

## Oakland Athletics: A Sentiment Analysis of MLB Performance

```

pip install pandas
pip install matplotlib
pip install lxml
pip install nltk
pip install wordcloud
pip install tqdm

Requirement already satisfied: beautifulsoup4 in c:\users\andre\downloads\new folder\envs\rl\lib\site-packages (4.12.2)
Requirement already satisfied: requests in c:\users\andre\downloads\new folder\envs\rl\lib\site-packages (2.27.1)
Requirement already satisfied: soupsieve>1.2 in c:\users\andre\downloads\new folder\envs\rl\lib\site-packages (from beautifulsoup4) (2.4.1)
Requirement already satisfied: urllib3<1.27,=>1.21.1 in c:\users\andre\downloads\new folder\envs\rl\lib\site-packages (from requests) (2021.10.8)
Requirement already satisfied: charset-normalizer<2.0.0 in c:\users\andre\downloads\new folder\envs\rl\lib\site-packages (from requests) (2.0.4)
Requirement already satisfied: idna<=2.5 in c:\users\andre\downloads\new folder\envs\rl\lib\site-packages (from requests) (3.3)
Requirement already satisfied: certifi>2017.4.17 in c:\users\andre\downloads\new folder\envs\rl\lib\site-packages (from requests) (2021.10.8)
Requirement already satisfied: pandas in c:\users\andre\downloads\new folder\envs\rl\lib\site-packages (2.0.2)
Requirement already satisfied: tzdata>2022.1 in c:\users\andre\downloads\new folder\envs\rl\lib\site-packages (from pandas) (2022.3)
Requirement already satisfied: pytz>2020.1 in c:\users\andre\downloads\new folder\envs\rl\lib\site-packages (from pandas) (2023.3)
Requirement already satisfied: python-dateutil>=2.8.2 in c:\users\andre\downloads\new folder\envs\rl\lib\site-packages (from pandas) (2.8.2)
Requirement already satisfied: numpy>=1.20.3 in c:\users\andre\downloads\new folder\envs\rl\lib\site-packages (from python-dateutil>=2.8.2-pandas) (1.16.0)
Requirement already satisfied: matplotlib in c:\users\andre\downloads\new folder\envs\rl\lib\site-packages (3.5.1)
Requirement already satisfied: numpy>=1.20 in c:\users\andre\downloads\new folder\envs\rl\lib\site-packages (from matplotlib) (1.21.5)
Requirement already satisfied: kiwisolver>=1.0.1 in c:\users\andre\downloads\new folder\envs\rl\lib\site-packages (from matplotlib) (1.4.4)
Requirement already satisfied: fonttools>=4.22.0 in c:\users\andre\downloads\new folder\envs\rl\lib\site-packages (from matplotlib) (4.40.0)
Requirement already satisfied: pyparsing>=2.1.1 in c:\users\andre\downloads\new folder\envs\rl\lib\site-packages (from matplotlib) (1.21.3)
Requirement already satisfied: cycler>=0.10 in c:\users\andre\downloads\new folder\envs\rl\lib\site-packages (from matplotlib) (0.11.0)
Requirement already satisfied: pillow>=6.2.0 in c:\users\andre\downloads\new folder\envs\rl\lib\site-packages (from matplotlib) (9.5.0)
Requirement already satisfied: importlib-resources>=3.2.0 in c:\users\andre\downloads\new folder\envs\rl\lib\site-packages (from matplotlib) (5.12.0)
Requirement already satisfied: python-dateutil>=2.7 in c:\users\andre\downloads\new folder\envs\rl\lib\site-packages (from matplotlib) (2.8.2)
Requirement already satisfied: packaging>=20.0 in c:\users\andre\downloads\new folder\envs\rl\lib\site-packages (from matplotlib) (21.3)
Requirement already satisfied: contourpy>=1.0.1 in c:\users\andre\downloads\new folder\envs\rl\lib\site-packages (from matplotlib) (1.1.0)
Requirement already satisfied: zipp>=3.1.0 in c:\users\andre\downloads\new folder\envs\rl\lib\site-packages (from importlib-resources>=3.2.0-matplotlib) (3.7.0)
Requirement already satisfied: six>=1.5 in c:\users\andre\downloads\new folder\envs\rl\lib\site-packages (from python-dateutil>=2.7-matplotlib) (1.16.0)
Requirement already satisfied: lxml in c:\users\andre\downloads\new folder\envs\rl\lib\site-packages (4.9.2)
Requirement already satisfied: nltk in c:\users\andre\downloads\new folder\envs\rl\lib\site-packages (3.8.1)
Requirement already satisfied: joblib in c:\users\andre\downloads\new folder\envs\rl\lib\site-packages (from nltk) (1.2.0)
Requirement already satisfied: regex>=2021.8.3 in c:\users\andre\downloads\new folder\envs\rl\lib\site-packages (from nltk) (2022.6.3)
Requirement already satisfied: tqdm in c:\users\andre\downloads\new folder\envs\rl\lib\site-packages (from nltk) (4.65.0)
Requirement already satisfied: click in c:\users\andre\downloads\new folder\envs\rl\lib\site-packages (from nltk) (8.0.4)
Requirement already satisfied: colorama in c:\users\andre\downloads\new folder\envs\rl\lib\site-packages (from nltk>=3.8.1) (0.4.4)
