32
63 losing_team                  15075 non-null object
64 winning_team                 15075 non-null object
65 run_differential             15075 non-null int64
66 total_runs                   15075 non-null int64
67 obp                           15075 non-null float64
68 ID                            15075 non-null object
69 win                           15075 non-null int64
dtypes: float64(4), int32(1), int64(59), object(6)
memory usage: 8.1+ MB
```

In [50]: visiting_lose_df.info()

```

<class 'pandas.core.frame.DataFrame'>
Int64Index: 17359 entries, 32434 to 64866
Data columns (total 70 columns):
 #   Column           Non-Null Count Dtype  
 ---  -- 
 0   date             17359 non-null  object  
 1   team_visiting    17359 non-null  object  
 2   team_home         17359 non-null  object  
 3   score_visiting   17359 non-null  int64   
 4   score_home        17359 non-null  int64   
 5   outs_in_game     17359 non-null  int64   
 6   at_bats_visiting 17359 non-null  float64 
 7   hits_visiting    17359 non-null  int64   
 8   double_visiting   17359 non-null  float64 
 9   triple_visiting   17359 non-null  float64 
 10  home_run_visiting 17359 non-null  int64   
 11  rbi_visiting     17359 non-null  int64   
 12  sacrifice_hit_visiting 17359 non-null  int64   
 13  sacrifice_fly_visiting 17359 non-null  int64   
 14  hit_by_pitch_visiting 17359 non-null  int64   
 15  walk_visiting     17359 non-null  int64   
 16  intent_walk_visiting 17359 non-null  int64   
 17  strikeout_visiting 17359 non-null  int64   
 18  stolen_base_visiting 17359 non-null  int64   
 19  caught_stealing_visiting 17359 non-null  int64   
 20  grounded_into_double_plays_visiting 17359 non-null  int64   
 21  first_catcher_interfere_visiting 17359 non-null  int64   
 22  left_on_base_visiting 17359 non-null  int64   
 23  pitchers_used_visiting 17359 non-null  int64   
 24  individual_earned_runs_visiting 17359 non-null  int64   
 25  team_earned_runs_visiting 17359 non-null  int64   
 26  wild_pitches_visiting 17359 non-null  int64   
 27  balks_visiting     17359 non-null  int64   
 28  putouts_visiting   17359 non-null  int64   
 29  assists_visiting   17359 non-null  int64   
 30  errors_visiting    17359 non-null  int64   
 31  passed_balls_visiting 17359 non-null  int64   
 32  double_def_visiting 17359 non-null  int64   
 33  triple_def_visiting 17359 non-null  int64   
 34  at_bats_home        17359 non-null  int64   
 35  hits_home          17359 non-null  int64   
 36  double_home         17359 non-null  int64   
 37  triple_home         17359 non-null  int64   
 38  home_run_home       17359 non-null  int64   
 39  rbi_home            17359 non-null  int64   
 40  sacrifice_hit_home 17359 non-null  int64   
 41  sacrifice_fly_home 17359 non-null  int64   
 42  hit_by_pitch_home   17359 non-null  int64   
 43  walk_home           17359 non-null  int64   
 44  intent_walk_home    17359 non-null  int64   
 45  strikeout_home       17359 non-null  int64   
 46  stolen_base_home    17359 non-null  int64   
 47  caught_stealing_home 17359 non-null  int64   
 48  grounded_into_double_plays_home 17359 non-null  int64   
 49  first_catcher_interfere_home 17359 non-null  int64   
 50  left_on_base_home   17359 non-null  int64   
 51  pitchers_used_home 17359 non-null  int64

```

```
52 individual_earned_runs_home      17359 non-null int64
53 team_earned_runs_home          17359 non-null int64
54 wild_pitches_home            17359 non-null int64
55 balks_home                   17359 non-null int64
56 putouts_home                  17359 non-null int64
57 assists_home                  17359 non-null int64
58 errors_home                  17359 non-null int64
59 passed_balls_home            17359 non-null int64
60 double_def_home              17359 non-null int64
61 triple_def_home              17359 non-null int64
62 year                          17359 non-null int32
63 losing_team                  17359 non-null object
64 winning_team                 17359 non-null object
65 run_differential             17359 non-null int64
66 total_runs                   17359 non-null int64
67 obp                           17359 non-null float64
68 ID                            17359 non-null object
69 win                           17359 non-null int64
dtypes: float64(4), int32(1), int64(59), object(6)
memory usage: 9.3+ MB
```

```
In [51]: team_home_df = pd.concat([home_win_df, home_lose_df], ignore_index=True)
```

```
In [52]: team_visiting_df = pd.concat([visiting_win_df, visiting_lose_df], ignore_index=True)
```

```
In [53]: team_home_df.info()
```

```

<class 'pandas.core.frame.DataFrame'>
RangeIndex: 32434 entries, 0 to 32433
Data columns (total 70 columns):
 #   Column           Non-Null Count Dtype
 ---  ----
 0   date             32434 non-null  object
 1   team_visiting    32434 non-null  object
 2   team_home         32434 non-null  object
 3   score_visiting   32434 non-null  int64
 4   score_home        32434 non-null  int64
 5   outs_in_game      32434 non-null  int64
 6   at_bats_visiting 32434 non-null  float64
 7   hits_visiting     32434 non-null  int64
 8   double_visiting   32434 non-null  float64
 9   triple_visiting   32434 non-null  float64
 10  home_run_visiting 32434 non-null  int64
 11  rbi_visiting      32434 non-null  int64
 12  sacrifice_hit_visiting 32434 non-null  int64
 13  sacrifice_fly_visiting 32434 non-null  int64
 14  hit_by_pitch_visiting 32434 non-null  int64
 15  walk_visiting     32434 non-null  int64
 16  intent_walk_visiting 32434 non-null  int64
 17  strikeout_visiting 32434 non-null  int64
 18  stolen_base_visiting 32434 non-null  int64
 19  caught_stealing_visiting 32434 non-null  int64
 20  grounded_into_double_plays_visiting 32434 non-null  int64
 21  first_catcher_interfere_visiting 32434 non-null  int64
 22  left_on_base_visiting 32434 non-null  int64
 23  pitchers_used_visiting 32434 non-null  int64
 24  individual_earned_runs_visiting 32434 non-null  int64
 25  team_earned_runs_visiting 32434 non-null  int64
 26  wild_pitches_visiting 32434 non-null  int64
 27  balks_visiting     32434 non-null  int64
 28  putouts_visiting   32434 non-null  int64
 29  assists_visiting   32434 non-null  int64
 30  errors_visiting    32434 non-null  int64
 31  passed_balls_visiting 32434 non-null  int64
 32  double_def_visiting 32434 non-null  int64
 33  triple_def_visiting 32434 non-null  int64
 34  at_bats_home        32434 non-null  int64
 35  hits_home          32434 non-null  int64
 36  double_home         32434 non-null  int64
 37  triple_home         32434 non-null  int64
 38  home_run_home       32434 non-null  int64
 39  rbi_home            32434 non-null  int64
 40  sacrifice_hit_home 32434 non-null  int64
 41  sacrifice_fly_home 32434 non-null  int64
 42  hit_by_pitch_home   32434 non-null  int64
 43  walk_home           32434 non-null  int64
 44  intent_walk_home    32434 non-null  int64
 45  strikeout_home       32434 non-null  int64
 46  stolen_base_home    32434 non-null  int64
 47  caught_stealing_home 32434 non-null  int64
 48  grounded_into_double_plays_home 32434 non-null  int64
 49  first_catcher_interfere_home 32434 non-null  int64
 50  left_on_base_home   32434 non-null  int64
 51  pitchers_used_home 32434 non-null  int64

```

```
52 individual_earned_runs_home      32434 non-null int64
53 team_earned_runs_home          32434 non-null int64
54 wild_pitches_home            32434 non-null int64
55 balks_home                   32434 non-null int64
56 putouts_home                  32434 non-null int64
57 assists_home                  32434 non-null int64
58 errors_home                  32434 non-null int64
59 passed_balls_home            32434 non-null int64
60 double_def_home               32434 non-null int64
61 triple_def_home              32434 non-null int64
62 year                          32434 non-null int32
63 losing_team                   32434 non-null object
64 winning_team                  32434 non-null object
65 run_differential              32434 non-null int64
66 total_runs                    32434 non-null int64
67 obp                           32434 non-null float64
68 ID                            32434 non-null object
69 win                           32434 non-null int64
dtypes: float64(4), int32(1), int64(59), object(6)
memory usage: 17.2+ MB
```

```
In [54]: team_visiting_df.info()
```

```

<class 'pandas.core.frame.DataFrame'>
RangeIndex: 32434 entries, 0 to 32433
Data columns (total 70 columns):
 #   Column           Non-Null Count Dtype
 ---  ----
 0   date             32434 non-null  object
 1   team_visiting    32434 non-null  object
 2   team_home         32434 non-null  object
 3   score_visiting   32434 non-null  int64
 4   score_home        32434 non-null  int64
 5   outs_in_game      32434 non-null  int64
 6   at_bats_visiting 32434 non-null  float64
 7   hits_visiting     32434 non-null  int64
 8   double_visiting   32434 non-null  float64
 9   triple_visiting   32434 non-null  float64
 10  home_run_visiting 32434 non-null  int64
 11  rbi_visiting      32434 non-null  int64
 12  sacrifice_hit_visiting 32434 non-null  int64
 13  sacrifice_fly_visiting 32434 non-null  int64
 14  hit_by_pitch_visiting 32434 non-null  int64
 15  walk_visiting     32434 non-null  int64
 16  intent_walk_visiting 32434 non-null  int64
 17  strikeout_visiting 32434 non-null  int64
 18  stolen_base_visiting 32434 non-null  int64
 19  caught_stealing_visiting 32434 non-null  int64
 20  grounded_into_double_plays_visiting 32434 non-null  int64
 21  first_catcher_interfere_visiting 32434 non-null  int64
 22  left_on_base_visiting 32434 non-null  int64
 23  pitchers_used_visiting 32434 non-null  int64
 24  individual_earned_runs_visiting 32434 non-null  int64
 25  team_earned_runs_visiting 32434 non-null  int64
 26  wild_pitches_visiting 32434 non-null  int64
 27  balks_visiting     32434 non-null  int64
 28  putouts_visiting   32434 non-null  int64
 29  assists_visiting   32434 non-null  int64
 30  errors_visiting    32434 non-null  int64
 31  passed_balls_visiting 32434 non-null  int64
 32  double_def_visiting 32434 non-null  int64
 33  triple_def_visiting 32434 non-null  int64
 34  at_bats_home        32434 non-null  int64
 35  hits_home          32434 non-null  int64
 36  double_home         32434 non-null  int64
 37  triple_home         32434 non-null  int64
 38  home_run_home       32434 non-null  int64
 39  rbi_home            32434 non-null  int64
 40  sacrifice_hit_home 32434 non-null  int64
 41  sacrifice_fly_home 32434 non-null  int64
 42  hit_by_pitch_home   32434 non-null  int64
 43  walk_home           32434 non-null  int64
 44  intent_walk_home    32434 non-null  int64
 45  strikeout_home       32434 non-null  int64
 46  stolen_base_home    32434 non-null  int64
 47  caught_stealing_home 32434 non-null  int64
 48  grounded_into_double_plays_home 32434 non-null  int64
 49  first_catcher_interfere_home 32434 non-null  int64
 50  left_on_base_home   32434 non-null  int64
 51  pitchers_used_home 32434 non-null  int64

```

