**Assignment 3**

Rahul Chowdary Velaga - 120377273

```
1 import sqlite3
2 from google.colab import drive
3 drive.mount('/content/gdrive')
4 path = '/content/gdrive/MyDrive/lahman_1871-2022.sqlite'
5 conn = sqlite3.connect(path)
6
```

    Drive already mounted at /content/gdrive; to attempt to forcibly remount, call drive.mount("/content/gdrive", force_remount=True).

Part 1:

Problem 1:

```
 1 # Import the sqlite3 library
 2 import sqlite3
 3 import pandas as pd
 4
 5 # Set the path to your SQLite database
 6 path = '/content/gdrive/MyDrive/lahman_1871-2022.sqlite'
 7
 8 # Connect to the SQLite database
 9 conn = sqlite3.connect(path)
10
11 # Write the SQL query
12 sql_query = '''
13 SELECT
14     t.teamID,
15     t.yearID,
16     t.franchID,
17     t.W,
18     t.G,
19     s.total_payroll,
20     t.W * 100.0 / t.G AS winning_percentage
21 FROM
22     Teams t
23 JOIN
24     (
25         SELECT
26             teamID,
27             yearID,
28             SUM(salary) AS total_payroll
29         FROM
30             Salaries
31         GROUP BY
32             teamID, yearID
33     ) s
34 ON
35     t.teamID = s.teamID AND t.yearID = s.yearID
36 WHERE
37     t.G > 0;
38 '''
39 data = pd.read_sql_query(sql_query, conn)
40 data
41
```

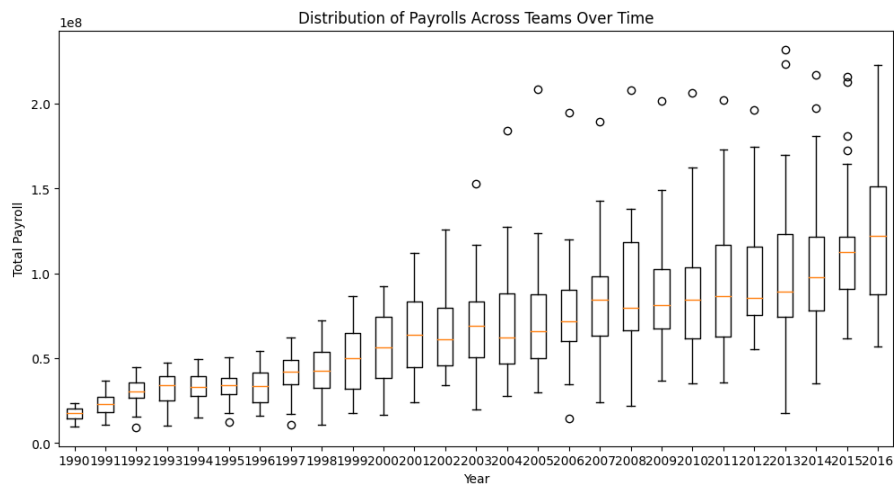|   | teamID | yearID | franchID | W | G | total_payroll | winning_percentage |
|---|--------|--------|----------|---|---|---------------|--------------------|
| 0 | ATL | 1985 | ATL | 66 | 162 | 14807000.0 | 40.740741 |
| 1 | BAL | 1985 | BAL | 83 | 161 | 11560712.0 | 51.552795 |
| 2 | BOS | 1985 | BOS | 81 | 163 | 10897560.0 | 49.693252 |

Part 2:

Problem 2:

```
1 import matplotlib.pyplot as plt
2 new_data = data[(data['yearID'] >= 1990) & (data['yearID'] <= 2022)]
3 plt.figure(figsize=(12, 6))
4 plt.boxplot([new_data[new_data['yearID'] == year]['total_payroll'] for year in new_data['yearID'].unique()], labels=new_data['yearID'].un
5 plt.xlabel('Year')
6 plt.ylabel('Total Payroll')
7 plt.title('Distribution of Payrolls Across Teams Over Time')
8 plt.show()
9 data
```