Requirement already satisfied: wordcloud in c:\users\andre\downloads\new folder\envs\rl\lib\site-packages (1.9.2)
Requirement already satisfied: pillow in c:\users\andre\downloads\new folder\envs\rl\lib\site-packages (from wordcloud) (9.5.0)
Requirement already satisfied: numpy>=1.6.1 in c:\users\andre\downloads\new folder\envs\rl\lib\site-packages (from wordcloud) (1.21.5)
Requirement already satisfied: matplotlib>=3.0.0 in c:\users\andre\downloads\new folder\envs\rl\lib\site-packages (from wordcloud) (3.7.1)
Requirement already satisfied: python-dateutil>=2.7 in c:\users\andre\downloads\new folder\envs\rl\lib\site-packages (from matplotlib>=3.0.0-wordcloud) (2.8.2)
Requirement already satisfied: cycler>=0.10 in c:\users\andre\downloads\new folder\envs\rl\lib\site-packages (from matplotlib>=3.0.0-wordcloud) (0.11.0)
Requirement already satisfied: kiwisolver>=1.0.1 in c:\users\andre\downloads\new folder\envs\rl\lib\site-packages (from matplotlib>=3.0.0-wordcloud) (1.4.4)
Requirement already satisfied: pyparsing>=2.1.1 in c:\users\andre\downloads\new folder\envs\rl\lib\site-packages (from matplotlib>=3.0.0-wordcloud) (3.0.4)
Requirement already satisfied: importlib-resources>=3.2.0 in c:\users\andre\downloads\new folder\envs\rl\lib\site-packages (from matplotlib>=3.0.0-wordcloud) (5.12.0)
Requirement already satisfied: packaging>=20.0 in c:\users\andre\downloads\new folder\envs\rl\lib\site-packages (from matplotlib>=3.0.0-wordcloud) (21.3)
Requirement already satisfied: contourpy>=1.0.1 in c:\users\andre\downloads\new folder\envs\rl\lib\site-packages (from matplotlib>=3.0.0-wordcloud) (1.1.0)
Requirement already satisfied: zipp>=3.1.0 in c:\users\andre\downloads\new folder\envs\rl\lib\site-packages (from importlib-resources>=3.2.0-matplotlib>=3.0.0-wordcloud) (3.7.0)
Requirement already satisfied: six>=1.5 in c:\users\andre\downloads\new folder\envs\rl\lib\site-packages (from python-dateutil>=2.7-matplotlib>=3.0.0-wordcloud) (1.16.0)
Requirement already satisfied: tqdm in c:\users\andre\downloads\new folder\envs\rl\lib\site-packages (4.65.0)
Requirement already satisfied: colorama in c:\users\andre\downloads\new folder\envs\rl\lib\site-packages (from tqdm) (0.4.4)

In [2]:
# Import required libraries
import requests
from bs4 import BeautifulSoup
import pandas as pd
import matplotlib.pyplot as plt
import lxml

import os
import re
import numpy as np
import string
import glob
import nltk
from wordcloud import WordCloud
from collections import Counter, defaultdict
from nltk.corpus import stopwords
from nltk.tokenize import word_tokenize
from nltk.stem import PorterStemmer
from string import punctuation
from tqdm import tqdm

sw = stopwords.words("english")

The first step is to scrape the text from the blog posts from the mlbrdatafurnomics.com website. The posts include text, dates posted, and the authors name, all of which will be scraped before storing into a dataframe to use and start the text analysis.

In [3]:
def scrape_website(url, file_path):
    # Send a GET request to the website
    response = requests.get(url)

    # Parse the HTML content using BeautifulSoup
    soup = BeautifulSoup(response.content, "html.parser")

    # Find all the news article elements
    articles = soup.find_all("article")

    # Create lists to store the extracted values
    titles = []
    authors = []
    dates = []
    contents = []

    # Iterate over the articles and extract the text
    with open(file_path, "w", encoding="utf-8") as file:
        for article in articles:
            # Extract the article title
            title_element = article.find("h2")
            title = title_element.text.strip() if title_element else ""
            titles.append(title)

            # Extract the article author
            author_element = article.find("span", class="entry-author")
            author = author_element.text.strip() if author_element else ""
            authors.append(author)

            # Extract the article date
            date_element = article.find("time", class="entry-time")
            date = date_element.text.strip() if date_element else ""
            dates.append(date)

            # Extract the article content
            content_element = article.find("div", class="entry-content")
            content = content_element.text.strip() if content_element else ""
            contents.append(content)

    # Write the title and content to the file
    file.write("\n".join([title + " " + content for title, content in zip(titles, contents)]))