```
52 individual_earned_runs_home      32434 non-null int64
53 team_earned_runs_home           32434 non-null int64
54 wild_pitches_home              32434 non-null int64
55 balks_home                     32434 non-null int64
56 putouts_home                   32434 non-null int64
57 assists_home                   32434 non-null int64
58 errors_home                    32434 non-null int64
59 passed_balls_home              32434 non-null int64
60 double_def_home                32434 non-null int64
61 triple_def_home                32434 non-null int64
62 year                           32434 non-null int32
63 losing_team                    32434 non-null object
64 winning_team                   32434 non-null object
65 run_differential               32434 non-null int64
66 total_runs                     32434 non-null int64
67 obp                            32434 non-null float64
68 ID                             32434 non-null object
69 win                            32434 non-null int64
dtypes: float64(4), int32(1), int64(59), object(6)
memory usage: 17.2+ MB
```

In [55]: *# When I stacked the dataframes on top of one another they all had to have the # The _home and _visiting were removed from the features.*

```
team_home_df.drop(columns=["at_bats_visiting", "hits_visiting", "double_visiting",
                           "sacrifice_hit_visiting", "sacrifice_fly_visiting", "hit_by_pitch_visiting",
                           "strikeout_visiting", "stolen_base_visiting", "caught_stealing_visiting",
                           "first_catcher_interfere_visiting", "left_on_base_visiting", "pitchers_used_visiting",
                           "team_earned_runs_visiting", "wild_pitches_visiting", "balks_visiting",
                           "errors_visiting", "passed_balls_visiting", "double_def_visiting",
                           "triple_def_visiting"])

team_visiting_df.drop(columns=["at_bats_home", "hits_home", "double_home", "triple_home",
                               "sacrifice_hit_home", "sacrifice_fly_home", "hit_by_pitch_home",
                               "strikeout_home", "stolen_base_home", "caught_stealing_home", "grounded_out_home",
                               "first_catcher_interfere_home", "left_on_base_home", "pitchers_used_home",
                               "team_earned_runs_home", "wild_pitches_home", "balks_home", "putouts_home",
                               "errors_home", "passed_balls_home", "double_def_home", "triple_def_home"])
```

```
In [56]: rename_dict_home = {  
  
    'at_bats_home': 'at_bats',  
    'hits_home': 'hits',  
    'double_home': 'double',  
    'triple_home': 'triple',  
    'home_run_home': 'home_run',  
    'rbi_home': 'rbi',  
    'sacrifice_hit_home': 'sacrifice_hit',  
    'sacrifice_fly_home': 'sacrifice_fly',  
    'hit_by_pitch_home': 'hit_by_pitch',  
    'walk_home': 'walk',  
    'intent_walk_home': 'intent_walk',  
    'strikeout_home': "strikeout",  
    'stolen_base_home': 'stolen_base',  
    'caught_stealing_home': 'caught_stealing',  
    'grounded_into_double_plays_home': 'grounded_into_double_plays',  
    'first_catcher_interfere_home': 'first_catcher_interfere',  
    'left_on_base_home': 'left_on_base',  
    'pitchers_used_home': 'pitchers_used',  
    'individual_earned_runs_home': 'individual_earned_runs',  
    'team_earned_runs_home': 'team_earned_runs',  
    'wild_pitches_home': 'wild_pitches',  
    'balks_home': 'balks',  
    'putouts_home': 'putouts',  
    'assists_home': 'assists',  
    'errors_home': 'errors',  
    'passed_balls_home': 'passed_balls',  
    'double_def_home': 'double_def',  
    'triple_def_home': 'triple_def'  
}  
  
# Rename the columns using the rename() function  
team_home_df.rename(columns=rename_dict_home, inplace=True)
```

```
In [57]: rename_dict_visiting = {

    "at_bats_visiting": "at_bats",
    "hits_visiting": "hits",
    "double_visiting": "double",
    "triple_visiting": "triple",
    "home_run_visiting": "home_run",
    "rbi_visiting": "rbi",
    "sacrifice_hit_visiting": "sacrifice_hit",
    'sacrifice_fly_visiting': "sacrifice_fly",
    'hit_by_pitch_visiting': "hit_by_pitch",
    'walk_visiting': 'walk',
    'intent_walk_visiting': 'intent_walk',
    'strikeout_visiting': 'strikeout',
    'stolen_base_visiting': 'stolen_base',
    'caught_stealing_visiting': 'caught_stealing',
    'grounded_into_double_plays_visiting': 'grounded_into_double_plays',
    'first_catcher_interfere_visiting': 'first_catcher_interfere',
    'left_on_base_visiting': 'left_on_base',
    'pitchers_used_visiting': "pitchers_used",
    'individual_earned_runs_visiting': "individual_earned_runs",
    'team_earned_runs_visiting': 'team_earned_runs',
    'wild_pitches_visiting': 'wild_pitches',
    'balks_visiting': 'balks',
    'putouts_visiting': 'putouts',
    'assists_visiting': 'assists',
    'errors_visiting': 'errors',
    'passed_balls_visiting': 'passed_balls',
    'double_def_visiting': 'double_def',
    'triple_def_visiting': 'triple_def'
}

# Rename the columns using the rename() function
team_visiting_df.rename(columns=rename_dict_visiting, inplace=True)
```

```
In [58]: combo_df = pd.concat([team_home_df, team_visiting_df], ignore_index=True)
```

```
In [59]: # this is so that when I train test split I can do an exact match on ID to ens
combo_df['ID'] = combo_df['ID'].str[:-2]
```

```
In [60]: combo_df.info()
```

```
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 64868 entries, 0 to 64867
Data columns (total 41 columns):
 #   Column           Non-Null Count  Dtype  
--- 
 0   date             64868 non-null   object  
 1   team_visiting    64868 non-null   object  
 2   team_home         64868 non-null   object  
 3   score_visiting   64868 non-null   int64  
 4   score_home        64868 non-null   int64  
 5   outs_in_game     64868 non-null   int64  
 6   at_bats          64868 non-null   float64 
 7   hits             64868 non-null   int64  
 8   double            64868 non-null   float64 
 9   triple            64868 non-null   float64 
 10  home_run          64868 non-null   int64  
 11  rbi              64868 non-null   int64  
 12  sacrifice_hit    64868 non-null   int64  
 13  sacrifice_fly    64868 non-null   int64  
 14  hit_by_pitch     64868 non-null   int64  
 15  walk              64868 non-null   int64  
 16  intent_walk       64868 non-null   int64  
 17  strikeout         64868 non-null   int64  
 18  stolen_base       64868 non-null   int64  
 19  caught_stealing   64868 non-null   int64  
 20  grounded_into_double_plays 64868 non-null   int64  
 21  first_catcher_interfere 64868 non-null   int64  
 22  left_on_base      64868 non-null   int64  
 23  pitchers_used     64868 non-null   int64  
 24  individual_earned_runs 64868 non-null   int64  
 25  team_earned_runs   64868 non-null   int64  
 26  wild_pitches      64868 non-null   int64  
 27  balks             64868 non-null   int64  
 28  putouts            64868 non-null   int64  
 29  assists            64868 non-null   int64  
 30  errors             64868 non-null   int64  
 31  passed_balls      64868 non-null   int64  
 32  double_def         64868 non-null   int64  
 33  triple_def         64868 non-null   int64  
 34  year               64868 non-null   int32  
 35  losing_team        64868 non-null   object  
 36  winning_team       64868 non-null   object  
 37  run_differential   64868 non-null   int64  
 38  obp                64868 non-null   float64 
 39  ID                 64868 non-null   object  
 40  win                64868 non-null   int64  
dtypes: float64(4), int32(1), int64(30), object(6)
memory usage: 20.0+ MB
```

```
In [61]: combo_df.head()
```

Out[61]:

	date	team_visiting	team_home	score_visiting	score_home	outs_in_game	at_bats	hits
0	20100404	NYA	BOS	7	9	51	34.0	12
1	20100405	MIN	ANA	3	6	51	33.0	8
2	20100405	CLE	CHA	0	6	51	28.0	6
3	20100405	TOR	TEX	4	5	52	31.0	6
4	20100405	SDN	ARI	3	6	51	33.0	8

Exploratory Data Analysis

This section runs through a set of features looking at counts, descriptive statistics, crosstabs, and histograms.

```
In [62]: # This function runs through a list of features of interest and provides the value counts for each
def get_value_counts(combo_df, features):
    value_counts = {}
    for feature in features:
        value_counts[feature] = combo_df[feature].value_counts()
    return value_counts

# List of features to get value counts for
features_to_check = ['at_bats',
    'hits',
    'double',
    'triple',
    'home_run',
    'sacrifice_hit',
    'sacrifice_fly',
    'hit_by_pitch',
    'walk',
    'intent_walk',
    'strikeout',
    'stolen_base',
    'caught_stealing',
    'grounded_into_double_plays',
    'first_catcher_interfere',
    'left_on_base',
    'pitchers_used',
    'wild_pitches',
    'balks',
    'assists',
    'errors',
    'passed_balls',
    'double_def',
    'triple_def'
]
# Get the value counts
value_counts = get_value_counts(combo_df, features_to_check)

# Print the value counts
for feature, counts in value_counts.items():
    print(f"\nValue counts for '{feature}':")
    print(counts)
```

Value counts for 'at_bats':

33.0	7448
32.0	7213
34.0	6961
31.0	6336
35.0	5872
30.0	4997
36.0	4779
37.0	3808
29.0	3345
38.0	2798
39.0	2045
28.0	1946
40.0	1500
41.0	1025
27.0	922
42.0	774
43.0	582
~ ~	~ ~

```
In [63]: # This function runs through a list of features of interest and provides the descriptive statistics for each
def get_descriptive_statistics(df, features):
    descriptive_stats = {}
    for feature in features:
        descriptive_stats[feature] = df[feature].describe()
    return descriptive_stats

# List of features to get descriptive statistics for
features_to_check = ['at_bats',
                     'hits',
                     'double',
                     'triple',
                     'home_run',
                     'sacrifice_hit',
                     'sacrifice_fly',
                     'hit_by_pitch',
                     'walk',
                     'intent_walk',
                     'strikeout',
                     'stolen_base',
                     'caught_stealing',
                     'grounded_into_double_plays',
                     'first_catcher_interfere',
                     'left_on_base',
                     'pitchers_used',
                     'wild_pitches',
                     'balks',
                     'assists',
                     'errors',
                     'passed_balls',
                     'double_def',
                     'triple_def',
                     'obp']

# Get the descriptive statistics
descriptive_stats = get_descriptive_statistics(combo_df, features_to_check)

# Print the descriptive statistics
for feature, stats in descriptive_stats.items():
    print(f"\nDescriptive statistics for '{feature}':")
    print(stats)
```

```
Descriptive statistics for 'at_bats':  
count    64868.000000  
mean     33.946353  
std      4.271191  
min     16.000000  
25%     31.000000  
50%     33.000000  
75%     36.000000  
max     76.000000  
Name: at_bats, dtype: float64
```

```
Descriptive statistics for 'hits':  
count    64868.000000  
mean     8.537430  
std      3.431466  
min     0.000000  
25%     6.000000  
50%     8.000000  
75%     11.000000
```

In [64]: `combo_df['run_differential'].describe()`

Out[64]:

```
count    64868.000000  
mean     3.454770  
std      2.667129  
min     1.000000  
25%     1.000000  
50%     3.000000  
75%     5.000000  
max     24.000000  
Name: run_differential, dtype: float64
```

In [65]: `#combo_df[''].describe()`

In [66]: `#combo_df['hits'].describe()`

In [67]: `#combo_df['double'].describe()`

In [68]: `#combo_df['triple'].describe()`

In [69]: `#combo_df['home_run'].describe()`

In [70]: `#combo_df['rbi'].describe()`

In [71]: `#combo_df['walk'].describe()`

```
In [72]: #combo_df['strikeout'].describe()
```

```
In [73]: #combo_df['stolen_base'].describe()
```

```
In [74]: cross_tab = pd.crosstab(combo_df['run_differential'], combo_df['year'])
print(cross_tab)
```

run_differential	year	2010	2011	2012	2013	2014	2015	2016	2017	2018	2019
	2020	2021	2022	2023							
1				1458	1502	1394	1506	1472	1410	1368	1292
494	494	1352	1400	1346	894	878	912	908	898	940	886
2					846					846	876
290	290	886	904		596	708	692	662	694	696	728
3					716					686	728
262	262	694	688			570	590	638	554	544	536
4					620					552	570
216	216	564	564			400	362	398	416	434	414
5					352					404	418
176	176	398	362			406				430	
6						264					288
108	108	288	314			168	216	164	194		
7						192					236
70	70	220	198				118	110	122	132	146
8						154				136	162
42	42	136	134			122				176	
9						90	100	70	90	82	104
44	44	102	108			124				90	110
10						62	56	46	66	58	56
20	20	72	76			70				70	72
11						40	30	28	34	32	48
22	22	34	30			50				48	52
12						34	24	26	16	16	14
16	16	34	28			22				14	38
13						12	14	8	18	8	20
14	14	22	28			36				26	22
14						10	12	0	2	6	16
6	6	20	10			12				6	12
15						4	2	6	8	2	10
4	4	8	2			10				0	10
16						0	4	4	4	2	6
0	0	2	6			2				4	10
17						2	0	2	0	2	2
2	2	4	2			4				4	2
18						2	2	2	0	0	0
2	2	2	0			2				4	2
19						0	0	0	0	0	0
2	2	2	0			2				2	0
20						0	0	0	0	0	2
2	2	0	0			0				0	2
21						0	0	0	0	0	2
0	0	2	0			0				0	0
23						0	0	0	0	0	0
0	0	0	2			0				0	0
24						0	0	0	0	0	0
0	0	0	0			2				0	0

```
In [75]: cross_tab2 = pd.crosstab(combo_df['strikeout'], combo_df['year'])
print(cross_tab2)
```



```
26      0      0      0      0      0      0      0      0      2      0      1      0
0      0      0
```

```
In [76]: #grouped_df = df.groupby('').sum()
#grouped_df.head()
```

```
In [77]: #grouped_df = df.groupby(['', 'year']).sum()
#grouped_df.head(20)
```

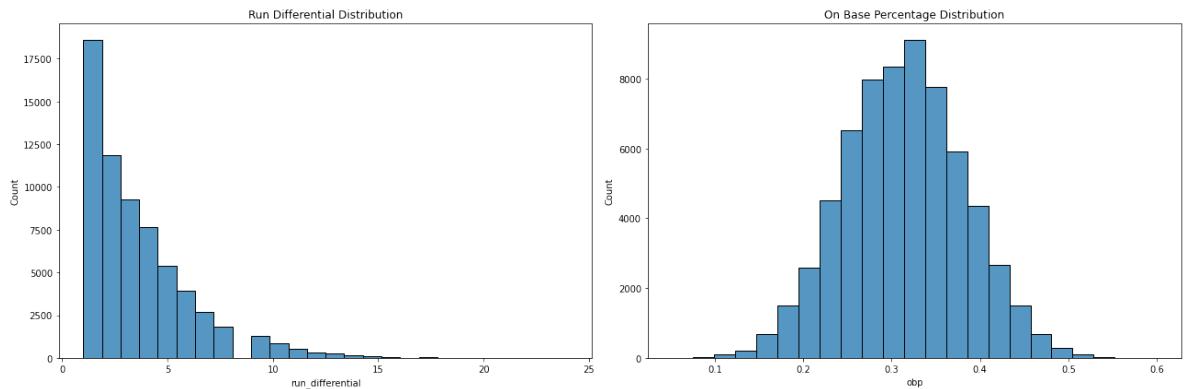
```
In [78]: fig, axes = plt.subplots(1, 2, figsize=(18, 6))

# Plot 1:
sns.histplot(combo_df['run_differential'], bins=26, ax=axes[0])
axes[0].set_title('Run Differential Distribution')

# Plot 2:
sns.histplot(combo_df['obp'], bins=23, ax=axes[1])
axes[1].set_title('On Base Percentage Distribution')

# Adjust Layout
plt.tight_layout()

# Display the plots
plt.show()
```



```
In [79]: fig, axes = plt.subplots(1, 3, figsize=(18, 6))

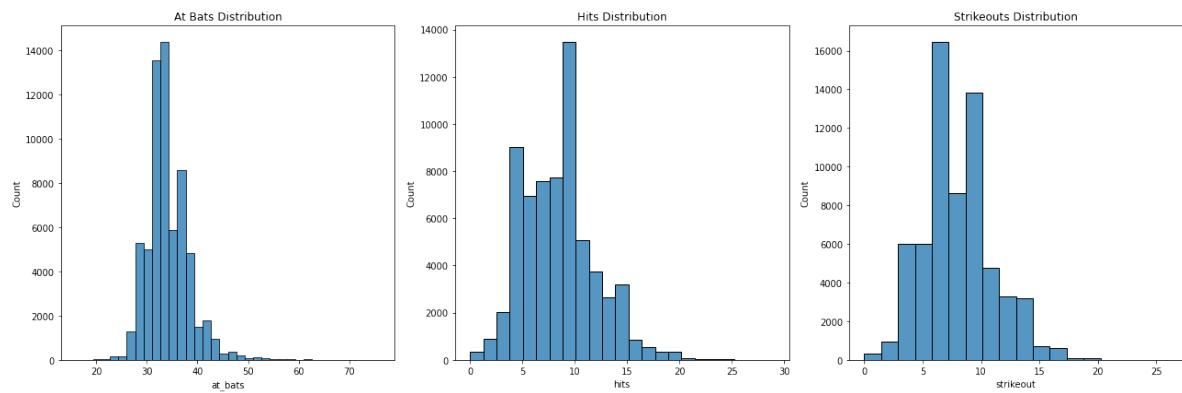
# Plot 1:
sns.histplot(combo_df['at_bats'], bins=36, ax=axes[0])
axes[0].set_title('At Bats Distribution')

# Plot 2:
sns.histplot(combo_df['hits'], bins=23, ax=axes[1])
axes[1].set_title('Hits Distribution')

# Plot 3:
sns.histplot(combo_df['strikeout'], bins=18, ax=axes[2])
axes[2].set_title('Strikeouts Distribution')

# Adjust Layout
plt.tight_layout()

# Display the plots
plt.show()
```



```
In [80]: fig, axes = plt.subplots(1, 3, figsize=(18, 6))

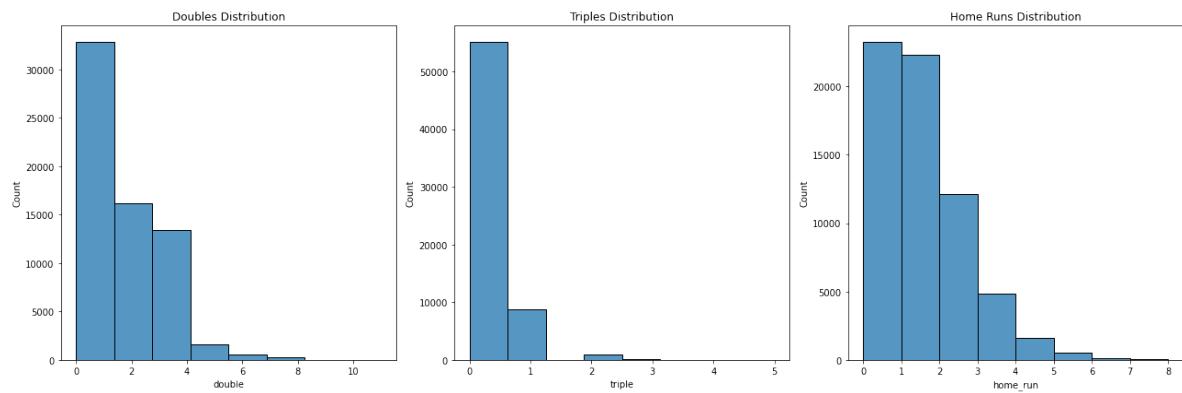
# Plot 1:
sns.histplot(combo_df['double'], bins=8, ax=axes[0])
axes[0].set_title('Doubles Distribution')

# Plot 2:
sns.histplot(combo_df['triple'], bins=8, ax=axes[1])
axes[1].set_title('Triples Distribution')

# Plot 3:
sns.histplot(combo_df['home_run'], bins=8, ax=axes[2])
axes[2].set_title('Home Runs Distribution')

# Adjust Layout
plt.tight_layout()

# Display the plots
plt.show()
```



```
In [81]: fig, axes = plt.subplots(1, 3, figsize=(18, 6))

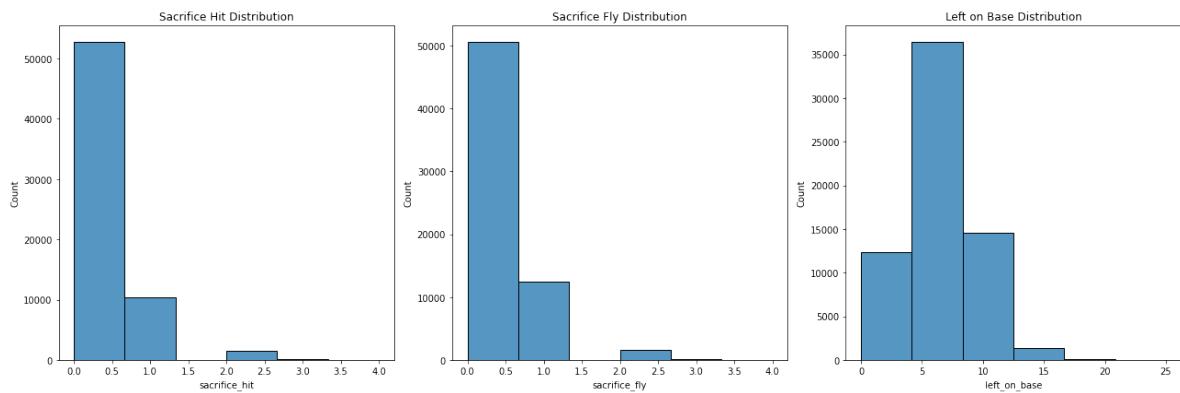
# Plot 1:
sns.histplot(combo_df['sacrifice_hit'], bins=6, ax=axes[0])
axes[0].set_title('Sacrifice Hit Distribution')