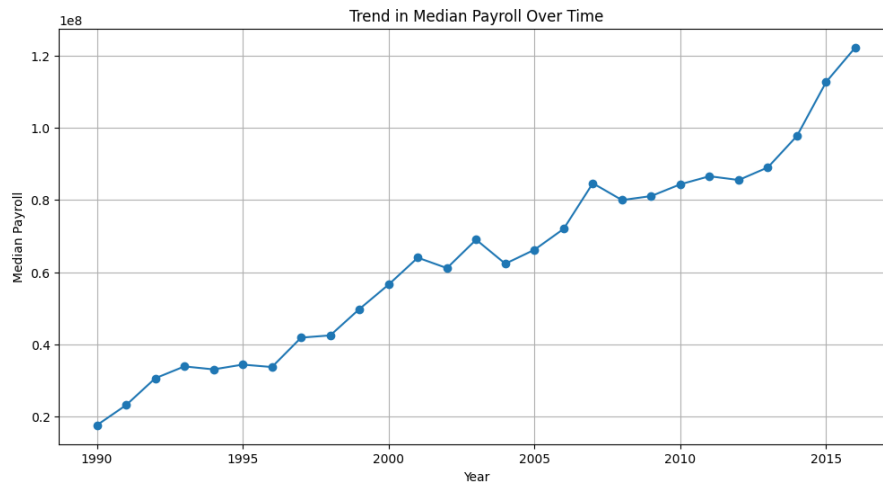|   | teamID | yearID | franchID | W | G | total_payroll | winning_percentage |
|---|--------|--------|----------|---|---|---------------|--------------------|
| 0 | ATL | 1985 | ATL | 66 | 162 | 14807000.0 | 40.740741 |
| 1 | BAL | 1985 | BAL | 83 | 161 | 11560712.0 | 51.552795 |
| 2 | BOS | 1985 | BOS | 81 | 163 | 10897560.0 | 49.693252 |
| 3 | CAL | 1985 | ANA | 90 | 162 | 14427894.0 | 55.555556 |
| 4 | CHA | 1985 | CHW | 85 | 163 | 9846178.0 | 52.147239 |
| ... | ... | ... | ... | ... | ... | ... | ... |
| 913 | SLN | 2016 | STL | 86 | 162 | 143053500.0 | 53.086420 |
| 914 | TBA | 2016 | TBD | 68 | 162 | 57097310.0 | 41.975309 |
| 915 | TEX | 2016 | TEX | 95 | 162 | 176038723.0 | 58.641975 |
| 916 | TOR | 2016 | TOR | 89 | 162 | 138701700.0 | 54.938272 |
| 917 | WAS | 2016 | WSN | 95 | 162 | 141652646.0 | 58.641975 |

918 rows × 7 columns

Question 1: As the years progressed, the median tends to increase gradually. The change in the width of the inter-quartile range does not increase or decrease linearly and is different for different years.

Also from 2003 we can see outliers above the inter-quartile range which indicates that those players are paid way more than an average player.

Problem 3:

```
1 payroll_median = new_data.groupby('yearID')['total_payroll'].median().reset_index()
2 plt.figure(figsize=(12, 6))
3 plt.plot(payroll_median['yearID'], payroll_median['total_payroll'], marker='o')
4 plt.xlabel('Year')
5 plt.ylabel('Median Payroll')
6 plt.title('Trend in Median Payroll Over Time')
7 plt.grid(True)
8 plt.show()
9 data.head()
```
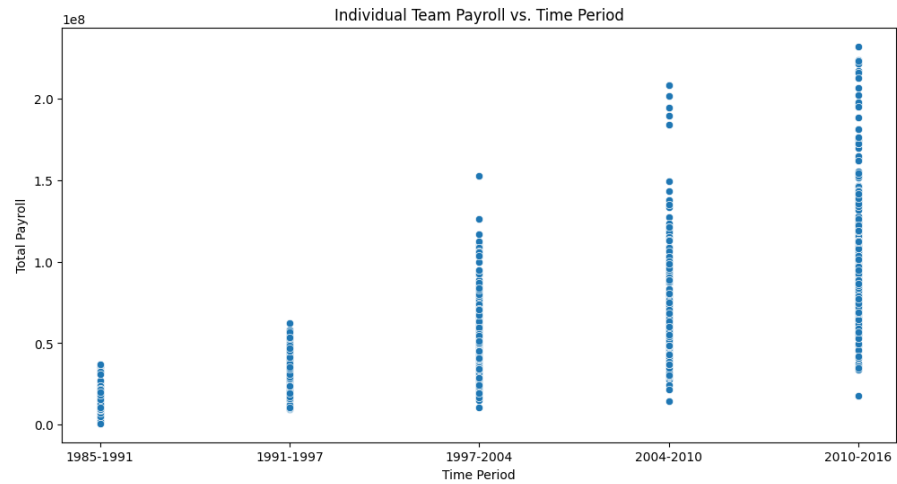