```

```

file.write("date: " + date + "\n")
file.write("Content: " + content + "\n")
file.write("----\n")

# Create a DataFrame from the extracted values
data = {
    "Title": titles,
    "Author": authors,
    "Date": dates,
    "Content": contents
}

df = pd.DataFrame(data)

# Drop rows without content
df.dropna()

return df

```

In [4]:

```

# URL of the website to scrape
url = "https://www.mlbtraderumors.com/oakland-athletics?tab=all"
file_path = "%i/Users/andre/OneDrive/Project/mlbtraderumors_oakland_athletics_articles.txt"

# Scrape the website, write the scraped text to a file, and get the head of the DataFrame
df = scrape_website(url, file_path)
df.head()

```

Out[4]:

	Title	Author	Date	Content
0	Dick Hall Passes Away	Darragh McDonald	June 19, 2023	The Orioles have announced that former major l...
1	A's Acquire Yacelski Rios	Nick Deeds	June 18, 2023	6:56 PM: As noted by Justin Toscano of the Atl...
2	Report: Red Sox Interested In Aledmys Diaz	Mark Polishuk	June 18, 2023	The Red Sox are looking to acquire an infiele...
3	A's To Stay In Oakland Through 2024 Season	Nick Deeds	June 17, 2023	While the Athletics seem more likely than ever...
4	A's Select Tyler Wade, Place Kevin Smith On 10...	Mark Polishuk	June 17, 2023	The Athletics placed infielder Kevin Smith on ...

In [5]:

```

# Change the Date column to a datetime object:
df['Date'] = pd.to_datetime(df['Date'], format='%d %b %d, %Y', errors='coerce')

```

In [6]:

```

df.dtypes

```

Out[6]:

```

Title      object
Author     object
Date       datetime64[ns]
Content    object
dtype: object

```

In [7]:

```

# Check for missing values
df.isna().sum()

```

Out[7]:	Title	0
	Author	0
	Date	30
	Content	0
	dtype:	int64
In [8]:	len(df)	
Out[8]:	1030	
In [9]:	df = df.dropna()	
	df	
Out[9]:	Title	Author
	Date	Content