# Plot 2:
sns.histplot(combo_df['sacrifice_fly'], bins=6, ax=axes[1])
axes[1].set_title('Sacrifice Fly Distribution')

# Plot 3:
sns.histplot(combo_df['left_on_base'], bins=6, ax=axes[2])
axes[2].set_title('Left on Base Distribution')

# Adjust Layout
plt.tight_layout()

# Display the plots
plt.show()
```



```
In [82]: fig, axes = plt.subplots(1, 4, figsize=(18, 4))

# Plot 1:
sns.histplot(combo_df['hit_by_pitch'], bins=6, ax=axes[0])
axes[0].set_title('Hit By Pitch Distribution')

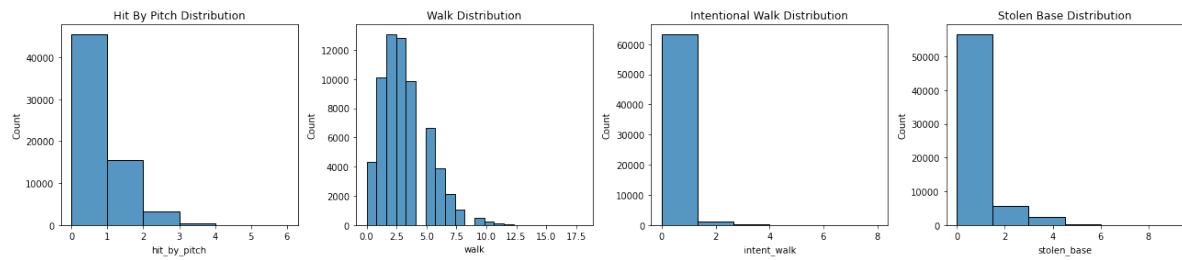
# Plot 2:
sns.histplot(combo_df['walk'], bins=22, ax=axes[1])
axes[1].set_title('Walk Distribution')

# Plot 3:
sns.histplot(combo_df['intent_walk'], bins=6, ax=axes[2])
axes[2].set_title('Intentional Walk Distribution')

# Plot 4:
sns.histplot(combo_df['stolen_base'], bins=6, ax=axes[3])
axes[3].set_title('Stolen Base Distribution')

# Adjust layout
plt.tight_layout()

# Display the plots
plt.show()
```



```
In [83]: fig, axes = plt.subplots(1, 3, figsize=(18, 6))

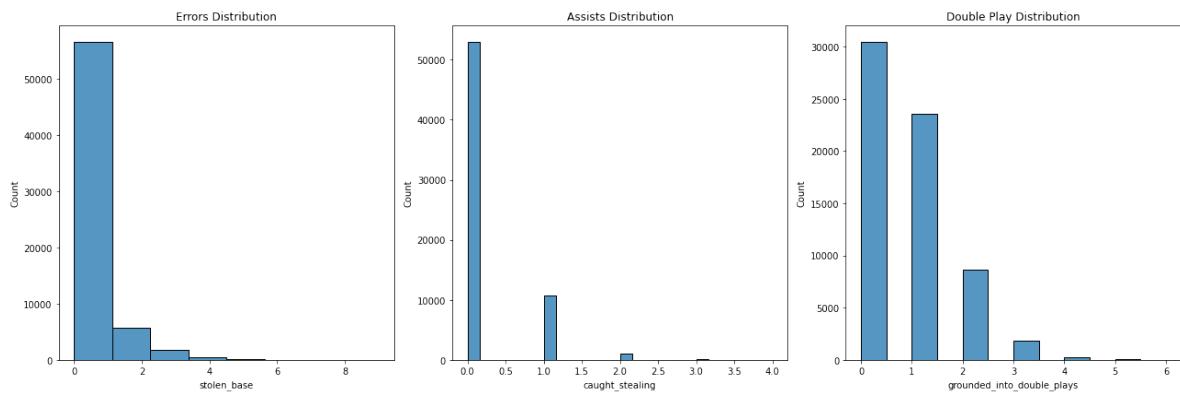
# Plot 1:
sns.histplot(combo_df['stolen_base'], bins=8, ax=axes[0])
axes[0].set_title('Errors Distribution')

# Plot 2:
sns.histplot(combo_df['caught_stealing'], bins=24, ax=axes[1])
axes[1].set_title('Assists Distribution')

# Plot 3:
sns.histplot(combo_df['grounded_into_double_plays'], bins=12, ax=axes[2])
axes[2].set_title('Double Play Distribution')

# Adjust Layout
plt.tight_layout()

# Display the plots
plt.show()
```



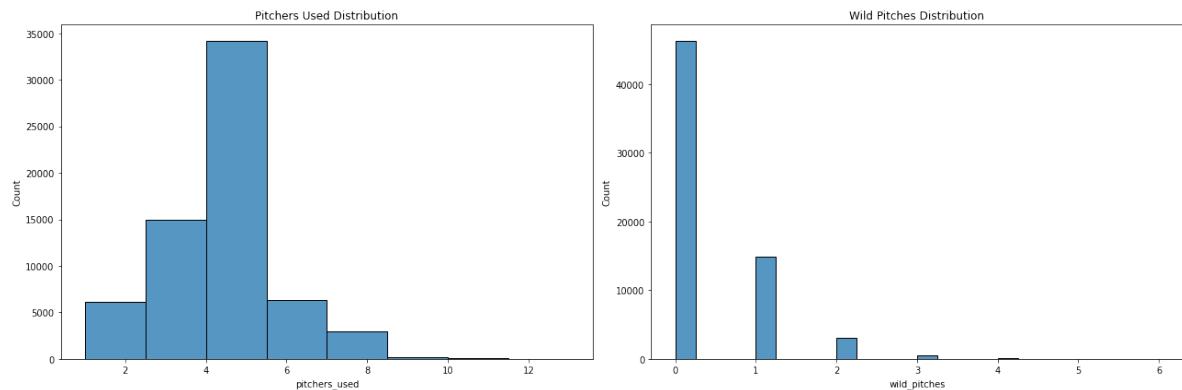
```
In [84]: fig, axes = plt.subplots(1, 2, figsize=(18, 6))

# Plot 1:
sns.histplot(combo_df['pitchers_used'], bins=8, ax=axes[0])
axes[0].set_title('Pitchers Used Distribution')

# Plot 2:
sns.histplot(combo_df['wild_pitches'], bins=24, ax=axes[1])
axes[1].set_title('Wild Pitches Distribution')

# Adjust layout
plt.tight_layout()

# Display the plots
plt.show()
```



```
In [85]: fig, axes = plt.subplots(1, 3, figsize=(18, 6))

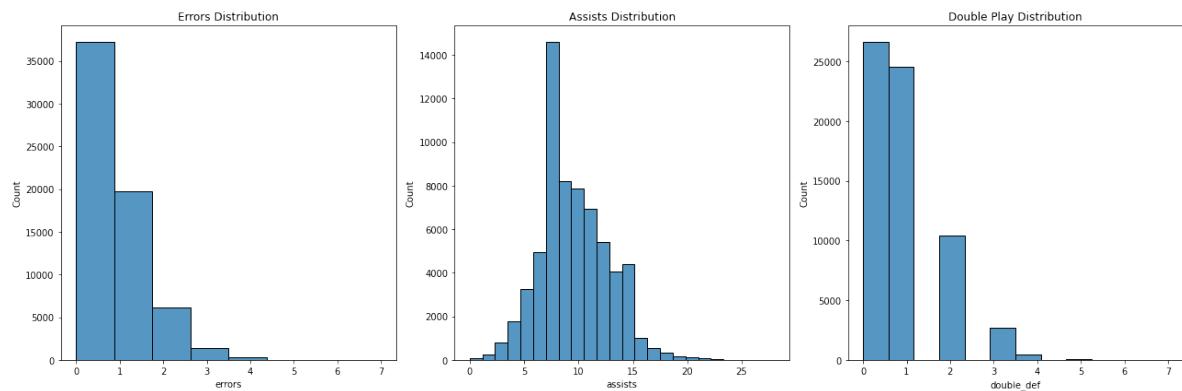
# Plot 1:
sns.histplot(combo_df['errors'], bins=8, ax=axes[0])
axes[0].set_title('Errors Distribution')

# Plot 2:
sns.histplot(combo_df['assists'], bins=24, ax=axes[1])
axes[1].set_title('Assists Distribution')

# Plot 3:
sns.histplot(combo_df['double_def'], bins=12, ax=axes[2])
axes[2].set_title('Double Play Distribution')

# Adjust Layout
plt.tight_layout()

# Display the plots
plt.show()
```



```
In [86]: fig, axes = plt.subplots(1, 2, figsize=(18, 6))

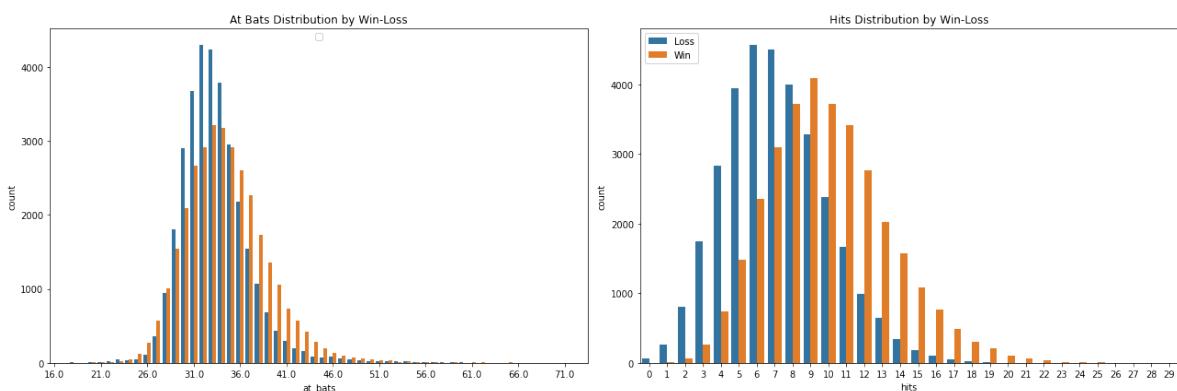
# Plot 1:
sns.countplot(x='at_bats', hue='win', data= combo_df, ax=axes[0])
axes[0].set_title('At Bats Distribution by Win-Loss')
handles, labels = axes[1].get_legend_handles_labels()
axes[0].legend(handles, ['Loss', 'Win'])

ticks = axes[0].get_xticks()
axes[0].set_xticks(ticks[::5])

# Plot 2:
sns.countplot(x='hits', hue='win', data= combo_df, ax=axes[1])
axes[1].set_title('Hits Distribution by Win-Loss')
handles, labels = axes[1].get_legend_handles_labels()
axes[1].legend(handles, ['Loss', 'Win'])