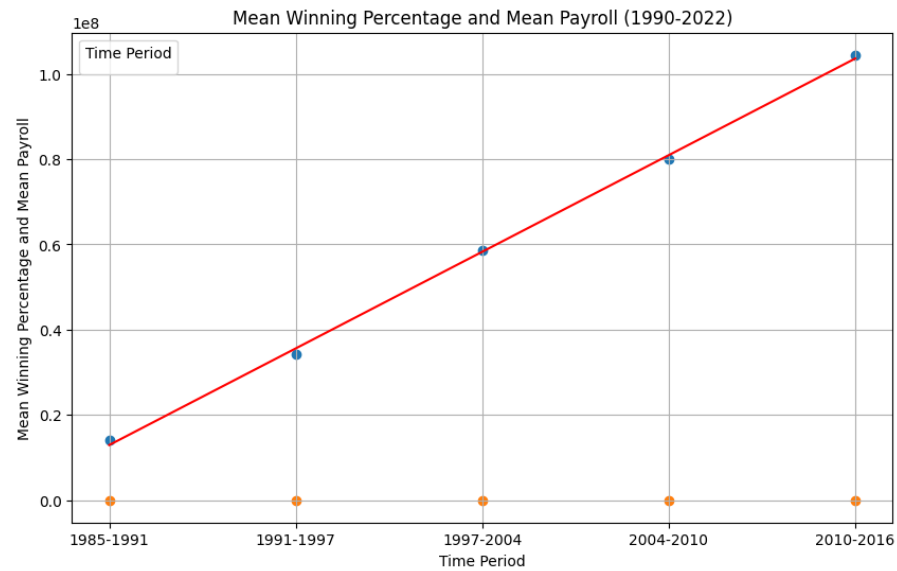| | teamID | yearID | franchID | W | G | total_payroll | winning_percentage |
|---|--------|--------|----------|----|-----|---------------|--------------------|
| 0 | ATL | 1985 | ATL | 66 | 162 | 14807000.0 | 40.740741 |
| 1 | BAL | 1985 | BAL | 83 | 161 | 11560712.0 | 51.552795 |
| 2 | BOS | 1985 | BOS | 81 | 163 | 10897560.0 | 49.693252 |
| 3 | CAL | 1985 | ANA | 90 | 162 | 14427894.0 | 55.555556 |
| 4 | CHA | 1985 | CHW | 85 | 163 | 9846178.0 | 52.147239 |