0	Dick Hill Passes Away	Darragh McDonald	2023-06-19	The Orioles have announced that former major l...
1	A's Acquire Tackles Rios	Nick Dees	2023-06-18	6:56 PM: As noted by Justin Tosome of the Atl...
2	Report: Red Sox Interested In Aleksey D...	Mark Polushko	2023-06-18	The Red Sox are looking to acquire an infiele...
3	A's To Stay In Oakland Through 2024 Season	Nick Dees	2023-06-17	While the Athletics seem more likely than ever...
4	A's Select Tyler Wade, Place Kevin Smith On 10...	Mark Polushko	2023-06-17	The Athletics placed infielder Kevin Smith on ...
...				
995	How Jonathan Lucroy Has Helped The A's	Mark Polushko	2018-10-10	Jonathan Lucroy didn't contribute much at the...
996	A's To Start Lian Hendriks In Wild-Card Game	Connor Byrne	2018-10-02	The Yankees and Athletics have named their st...
997	Andrew Triggs Undergoes Thoracic Outlet Surgery	Jeff Todd	2018-09-28	Sent: The A's announced tonight that Trigg...
998	Athletics Plan To Discuss Bob Melvin Extension...	Steve Adams	2018-09-28	The Athletics have been announced that's unexp...

```
# Identify any noise in the data
RE.SUBSTITUTIONS = re.compile(r'[^\w\d@+]{4}|[^\w\d@+]{8}')

def impurity(text, min_len=10):
    """returns the share of suspicious characters in a text"""
    if text == None or len(text) < min_len:
        return 0
    else:
        return len(RE.SUBSTITUTIONS.findall(text))/len(text)
```

```
[In [11]: df['Content'].apply

Out[11]:
<bound method Series.apply of 0 The Orioles have announced that former major l...
1  6:56 PM: As noted by Justin Toscano of the Atl...
2 The Red Sox are looking to acquire an infielder...
3 While the Athletics seem more likely than ever...
4 The Athletics placed infielder Kevin Smith on ...
...
995 Jonathan Lucroy didn't contribute much at the ...
996 The Yankees and Athletics have named their sta...
997 Sept. 28: The A's announced tonight that Trigu...
998 The Athletics have been baseball's most unapp...
999 Ramon Laureano's brilliant play with the Athle...
Name: Content, Length: 1000, dtype: object>
```

```
pd.options.display.max_columns = 100 #??
# Add new column to data frame
df['Impurity'] = df['Content'].apply(lambda x: impurity(x, min_level=1))

# Acquire the top 3 records
df[['Content', 'Impurity']].sort_values(by='Impurity', ascending=False).head(3)
```

		Content	Impurity
111	Major League Baseball's international signing period for 2023 has officially opened up today, w...		0.005563
944	Left-handed pitching prospect Jesus Luzardo is drawing fans from both inside and outside the A...		0.004474
531	Kington fit the mold of many of the Angels' candidates – well-regarded younger executives who a...		0.003306

The steps above displays the highest impurity levels for the Athletics' data set. The impurity levels above are extremely low, being below 1% of all characters in each bin, so not always exciting.

## Character Normalization & Tokenization

```
[in 13]: # Example function to normalize the text in the "Contents" column
def normalize_text(text):
    # Lowercase the text
    text = text.lower()

    # Remove punctuation
    text = text.translate(str.maketrans("", "", string.punctuation))

    # Remove special characters and digits
    text = re.sub(r"[^a-zA-Z]", " ", text)

    # Tokenize the text
    tokens = word_tokenize(text)
```

```
# Remove stopwords
stop_words = set(stopwords.words("English"))
tokens = [token for token in tokens if token not in stop_words]

# Perform stemming
stemmer = PorterStemmer()
tokens = [stemmer.stem(token) for token in tokens]

return tokens

# Apply normalization to the "contents" column
df["Contents_Normalized"] = df["Content"].apply(normalize_text)

# Print the head of the DataFrame with the normalized contents
df[["Content", "Contents_Normalized"]].head()
```

0	The Orioles have announced that former major leaguer Dick Hall has passed away. He was 92 years old.	[or]ol, announc, former, major, leaguer, dick, hall, pass, away, year, old, hall, underw, main, .
1	6:56 PM: As noted by Justin Toscano of the Atlanta Journal-Constitution, a clause in R's last pass...	[pm, note, toscan, atlanta, atlanta_journalconstitution, clause, ri, pass, brave, requir, claus, .
2	The Red Sox are looking to acquire an infielder who can play multiple positions. <a href="#">Massive.com's</a> ...	[red, sox, look, acquire, infield, play, multipl, posit, massivew, come, mcadam, writte, focs, .
3	While the Athletics seem more likely than ever to relocate to Las Vegas after Nevada's governor...	[athlet, seem, like, ever, reloc, la, vega, nevada, governor, sign, bill, provid, mm, public, la, .
4	The Athletics placed infielder Kevin Smith on the 10-day injured list due to a back strain with...	[athlet, place, infield, kevin, smith, day, injur, list, due, back, strain, placemen, reposit, .