# Adjust Layout
plt.tight_layout()

# Display the plots
plt.show()
```



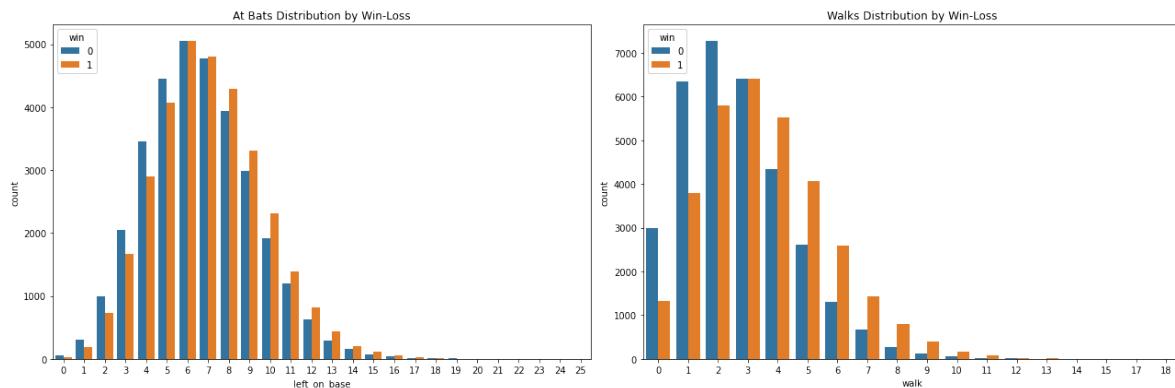
```
In [87]: fig, axes = plt.subplots(1, 2, figsize=(18, 6))

# Plot 1:
sns.countplot(x='left_on_base', hue='win', data=combo_df, ax=axes[0])
axes[0].set_title('At Bats Distribution by Win-Loss')

# Plot 2:
sns.countplot(x='walk', hue='win', data=combo_df, ax=axes[1])
axes[1].set_title('Walks Distribution by Win-Loss')

# Adjust layout
plt.tight_layout()

# Display the plots
plt.show()
```

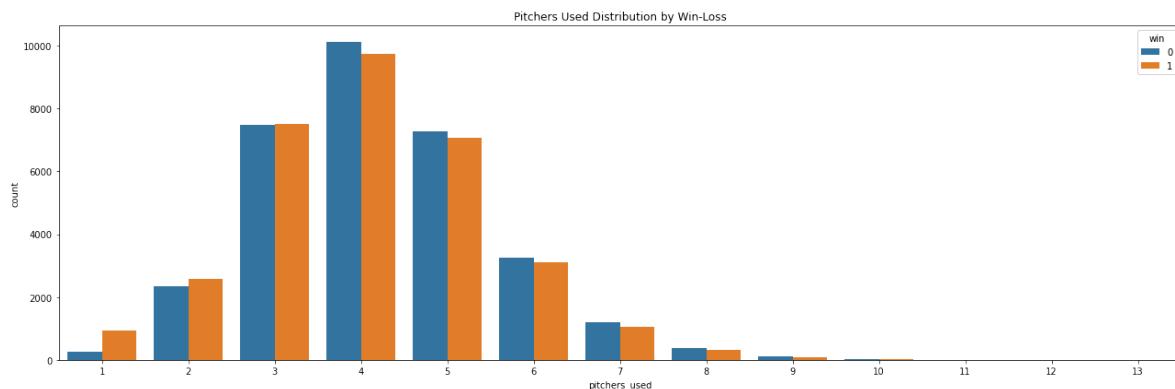


```
In [88]: fig, ax = plt.subplots(1, figsize=(18, 6))

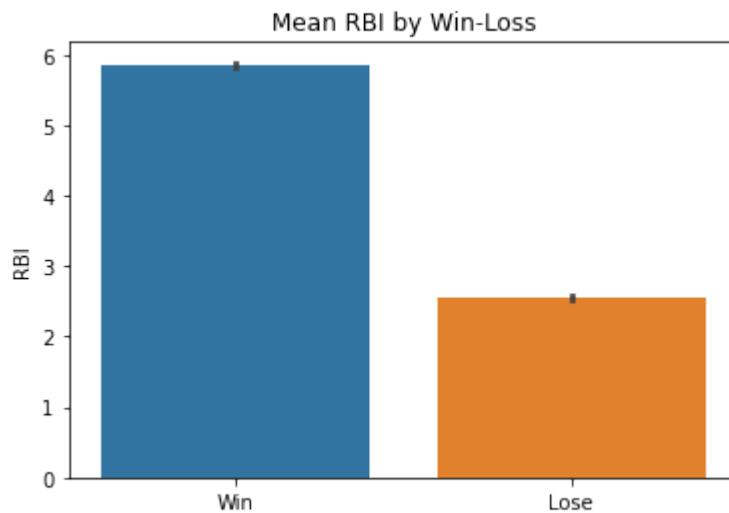
# Plot 1:
sns.countplot(x='pitchers_used', hue='win', data=combo_df, ax=ax)
ax.set_title('Pitchers Used Distribution by Win-Loss')

# Adjust layout
plt.tight_layout()

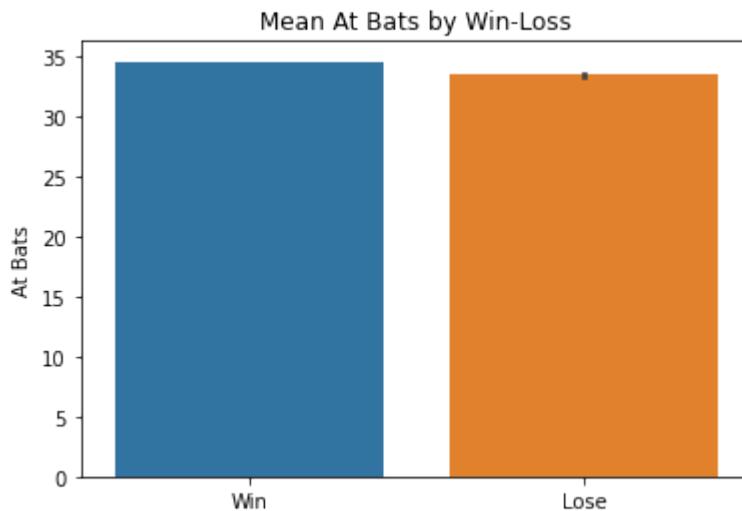
# Display the plots
plt.show()
```



```
In [89]: combo_df['win'] = combo_df['win'].map({0: 'Lose', 1: 'Win'})  
sns.barplot(x='win', y='rbi', data=combo_df)  
  
# Add Labels and title  
plt.xlabel('')  
plt.ylabel('RBI')  
plt.title('Mean RBI by Win-Loss')  
  
# Show plot  
plt.show()
```



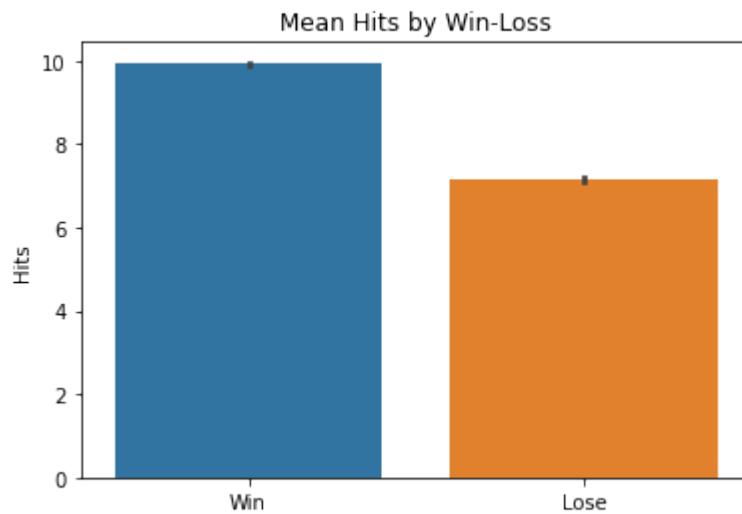
```
In [90]: sns.barplot(x='win', y='at_bats', data=combo_df)  
  
# Add Labels and title  
plt.xlabel('')  
plt.ylabel('At Bats')  
plt.title('Mean At Bats by Win-Loss')  
  
# Show plot  
plt.show();
```



```
In [91]: sns.barplot(x='win', y='hits', data=combo_df)

# Add Labels and title
plt.xlabel('')
plt.ylabel('Hits')
plt.title('Mean Hits by Win-Loss')

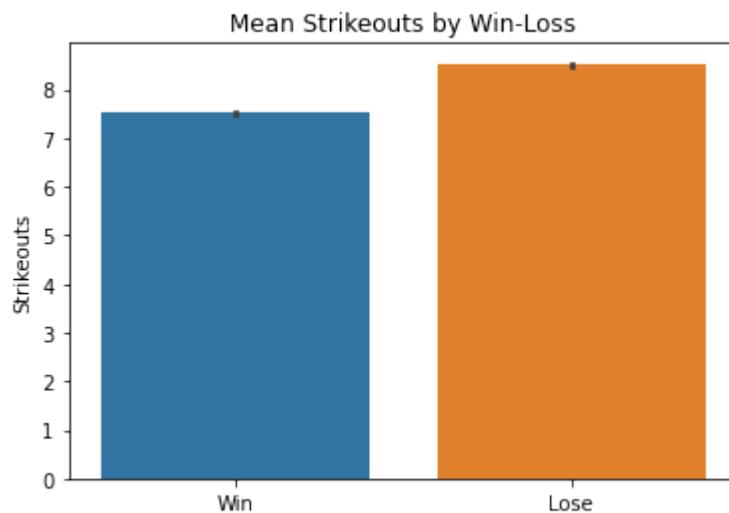
# Show plot
plt.show()
```



```
In [92]: sns.barplot(x='win', y='strikeout', data=combo_df)

# Add Labels and title
plt.xlabel('')
plt.ylabel('Strikeouts')
plt.title('Mean Strikeouts by Win-Loss')

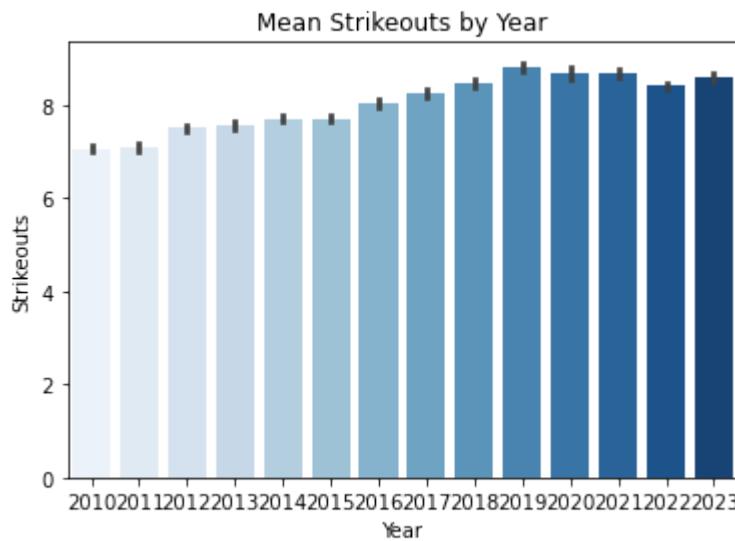
# Show plot
plt.show()
```



```
In [93]: sns.barplot(x='year', y='strikeout', data=combo_df, palette='Blues')
```

```
# Add Labels and title
plt.xlabel('Year')
plt.ylabel('Strikeouts')
plt.title('Mean Strikeouts by Year')