Out [14]:	Title	Author	Date	Content	Impurity	Contents Normalized	length
0	Dick Hall Passes Away	Darragh McDonald	06-19	The Orioles have announced that former major league Dick Hall has passed away. He was 92 years...	0.0	lorial, edmond, former, major, hall, pass, away, year, old, hall, underw, tenn...	173
1	A's Acquire Yackel Rios	Nick Deeds	06-18	6:56 PM: As noted by Justin Tascio of the Atlanta Journal-Constitution, a clause in Rios's...	0.0	[pm, note, tascio, rios, part, journalconstru, claus, in, clau...	206
2	Report: Red Sox Interested in Alexs Diaz	Polshak	06-18	The Red Sox are looking to acquire an infielder who can play multiple positions, ...	0.0	[red, sox, look, acquir, inflay, multipl, posit, massloveson, sen...	84
3	A's To Stay in Oakland Through 2014 Season	Nick Deeds	06-17	While the Athletics seem more likely than they would to relocate to Las Vegas after Nevada's...	0.0	[athlet, seem, the, ever, relic, las, nevada, govern, sign, provid...	102

```
4 A's Select Tyler Wade. Place Kevin Smith On  
10-day IL Mark 2023-06-17 The Athletics placed infielder Kevin Smith on  
the 10-day injured list due to a back strain,...
```

0.0 [athlet, place, infield, kevin, smith,  
day, injur, list, due, back, strain,  
placement, retroact...


127

```
In [15]: df['length'].plot(kind='box', vert=False, color='red', figsize=(8, 1))

Out[15]: <Axes: >
```

```
df['Title_length'].plot(kind='box', vert=False, figsize=(8, 1))
```


Out[16]: <Axes: >




A horizontal box plot showing the distribution of 'Title\_length'. The x-axis is labeled 'Title\_length' and ranges from 20 to 100. The plot shows a median around 40, with the interquartile range (IQR) spanning from approximately 35 to 50. Whiskers extend from 20 to 80. There are several outliers represented by open circles at approximately 85, 90, 95, 100, and 105.

```
df['length'].plot(kind='hist', bins=30, figsize=(8,2))
```

Out[17]: <Axes: ylabel='Frequency'>




A histogram showing the frequency distribution of 'length'. The x-axis is labeled 'length' and ranges from 0 to 100. The y-axis is labeled 'Frequency' and ranges from 0 to 10. The distribution is unimodal and slightly right-skewed, with a peak frequency of 10 at a length of approximately 10.



```
[In 18]: # Extract the month from the date and create a new column
df['Month'] = df['date'].dt.month

# Plot the average post length
df.groupby('Month').agg({'length': 'mean'}) \
    .plot(title='Avg. Post Length', ylim=(0,500), figsize=(6,2))

Out[18]: <Axes: title='center': 'Avg. Post Length', xlabel='Month'>
```



Month	Avg. Post Length
1	400
2	350
3	300
4	280
5	250
6	280
7	250
8	280
9	250
10	280
11	300
12	300

```
In [19]: df['count_words(df, column='Contents_Normalized', preprocess=None, min_freq=2):  
# process tokens and update counter  
df['update(doc)']:  
tokens = doc if preprocess is None else preprocess(doc)
```

```
# create counter and run through all data
counter = Counter()
tqdm.pandas() #initialise tqdm for progress bar
df[column].apply(update)

# transform counter into data frame
freq_df = pd.DataFrame.from_dict(counter, orient='index', columns=['freq'])
freq_df['freq_df.freq'] = freq_df['freq']
freq_df.index.name = 'token'

return freq_df.sort_values('freq', ascending=False)
```

```

out[20]:
freq
token
season 3189
year 2288
leagu 1834
mn 1595
oakland 1573

In [21]: #How many tokens are in the df?
len(freq_df)

```

```
In [22]: # top words with 10+ characters
count_words(df, column='Content',
            preprocess=lambda Content: re.findall(r"([w]{10,})", Content)).head(5)

Out[22]:
```

freq	token
734	<b>appearances</b>
375	<b>organization</b>
367	<b>arbitration</b>
316	<b>outfielder</b>

```
assignment 172

In [23]: ax = freq_df.head(15).plot(kind='barh', width=0.95, figsize=(8,3))
          ax.invert_yaxis()
          ax.set(klabel='Frequency', ylabel='Token', title='Top Words')

Out[23]: [Text(0.5, 0, 'Frequency'), Text(0, 0.5, 'Token'), Text(0.5, 1.0, 'Top Words')]

Top Words
season 100
year 95
league 85
city 80
oakland 75
effect 70
trade 65
player 60
```

Import data on the Oakland Athletics from baseball-reference.com.

```

In [24]: https://www.baseball-reference.com/teams/OAK/attend.shtml

# URL of the webpage to scrape
url = "https://www.baseball-reference.com/teams/OAK/attend.shtml"

# Read the HTML table into a list of DataFrames
tables = pd.read_html(url)

```

```
table = tables[0]

# Save the DataFrame as a CSV file
file_path = "C:/Users/andre/OneDrive/Project/MSB_data/athletics_statistics.csv"
table.to_csv(file_path, index=False, header=True)

print("Data saved to", file_path)

Data saved to C:/Users/andre/OneDrive/Project/MSB_data/athletics_statistics.csv

In [25]: path = "C:/Users/andre/OneDrive/Project/MSB_data/"

athletics_df = pd.DataFrame(pd.read_csv(path + 'athletics_statistics.csv'))
athletics_df.head()
```

[illegible]

Year	Tm	Lg	W	L	Finish	Playoffs	Attendance	Attend/G	Rank	Est. Payroll	PPF	BPF	Stadium	
0	2023	Oakland Athletics	AL West	19	55	5	NaN	368146.0	9688.0	15th of 15	\$51,230,000	94	92	RingCentral Coliseum
1	2022	Oakland Athletics	AL West	60	102	5	NaN	787902.0	9849.0	15th of 15	\$50,248.334	95	93	RingCentral Coliseum
2	2021	Oakland Athletics	AL West	86	76	3	NaN	701340.0	8660.0	15th of 15	\$94,555,884	95	95	RingCentral Coliseum
3	2020	Oakland Athletics	AL West	36	24	1	ALDS (2-1)	NaN	NaN	11th of 15	\$85,683,333	94	96	Oakland-Alameda County Coliseum
4	2019	Oakland Athletics	AL West	97	65	2	AL AWC (1-0)	1670734.0	20626.0	10th of 15	\$102,935,833	93	94	Oakland-Alameda County Coliseum

[illegible]

12	2011	Oakland Athletics	Al West	74	88	3	Na/N	1476791.0	18232.0	14th of 14	\$67,094,000	98	98	O.co Coliseum
13	2010	Oakland Athletics	Al West	81	81	2	Na/N	1418391.0	17511.0	13th of 14	\$57,904,900	99	99	Oakland-Alameda County Coliseum
14	2009	Oakland Athletics	Al West	75	87	4	Na/N	1408783.0	17392.0	14th of 14	\$66,945,000	98	98	Oakland-Alameda County Coliseum
15	2008	Oakland Athletics	Al West	75	86	4	Na/N	1665256.0	20559.0	13th of 14	\$47,967,126	94	94	McAfee Coliseum, Tokyo Dome
16	2007	Oakland Athletics	Al West	76	86	3	Na/N	1921844.0	23762.0	12th of 14	\$79,366,940	94	94	McAfee Coliseum
17	2006	Oakland Athletics	Al West	93	69	1	Lost ALC (4-0)	1976625.0	24403.0	12th of 14	\$64,843,079	97	97	McAfee Coliseum

```

18 2005 Athletics West 88 74 2 NaN 21091180 260380 14 555425.762 100 100 McAfee Coliseum
[In [27]:]
# Remove the dollar and comma signs from the following columns
athletics_df['Est_Payroll'] = athletics_df['Est_Payroll'].replace('\$', ''), regex = True
athletics_df['Attendance'] = athletics_df['Attendance'].replace(',', ''), regex = True
athletics_df['Attendance'] = athletics_df['Attendance'].replace('\,', ''), regex = True
athletics_df['Attendance'] = athletics_df['Attendance'].replace(' ', ''), regex = True
athletics_df['Stadium'] = athletics_df['Stadium'].replace('\(', ''), regex = True
athletics_df.head()
Out[27]:
```

	Year	Tm	Lg W	L Finish	Playoffs	Attendance	Attend/G	Rank	Est_Payroll	PFF	BPF	Stadium		
0	2023	Oakland Athletics	AL	19	55	5	NaN	368146.0	96880	15th of 30	\$1230000	94	92	RingCentral Coliseum

[illegible]

```

Lg      object
W       int64
L       int64
Finish  int64
Playoffs object
Attendance float64
Attends/G float64
Rank    object
Est. Payroll float64
PF      int64
BPF     int64
Stadium object
dtype: object