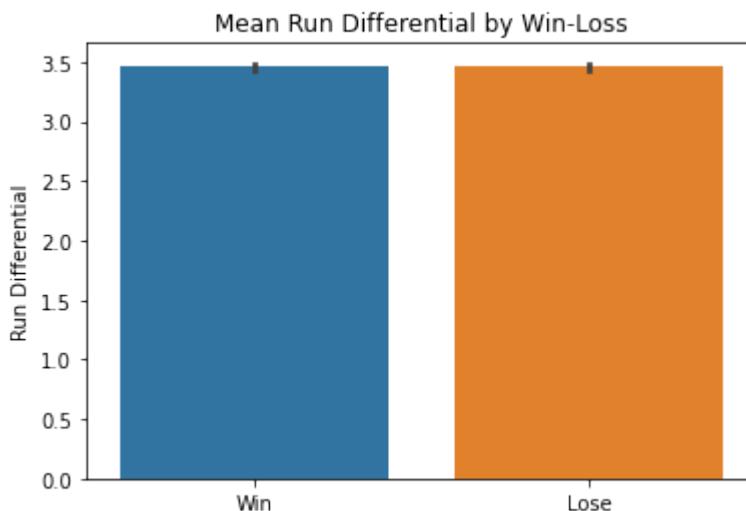
# Show plot
plt.show()
```



```
In [94]: sns.barplot(x='win', y='run_differential', data=combo_df)
```

```
# Add Labels and title
plt.xlabel('')
plt.ylabel('Run Differential')
plt.title('Mean Run Differential by Win-Loss')

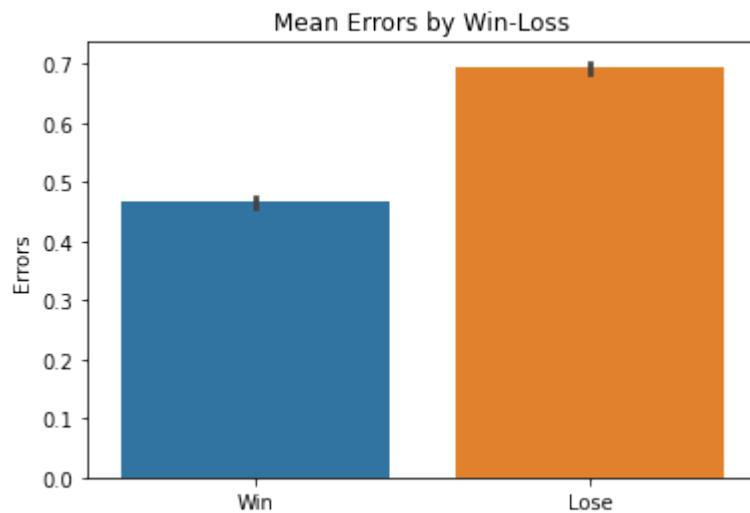
# Show plot
plt.show()
```



```
In [95]: # Runs by Year
sns.barplot(x='win', y='errors', data=combo_df)

# Add Labels and title
plt.xlabel('')
plt.ylabel('Errors')
plt.title('Mean Errors by Win-Loss')

# Show plot
plt.show()
```



```
In [96]: # Calculate the mean 'score_home' for each 'team_home'
mean_score_home = combo_df.groupby('team_home')['score_home'].mean().reset_index()

# Sort the data by mean 'score_home'
mean_score_home_sorted = mean_score_home.sort_values(by='score_home', ascending=False)

# Calculate the mean 'score_visiting' for each 'team_visiting'
mean_score_visiting = combo_df.groupby('team_visiting')['score_visiting'].mean()

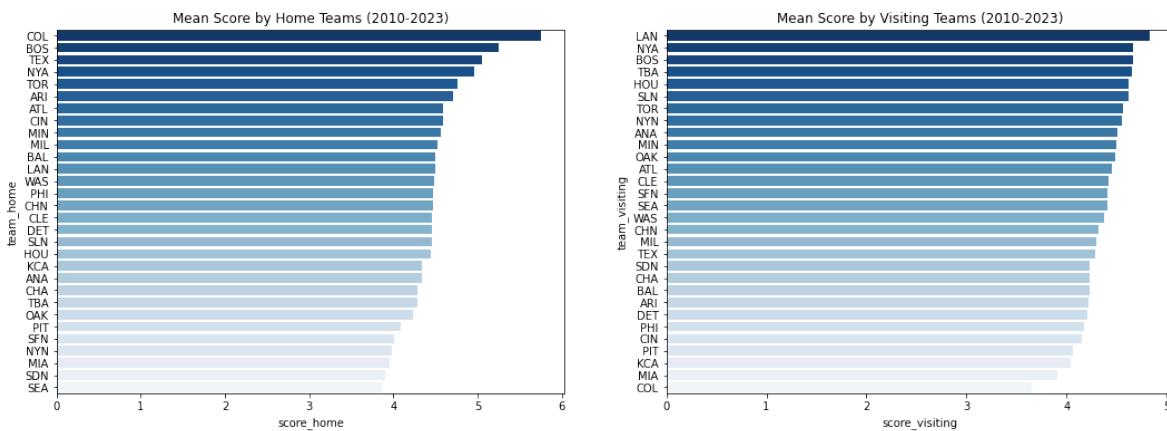
# Sort the data by mean 'score_visiting'
mean_score_visiting_sorted = mean_score_visiting.sort_values(by='score_visiting', ascending=False)

# Create a figure with 1 row and 2 columns of subplots
fig, axes = plt.subplots(1, 2, figsize=(18, 6))

# Plot 1: Bar plot of mean 'score_home' by 'team_home'
sns.barplot(x='score_home', y='team_home', data=mean_score_home_sorted, ax=axes[0])
axes[0].set_title('Mean Score by Home Teams (2010-2023)')

# Plot 2: Bar plot of mean 'score_visiting' by 'team_visiting'
sns.barplot(x='score_visiting', y='team_visiting', data=mean_score_visiting_sorted, ax=axes[1])
axes[1].set_title('Mean Score by Visiting Teams (2010-2023)')

# Show plot
plt.show()
```



The above demonstrates park factor--

```
In [97]: combo_df.head()
```

Out[97]:

	date	team_visiting	team_home	score_visiting	score_home	outs_in_game	at_bats	hits
0	20100404	NYA	BOS	7	9	51	34.0	12
1	20100405	MIN	ANA	3	6	51	33.0	8
2	20100405	CLE	CHA	0	6	51	28.0	6
3	20100405	TOR	TEX	4	5	52	31.0	6
4	20100405	SDN	ARI	3	6	51	33.0	8



```
In [98]: combo_df['win'] = combo_df['win'].map({'Lose':0, 'Win': 1})
```

```
In [99]: combo_df.to_csv('data/Combo_DF.csv')
```

Modeling

Data are split into a training set and test set for modeling. Games are paired so that features for the winning team and losing team of a game are not split between the training and test set to prevent data leakage. Features are standard scale standardized. No additional preprocessing is used. There are no categorical features, no missing data, and the target is 50:50 balanced.

The primary evaluation metric used for these models is accuracy. With a balanced dataset and equal interest in classification of wins as much as losses, accuracy is an appropriate evaluation metric.

X features varied by model and included: 'at_bats', 'hits', 'double', 'triple', 'home_run', 'sacrifice_hit', 'sacrifice_fly', 'hit_by_pitch', 'walk', 'intent_walk', 'strikeout', 'stolen_base', 'caught_stealing', 'grounded_into_double_plays', 'left_on_base', 'pitchers_used', 'wild_pitches', 'assists', 'errors', 'double_def'.

```
In [100]: # Group by Game_ID
grouped = combo_df.groupby('ID')

# Convert groups to a List of DataFrames
game_list = [group for _, group in grouped]

# Perform train-test split on the List of games
train_games, test_games = train_test_split(game_list, test_size=0.2, random_st

# Concatenate the DataFrames back into train and test sets
train_df = pd.concat(train_games).reset_index(drop=True)
test_df = pd.concat(test_games).reset_index(drop=True)
```

```
In [101]: train_df.head()
```

Out[101]:

	date	team_visiting	team_home	score_visiting	score_home	outs_in_game	at_bats	hits
0	20110406	OAK	TOR	3	5	51	32.0	11
1	20110406	OAK	TOR	3	5	51	32.0	6
2	20190630	ARI	SFN	4	10	51	35.0	13
3	20190630	ARI	SFN	4	10	51	34.0	8
4	20210608	ATL	PHI	9	5	54	35.0	11



```
In [102]: # The X features were selected based on EDA and domain interest
X_train= train_df[[
    'at_bats',
    'hits',
    'double',
    'triple',
    'home_run',
    'sacrifice_hit',
    'sacrifice_fly',
    'hit_by_pitch',
    'walk',
    'intent_walk',
    'strikeout',
    'stolen_base',
    'caught_stealing',
    'grounded_into_double_plays',
    'left_on_base',
    'pitchers_used',
    'wild_pitches',
    'assists',
    'errors',
    'double_def'
]]
y_train= train_df['win']

X_test= test_df[[
    'at_bats',
    'hits',
    'double',
    'triple',
    'home_run',
    'sacrifice_hit',
    'sacrifice_fly',
    'hit_by_pitch',
    'walk',
    'intent_walk',
    'strikeout',
    'stolen_base',
    'caught_stealing',
    'grounded_into_double_plays',
    'left_on_base',
    'pitchers_used',
    'wild_pitches',
    'assists',
    'errors',
    'double_def'
]]
y_test= test_df['win']
```

```
In [103]: X_train.shape
```

```
Out[103]: (51894, 20)
```

```
In [104]: y_train.shape
```

```
Out[104]: (51894,)
```

```
In [105]: X_test.shape
```

```
Out[105]: (12974, 20)
```

```
In [106]: y_test.shape
```

```
Out[106]: (12974,)
```

```
In [107]: X_train.info()
```

```
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 51894 entries, 0 to 51893
Data columns (total 20 columns):
 #   Column           Non-Null Count  Dtype  
--- 
 0   at_bats          51894 non-null   float64
 1   hits             51894 non-null   int64  
 2   double            51894 non-null   float64
 3   triple            51894 non-null   float64
 4   home_run          51894 non-null   int64  
 5   sacrifice_hit    51894 non-null   int64  
 6   sacrifice_fly    51894 non-null   int64  
 7   hit_by_pitch     51894 non-null   int64  
 8   walk              51894 non-null   int64  
 9   intent_walk      51894 non-null   int64  
 10  strikeout         51894 non-null   int64  
 11  stolen_base      51894 non-null   int64  
 12  caught_stealing  51894 non-null   int64  
 13  grounded_into_double_plays 51894 non-null   int64  
 14  left_on_base     51894 non-null   int64  
 15  pitchers_used    51894 non-null   int64  
 16  wild_pitches    51894 non-null   int64  
 17  assists           51894 non-null   int64  
 18  errors            51894 non-null   int64  
 19  double_def       51894 non-null   int64  
dtypes: float64(3), int64(17)
memory usage: 7.9 MB
```

Preprocessing

Data were standard scaled. X features were all numeric. There were no missing data. There was 50/50 class balance. No other preprocessors were used.

```
In [108]: scaler = StandardScaler()

X_train_scaled = scaler.fit_transform(X_train)
X_test_scaled = scaler.transform(X_test)
```

Baseline (Dummy) Model

The first model is a Dummy Classifier. As expected using the most frequent strategy the accuracy score of this balanced data set is 50%.

```
In [109]: # Instantiate the model using most_frequent
dummy_model = DummyClassifier(strategy="most_frequent")

# fitting the model
dummy_model.fit(X_train_scaled, y_train)

# predicting
y_pred_dummy = dummy_model.predict(X_test_scaled)
dummy_model.score(X_train_scaled, y_train)
```

Out[109]: 0.5

If the model were to predict 1 (a win) each time, it would be accurate 50% of the time. This serves as a basic baseline to use as a comparison for models.