```

```
In [29]: athletics_df.describe()
```

	Year	W	L	Finish	Attendance	Attendance%	Est. Payroll	PPP	SPR
count	19.000000	19.000000			19.000000	1.800000e+01	18.000000	1.900000e+01	19.000000
mean	20.000000	76.718759	75.789474	2.894737	1.518666e+06	1904.338889	6.754950e+07	96.105263	96.105263
std	5.627314	20.456629	17.915118	1.523692	6.682556e+05	157.226100	1.597226e+07	1.940640	2.157538
min	2005.000000	20.000000	67.000000	1.000000	3.681460e+05	8666.000000	4.796713e+07	93.000000	92.000000
25%	2009.000000	71.500000	70.000000	2.000000	1.432724e+06	17688.000000	5.519741e+07	94.500000	94.500000
50%	2014.000000	76.000000	76.000000	3.000000	1.619436e+06	19993.000000	6.484308e+07	97.000000	96.000000
75%	2018.000000	90.500000	87.000000	4.500000	1.799020e+06	22626.000000	7.462514e+07	97.000000	98.000000
max	2023.000000	100.000000	102.000000	5.000000	2.199118e+06	26038.000000	1.023958e+08	100.000000	100.000000

```
plt.figure(figsize=(10, 6))

# Plotting Attendance over the years
plt.plot(athletics_df['Year'], athletics_df['Attendance'], markers='o')
plt.xlabel('Year')
plt.ylabel('Attendance')
plt.title('Attendance of Oakland Athletics over the Years')

Out[31]: Text(0.5, 1.0, 'Attendance of Oakland Athletics over the Years')
```

The plot displays the annual attendance of the Oakland Athletics from 1961 to 2019. The x-axis represents the year, and the y-axis represents attendance in millions. The data points are connected by a blue line with circular markers. The attendance starts at approximately 1.6 million in 1961, peaks at nearly 2.1 million in 1962, and then fluctuates with a notable decline in the mid-1990s, reaching a low of about 1.7 million in 1995. It then rises to a peak of about 2.0 million in 2000 before declining again to around 1.8 million by 2019.


Year	Attendance (Millions)
1961	1.6
1962	2.05
1963	1.95
1964	1.9
1965	1.85
1966	1.8
1967	1.75
1968	1.7
1969	1.65
1970	1.6
1971	1.55
1972	1.5
1973	1.45
1974	1.4
1975	1.35
1976	1.3
1977	1.25
1978	1.2
1979	1.15
1980	1.1
1981	1.05
1982	1.0
1983	0.95
1984	0.9
1985	0.85
1986	0.8
1987	0.75
1988	0.7
1989	0.65
1990	0.6
1991	0.55
1992	0.5
1993	0.45
1994	0.4
1995	0.35
1996	0.3
1997	0.25
1998	0.2
1999	0.15
2000	0.1
2001	0.05
2002	0.0
2003	0.05
2004	0.1
2005	0.15
2006	0.2
2007	0.25
2008	0.3
2009	0.35
2010	0.4
2011	0.45
2012	0.5
2013	0.55
2014	0.6
2015	0.65
2016	0.7
2017	0.75
2018	0.8
2019	0.85

The graph shows the attendance of the 2009 World Championships in Athletics from 2005 to 2025. The attendance was approximately 150,000 in 2009 and 2013, and dropped to approximately 40,000 in 2023.

Year	Attendance
2009	150,000
2013	150,000
2023	40,000

```
[32]: # Plotting Wins over the years
plt.figure(figsize=(10, 6))
plt.plot(athletics_df['year'], athletics_df['W'], marker='o')
plt.xlabel('Year')
plt.ylabel('Wins')
plt.title('Wins of Oakland Athletics over the Years')

Text(0.5, 1.0, 'Wins of Oakland Athletics over the Years')
```



The plot is a line graph with 'Year' on the x-axis (ranging from 1961 to 2012) and 'Wins' on the y-axis (ranging from 60 to 100). The data points are marked with circles and connected by a blue line. The title is 'Wins of Oakland Athletics over the Years'.

Year	Wins
1961	88
1962	92
1963	88
1964	78
1965	78
1966	78
1967	78
1968	82
1969	78
1970	82
1971	92
1972	95
1973	95
1974	95
1975	92
1976	88
1977	82
1978	78
1979	78
1980	78
1981	78
1982	78
1983	78
1984	78
1985	78
1986	78
1987	78
1988	78
1989	78
1990	78
1991	78
1992	78
1993	78
1994	78
1995	78
1996	78
1997	78
1998	78
1999	78
2000	78
2001	78
2002	78
2003	78
2004	78
2005	78
2006	78
2007	78
2008	78
2009	78
2010	78
2011	78
2012	78

Year	Index
2005	70
2006	75
2007	68
2015	68
2016	70
2019	75
2020	38
2022	60
2022.5	19

```
plt.figure(figsize=(10, 6))
plt.plot(athletics_df['Year'], athletics_df['Bat. Payroll'], marker='o')
plt.xlabel('Year')
plt.ylabel('Payroll')
plt.title('Payroll of Oakland Athletics over the Years')

Out[33]: Text(0.5, 1.0, 'Payroll of Oakland Athletics over the Years')
```

Year	Bat. Payroll
1961	0.00
1962	0.00
1963	0.00
1964	0.00
1965	0.00
1966	0.00
1967	0.00
1968	0.00
1969	0.00
1970	0.00
1971	0.00
1972	0.00
1973	0.00
1974	0.00
1975	0.00
1976	0.00
1977	0.00
1978	0.00
1979	0.00
1980	0.00
1981	0.00
1982	0.00
1983	0.00
1984	0.00
1985	0.00
1986	0.00
1987	0.00
1988	0.00
1989	0.00
1990	0.00
1991	0.00
1992	0.00
1993	0.00
1994	0.00
1995	0.00
1996	0.00
1997	0.00
1998	0.00
1999	0.00
2000	0.00
2001	0.00
2002	0.00
2003	0.00
2004	0.00
2005	0.00
2006	0.00
2007	0.00
2008	0.00
2009	0.00
2010	0.00
2011	0.00
2012	0.00
2013	0.00
2014	0.00
2015	0.00
2016	0.00
2017	0.00
2018	0.00
2019	0.00

Year	Dayen
2005	0.56
2006	0.65
2007	0.80
2008	0.48
2009	0.66
2010	0.58
2011	0.67
2012	0.61
2013	0.69
2014	0.82
2015	0.65
2016	0.55
2017	0.52
2018	0.70
2019	0.85
2020	0.51
2021	0.50
2022	0.52

```
correlation = athletics_df.corr()
print(correlation)

# Heatmap of correlation matrix
plt.figure(figsize=(10, 8))
plt.imshow(correlation, cmap="coolwarm", interpolation="none")
plt.colorbar()
plt.xticks(range(len(correlation)), correlation.columns, rotation=90)
plt.yticks(range(len(correlation)), correlation.columns)
plt.title("Correlation Matrix")

Year      W      Y      F      L      Finish      Attendance      Goals / \
Year      1.000000   -0.461394   -0.272585   0.356362   -0.739316   -0.740579
W          -0.461394   1.000000   -0.202504   -0.498125   -0.703231   0.603377
Y          -0.272585   0.292504   1.000000   0.585745   -0.085816   0.285642
F          -0.356362   -0.498125   0.585745   1.000000   -0.544828   0.001918
Finish     -0.739316   0.703231   0.085816   -0.549138   1.000000   0.983150
Attendance -0.740579   0.603377   0.285642   0.001918   0.983150   1.000000
```

	Est. Payroll	0.214157	0.281505	-0.380622	-0.439569	0.160458	0.123188
YFP	FP707114	0.232063	0.232121	-0.070948	0.289479	0.288072	
BFF		-0.516398	0.396448	0.056171	-0.359041	0.512159	0.487661
	Est. Payroll	YFP	BFF				
Year		0.214157	-0.477914	-0.516398			
W		0.281505	0.232063	0.396448			
L		-0.380622	0.252121	0.056171			
Finish		-0.439569	-0.070948	0.259041			
Attendance		0.160458	0.289479	0.512159			
Attend/Gen		0.123188	0.288072	0.487661			
Est. Payroll		1.000000	0.214157	-0.288093			
PFF		0.316857	1.000000	0.905660			
BFF		-0.028493	0.905660	1.000000			
Text(0.5, 1.0, "Correlation Matrix")							

Heatmap showing the correlation matrix for variables: Year, W, L, Finish, Attendance, and Attend/G. The color scale ranges from -0.8 (dark blue) to 0.8 (dark red), with 0.0 being white. The diagonal elements are all 1.0 (white).

Heatmap showing the correlation matrix for the variables: Est Payroll, PPF, BPF, Year, W, L, Finish, Attendance, AttendG, Est Payroll, PPF, BPF. The color scale ranges from -0.6 (dark blue) to 0.2 (dark red).