```
In [110]: print(y_train.value_counts(normalize=True))
print(y_test.value_counts(normalize=True))
```

```
1    0.5
0    0.5
Name: win, dtype: float64
0    0.5
1    0.5
Name: win, dtype: float64
```

Decision Tree and Random Forest

The second model is a simple Decision Tree Classifier using parameters for criterion, maximum depth, and the full set of X features listed above. The accuracy score for this model is 71%. The model identifies: hits, at_bats, home_run, intent_walk, pitchers_used, sacrifice_fly, and grounded_into_double_plays as the most important features.

Next a Random Forest Classifier uses a grid search for the estimator, and maximum depth and features. The best parameters are used in the model with the full set of X features listed above. The accuracy score for this model is 78%. The model identifies: hits, at_bats, home_run, intent_walk, walk, double, and errors as the seven most important features.

```
In [111]: # instantiate tree
tree = DecisionTreeClassifier(random_state=13)
```

```
In [112]: # Train the model on the training data
tree.fit(X_train_scaled, y_train)

tree = DecisionTreeClassifier(criterion='gini', max_depth=4, random_state=13)
tree.fit(X_train_scaled, y_train)
y_pred = tree.predict(X_test_scaled)
precision_score(y_test, y_pred, average="weighted")
```

```
Out[112]: 0.7221045915673954
```

```
In [113]: accuracy = accuracy_score(y_test, y_pred)
print("Accuracy:", accuracy)
```

```
Accuracy: 0.7099583782950516
```

```
In [114]: X_train.info()
```

```
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 51894 entries, 0 to 51893
Data columns (total 20 columns):
 #   Column           Non-Null Count  Dtype  
 ---  --  
 0   at_bats          51894 non-null   float64
 1   hits              51894 non-null   int64  
 2   double             51894 non-null   float64
 3   triple             51894 non-null   float64
 4   home_run           51894 non-null   int64  
 5   sacrifice_hit     51894 non-null   int64  
 6   sacrifice_fly      51894 non-null   int64  
 7   hit_by_pitch       51894 non-null   int64  
 8   walk               51894 non-null   int64  
 9   intent_walk        51894 non-null   int64  
 10  strikeout          51894 non-null   int64  
 11  stolen_base        51894 non-null   int64  
 12  caught_stealing    51894 non-null   int64  
 13  grounded_into_double_plays 51894 non-null   int64  
 14  left_on_base        51894 non-null   int64  
 15  pitchers_used       51894 non-null   int64  
 16  wild_pitches        51894 non-null   int64  
 17  assists             51894 non-null   int64  
 18  errors              51894 non-null   int64  
 19  double_def           51894 non-null   int64  
dtypes: float64(3), int64(17)
memory usage: 7.9 MB
```

```
In [115]: # Get feature importances
feature_importances = tree.feature_importances_

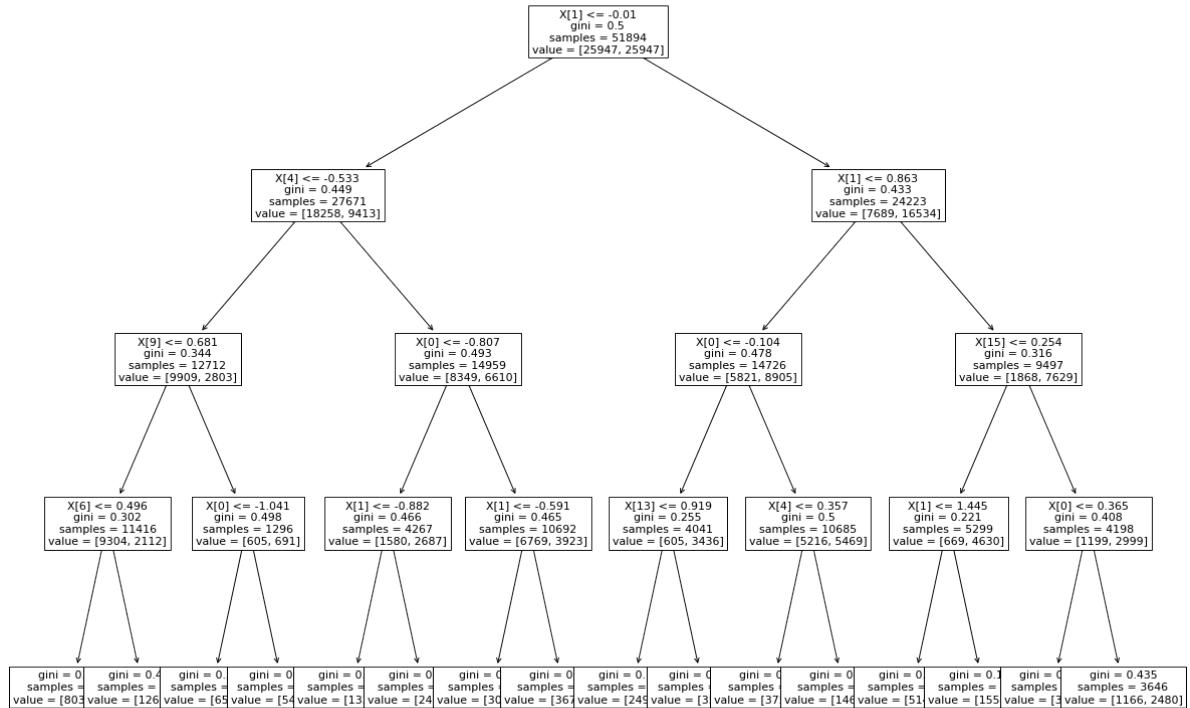
# Get the original feature names
feature_names = X_train.columns

# Map the indices to feature names
feature_indices = tree.tree_.feature
for idx in feature_indices:
    if idx != -2: # -2 is a special value in the tree structure indicating no
        print(f'Feature index {idx} corresponds to feature name: {feature_name}

# Alternatively, get the feature importance along with the feature names
importances = list(zip(feature_names, feature_importances))
importances = sorted(importances, key=lambda x: x[1], reverse=True)
print("Feature importances (sorted):")
for feature, importance in importances:
    print(f'{feature}: {importance:.4f}')
```

Feature index 1 corresponds to feature name: hits
 Feature index 4 corresponds to feature name: home_run
 Feature index 9 corresponds to feature name: intent_walk
 Feature index 6 corresponds to feature name: sacrifice_fly
 Feature index 0 corresponds to feature name: at_bats
 Feature index 0 corresponds to feature name: at_bats
 Feature index 1 corresponds to feature name: hits
 Feature index 1 corresponds to feature name: hits
 Feature index 1 corresponds to feature name: hits
 Feature index 0 corresponds to feature name: at_bats
 Feature index 13 corresponds to feature name: grounded_into_double_plays
 Feature index 4 corresponds to feature name: home_run
 Feature index 15 corresponds to feature name: pitchers_used
 Feature index 1 corresponds to feature name: hits
 Feature index 0 corresponds to feature name: at_bats
 Feature importances (sorted):
 hits: 0.6158
 at_bats: 0.1747
 home_run: 0.1218
 intent_walk: 0.0413
 pitchers_used: 0.0174
 sacrifice_fly: 0.0156
 grounded_into_double_plays: 0.0134
 double: 0.0000
 triple: 0.0000
 sacrifice_hit: 0.0000
 hit_by_pitch: 0.0000
 walk: 0.0000
 strikeout: 0.0000
 stolen_base: 0.0000
 caught_stealing: 0.0000
 left_on_base: 0.0000
 wild_pitches: 0.0000
 assists: 0.0000
 errors: 0.0000
 double_def: 0.0000

```
In [116]: # plotting the figure
plt.figure(figsize=(20,15)) # set plot size (denoted in inches)
plot_tree(tree, fontsize=11)
plt.show()
```



```
In [117]: # Initialize the Random Forest model
rf = RandomForestClassifier(random_state=13)
```

```
In [118]: # Creating a grid search to find the best hyperparameters to include in the model
grid = {
    'n_estimators': [1000, 2000],
    'max_depth': [8, 10, 12],
    'max_features': [8, 10, 12]
}
```

```
In [119]: #grid_search = GridSearchCV(estimator=rf, param_grid=grid, cv=5, scoring='accuracy')
#grid_search.fit(X_train_scaled, y_train)
```

```
In [121]: #grid_search.best_params_
```

```
In [ ]: #grid_search.best_score_
```

```
In [143]: # Train the model on the training data
rf.fit(X_train_scaled, y_train)

rf = RandomForestClassifier(n_estimators=1000, max_features=9, max_depth=8, random_state=42)
rf.fit(X_train_scaled, y_train)
y_pred = rf.predict(X_test_scaled)
precision_score(y_test, y_pred, average="weighted")
```

Out[143]: 0.7841078158789003

```
In [144]: accuracy = accuracy_score(y_test, y_pred)
print("Accuracy:", accuracy)
```

Accuracy: 0.784106674888238

```
In [145]: # Get feature importances
feature_importances = rf.feature_importances_

# Get the original feature names
feature_names = X_train.columns

# Alternatively, get the feature importance along with the feature names
importances = list(zip(feature_names, feature_importances))
importances = sorted(importances, key=lambda x: x[1], reverse=True)
print("Feature importances (sorted):")
for feature, importance in importances:
    print(f'{feature}: {importance:.4f}')
```

Feature importances (sorted):
hits: 0.3819
at_bats: 0.1832
home_run: 0.1568
intent_walk: 0.0684
walk: 0.0318
double: 0.0278
errors: 0.0256
grounded_into_double_plays: 0.0241
pitchers_used: 0.0219
left_on_base: 0.0211
stolen_base: 0.0149
strikeout: 0.0123
sacrifice_fly: 0.0102
sacrifice_hit: 0.0046
wild_pitches: 0.0044
assists: 0.0043
triple: 0.0019
caught_stealing: 0.0019
hit_by_pitch: 0.0016
double_def: 0.0014

Logistic Regression

Next, a logistic regression was performed to classify each game based on the predictors. Logistic Regression is a good model to use for classification. In this instance, a binary classification was performed to determine whether a game was won (1) or lost (0) for each team. The model uses the full set of X features listed above. The accuracy score for this model is 80%.

The second Logistic Regression model uses the subset of X features identified as the most

Log Model 1

```
In [122]: # instantiating the model  
logreg_full = LogisticRegression(random_state=13)  
  
# fitting the model to our scaled training set  
logreg_full.fit(X_train_scaled, y_train)
```

```
Out[122]: LogisticRegression(random_state=13)
```

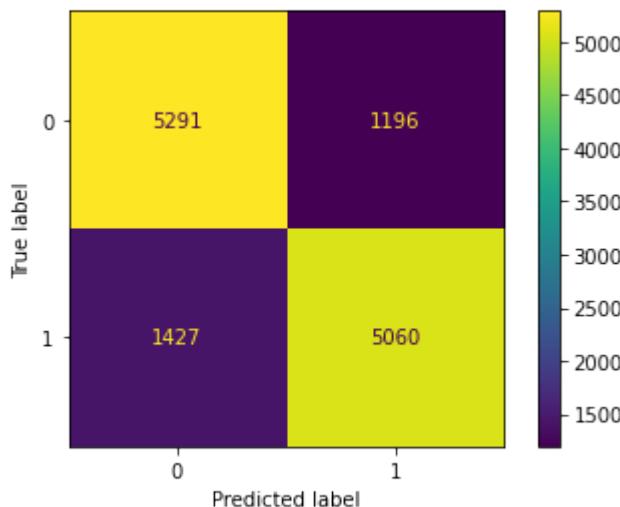
```
In [123]: # predicting our target  
y_pred = logreg_full.predict(X_test_scaled)  
y_pred
```

```
Out[123]: array([0, 1, 1, ..., 1, 1, 0], dtype=int64)
```

```
In [124]: accuracy = accuracy_score(y_test, y_pred)  
print(f'Accuracy: {accuracy:.3f}')
```

Accuracy: 0.798

```
In [125]: plot_confusion_matrix(logreg_full, X_test_scaled, y_test);
```



```
In [126]: # printing out a full classification report to check precision and recall
print(classification_report(y_test, y_pred))
```

	precision	recall	f1-score	support
0	0.79	0.82	0.80	6487
1	0.81	0.78	0.79	6487
accuracy			0.80	12974
macro avg	0.80	0.80	0.80	12974
weighted avg	0.80	0.80	0.80	12974

Log Model 2

```
In [127]: # The X features were selected based on the Random Forest Feature Importance
```

```
X_train_DT= train_df[[
    'hits',
    'at_bats',
    'home_run',
    'intent_walk',
    'walk',
    'double',
    'errors',
    ]]
y_train= train_df['win']

X_test_DT= test_df[[
    'hits',
    'at_bats',
    'home_run',
    'intent_walk',
    'walk',
    'double',
    'errors',
    ]]
y_test= test_df['win']
```

```
In [128]: scaler = StandardScaler()
```

```
X_train_DT_scaled = scaler.fit_transform(X_train_DT)
X_test_DT_scaled = scaler.transform(X_test_DT)
```

```
In [129]: # instantiating the model  
logreg_DT = LogisticRegression(random_state=13)  
  
# fitting the model to our scaled training set  
logreg_DT.fit(X_train_DT_scaled, y_train)
```

```
Out[129]: LogisticRegression(random_state=13)
```

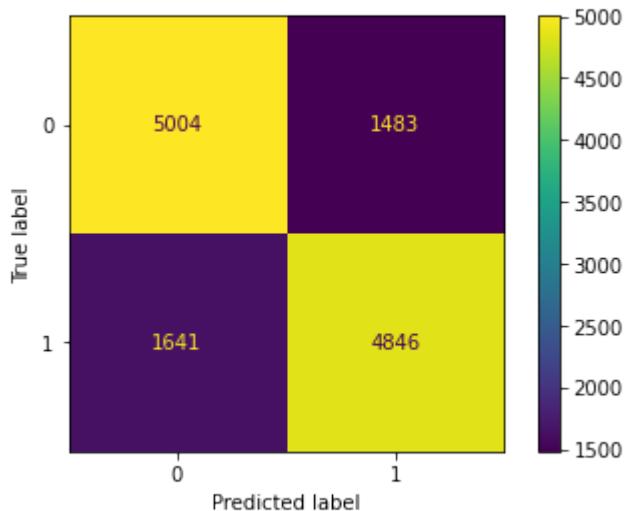
```
In [130]: # predicting our target  
y_pred = logreg_DT.predict(X_test_DT_scaled)  
y_pred
```

```
Out[130]: array([0, 1, 1, ..., 1, 1, 0], dtype=int64)
```

```
In [131]: accuracy = accuracy_score(y_test, y_pred)  
print(f'Accuracy: {accuracy:.3f}')
```

```
Accuracy: 0.759
```

```
In [132]: plot_confusion_matrix(logreg_DT, X_test_DT_scaled, y_test);
```



```
In [133]: # printing out a full classification report to check precision and recall  
print(classification_report(y_test, y_pred))
```

	precision	recall	f1-score	support
0	0.75	0.77	0.76	6487
1	0.77	0.75	0.76	6487
accuracy			0.76	12974
macro avg	0.76	0.76	0.76	12974
weighted avg	0.76	0.76	0.76	12974

```
In [134]: # printing out the coefficients
coefficients = logreg_DT.coef_
print("Coefficients:", coefficients)

Coefficients: [[ 1.71880188 -1.25877175  0.60519292  0.42112437  0.32118146
 0.16742738
 -0.36769609]]
```

```
In [135]: # converting from Log-odds to probabilities using a List comprehension from the
probabilities = [np.exp(x)/(1 + np.exp(x)) for x in coefficients[0]]
probabilities
```



```
Out[135]: [0.8479744468566349,
 0.22118540159720793,
 0.6468434646168807,
 0.6037522698715712,
 0.5796121568339767,
 0.5417593402661469,
 0.40909784376729935]
```

```
In [136]: # creating a dictionary of the probabilities that's easier to read
#dict_of_probs = {0.6809926151:'region_Arusha', 0.5213713979:'region_Kagera',
#dict_of_probs
```



Neural Network

The last model is a Neural Network. This model is built using extensive support from ChatGPT. The full set of X features listed above are used in the model. The accuracy score for this model is 83%; the best performing model. The seven most important feature identified in order of importance are: at_bats, hits, left_on_base, walk, pitchers_used, grounded_into_double_plays, and intent_walk.

```
In [137]: X = X_train_scaled
y = y_train.values

# Convert to PyTorch tensors
X_tensor = torch.tensor(X, dtype=torch.float32)
y_tensor = torch.tensor(y, dtype=torch.long)

# Create a TensorDataset and DataLoader
dataset = TensorDataset(X_tensor, y_tensor)
dataloader = DataLoader(dataset, batch_size=32, shuffle=True)
```

```
In [138]: # Print to verify
#for data in dataloader:
#    inputs, labels = data
#    print("Batch of inputs:\n", inputs)
#    print("Batch of labels:\n", labels)
#    break # Print only the first batch for brevity
```

In [139]:

```
# Define a simple neural network model
class SimpleNN(nn.Module):
    def __init__(self):
        super(SimpleNN, self).__init__()
        self.fc1 = nn.Linear(20, 50) # Adjust input size if necessary
        self.fc2 = nn.Linear(50, 2) # Assuming binary classification (2 output classes)

    def forward(self, x):
        x = torch.relu(self.fc1(x))
        x = self.fc2(x)
        return x

# Initialize the model, loss function, and optimizer
model = SimpleNN()
criterion = nn.CrossEntropyLoss()
optimizer = optim.Adam(model.parameters(), lr=0.001)

# Train the model
num_epochs = 20
for epoch in range(num_epochs):
    model.train()
    running_loss = 0.0
    for inputs, labels in dataloader:
        optimizer.zero_grad()
        outputs = model(inputs)
        loss = criterion(outputs, labels)
        loss.backward()
        optimizer.step()
        running_loss += loss.item()
    print(f"Epoch {epoch+1}, Loss: {running_loss/len(dataloader)}")

# Evaluate the model on the test set
model.eval()
X_test_tensor = torch.tensor(X_test_scaled, dtype=torch.float32)
y_test_tensor = torch.tensor(y_test.values, dtype=torch.long)
with torch.no_grad():
    outputs = model(X_test_tensor)
    probabilities = F.softmax(outputs, dim=1)
    _, predicted = torch.max(outputs, 1)
    #accuracy = (predicted == y_test_tensor).sum().item() / y_test_tensor.size[0]
    #print(f"Accuracy: {accuracy * 100:.2f}%")
print("Class probabilities:")
print(probabilities.numpy())
```

Epoch 1, Loss: 0.44767486626505115
Epoch 2, Loss: 0.40497730395833603
Epoch 3, Loss: 0.39542912543003716
Epoch 4, Loss: 0.38849830988486805
Epoch 5, Loss: 0.38449921404756837
Epoch 6, Loss: 0.3817154008173987
Epoch 7, Loss: 0.3799014177021792
Epoch 8, Loss: 0.3775792658540309
Epoch 9, Loss: 0.37458015781277204
Epoch 10, Loss: 0.37195216151443566
Epoch 11, Loss: 0.3683823105016707
Epoch 12, Loss: 0.364602299629907
Epoch 13, Loss: 0.360594648284477
Epoch 14, Loss: 0.3575948498013134
Epoch 15, Loss: 0.3546399397203646
Epoch 16, Loss: 0.3529066917992917
Epoch 17, Loss: 0.35146464203954036
Epoch 18, Loss: 0.35053367589849577
Epoch 19, Loss: 0.34991715737284335
Epoch 20, Loss: 0.34896557145601576
Class probabilities:
[[0.97051996 0.02948001]
 [0.38047126 0.6195288]
 [0.02450533 0.9754946]
 ...
 [0.3201521 0.67984784]
 [0.03240011 0.96759987]
 [0.9309934 0.06900655]]

```
In [140]: importances = model.fc1.weight.detach().numpy()
feature_importance = np.mean(np.abs(importances), axis=0)
print("Feature Importances:")
for name, importance in zip(X_train.columns, feature_importance):
    print(f"{name}: {importance}")
```

```
Feature Importances:
at_bats: 0.41823819279670715
hits: 0.31008386611938477
double: 0.11506428569555283
triple: 0.14872682094573975
home_run: 0.13251276314258575
sacrifice_hit: 0.14109162986278534
sacrifice_fly: 0.12816965579986572
hit_by_pitch: 0.12602120637893677
walk: 0.2221093326807022
intent_walk: 0.1498233675956726
strikeout: 0.10984161496162415
stolen_base: 0.12294166535139084
caught_stealing: 0.11900687962770462
grounded_into_double_plays: 0.15243425965309143
left_on_base: 0.2723269462585449
pitchers_used: 0.18633097410202026
wild_pitches: 0.11098752915859222
assists: 0.11575867980718613
errors: 0.10560189932584763
double_def: 0.12929143011569977
```

```
In [141]: # Extract feature importances from the first layer of the model  
importances = model.fc1.weight.detach().numpy()  
  
# Calculate the mean absolute importance for each feature  
feature_importance = np.mean(np.abs(importances), axis=0)  
  
# Create a list of (feature_name, importance) tuples  
feature_importance_tuples = [(name, importance) for name, importance in zip(X_  
  
# Sort the list of tuples by importance in descending order  
sorted_feature_importance = sorted(feature_importance_tuples, key=lambda x: x[  
  
# Print the sorted feature importances  
print("Feature Importances (most to least important):")  
for name, importance in sorted_feature_importance:  
    print(f"{name}: {importance:.4f}")
```

```
Feature Importances (most to least important):  
at_bats: 0.4182  
hits: 0.3101  
left_on_base: 0.2723  
walk: 0.2221  
pitchers_used: 0.1863  
grounded_into_double_plays: 0.1524  
intent_walk: 0.1498  
triple: 0.1487  
sacrifice_hit: 0.1411  
home_run: 0.1325  
double_def: 0.1293  
sacrifice_fly: 0.1282  
hit_by_pitch: 0.1260  
stolen_base: 0.1229  
caught_stealing: 0.1190  
assists: 0.1158  
double: 0.1151  
wild_pitches: 0.1110  
strikeout: 0.1098  
errors: 0.1056
```

```
In [142]: #!pip3 install torch torchvision torchaudio --index-url https://download.pytor
```