Problem 4:

```
1 import pandas as pd
2 import numpy as np
3 from sklearn.linear_model import LinearRegression
4 import seaborn as sns
5
6 data['time_period'] = pd.cut(data['yearID'], bins=5)
7 labels = data['time_period'].unique()
8
9 xlabels = []
10
11 for label in labels:
12     xlabels.append('{}-{}'.format(round(label.left), round(label.right)))
13 data['time_period'] = pd.cut(data['yearID'], bins=5, labels=xlabels)
14 labels = data['time_period'].unique()
15
16 mean_winning_percentage = data.groupby('time_period')['winning_percentage'].mean()
17 mean_payroll = data.groupby('time_period')['total_payroll'].mean()
18
```
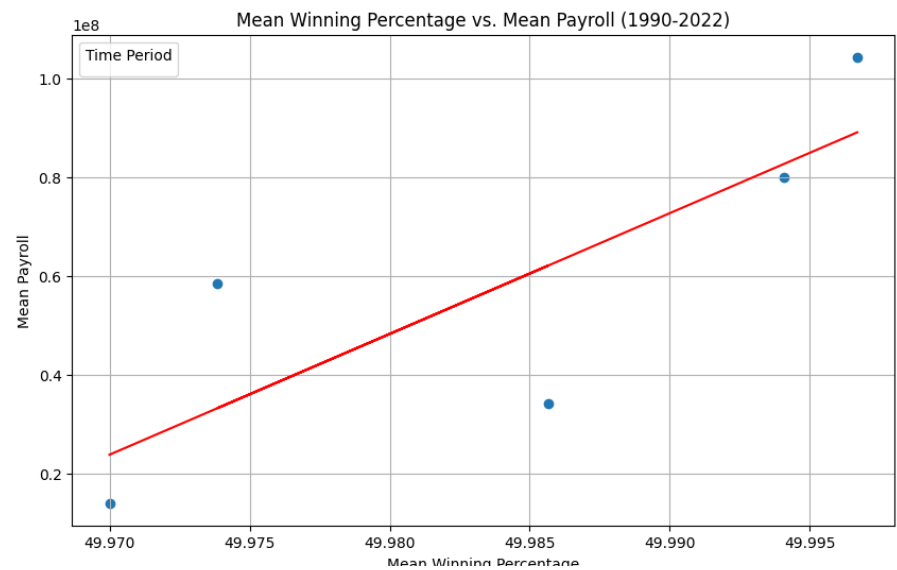
```
19 # payroll vs time
20 plt.figure(figsize=(12, 6))
21 sns.scatterplot(x='time_period', y='total_payroll', data=data)
22 plt.title('Individual Team Payroll vs. Time Period')
23 plt.xlabel('Time Period')
24 plt.ylabel('Total Payroll')
25 plt.show()
26
27
28 # payroll and win percentage vs time: regression
29 plt.figure(figsize=(10, 6))
30 plt.scatter(xlabels, list(mean_payroll))
31 plt.scatter(xlabels, list(mean_winning_percentage))
32
33
34 plt.title('Mean Winning Percentage and Mean Payroll (1990-2022)')
35 plt.xlabel('Time Period')
36 plt.ylabel('Mean Winning Percentage and Mean Payroll')
37 plt.legend(title='Time Period')
38 plt.grid(True)
39
40 x = np.array(range(len(mean_payroll)))
41 # x = mean_winning_percentage
42 y = mean_payroll.values
43 fit = np.polyfit(x, y, deg=1)
44 plt.plot(x, fit[0] * x + fit[1], color='red')
45 plt.show()
46
47 # win percentage vs payroll
48 plt.figure(figsize=(10, 6))
49 plt.scatter(list(mean_winning_percentage), list(mean_payroll))
50
51 plt.title('Mean Winning Percentage vs. Mean Payroll (1990-2022)')
52 plt.ylabel('Mean Payroll')
53 plt.xlabel('Mean Winning Percentage')
54 plt.legend(title='Time Period')
55 plt.grid(True)
56
57 x = mean_winning_percentage
58 y = mean_payroll.values
59 fit = np.polyfit(x, y, deg=1)
60 plt.plot(x, fit[0] * x + fit[1], color='red')
61 plt.show()
62 data.head()
```

Individual Team Payroll vs. Time Period

WARNING:matplotlib.legend:No artists with labels found to put in legend.   Note that arti

Mean Winning Percentage and Mean Payroll (1990-2022)

WARNING:matplotlib.legend:No artists with labels found to put in legend.   Note that arti

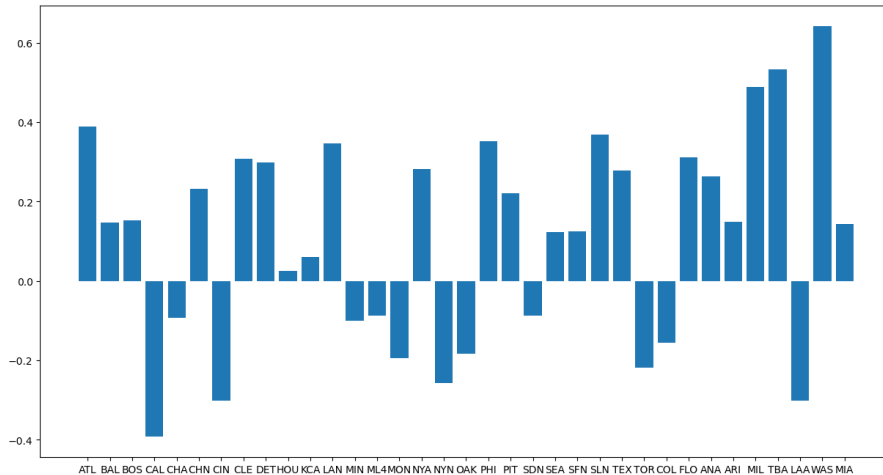Mean Winning Percentage vs. Mean Payroll (1990-2022)

Question 2: The mean payroll of each team has increased over the years but the win percentage has remained almost constant.

```
1  teams = data['teamID'].unique()
2  correlation = []
3  for team_id in teams:
4    team_data = data[data['teamID'] == team_id]
5    cor = team_data['winning_percentage'].corr(team_data['total_payroll'])
6    correlation.append([team_id, cor])
7
8  corr = pd.DataFrame(correlation, columns = ['team_id', 'corr'])
9
10 plt.figure(figsize = (15, 8))
11 plt.bar(corr['team_id'], corr['corr'])
12 plt.show()
```
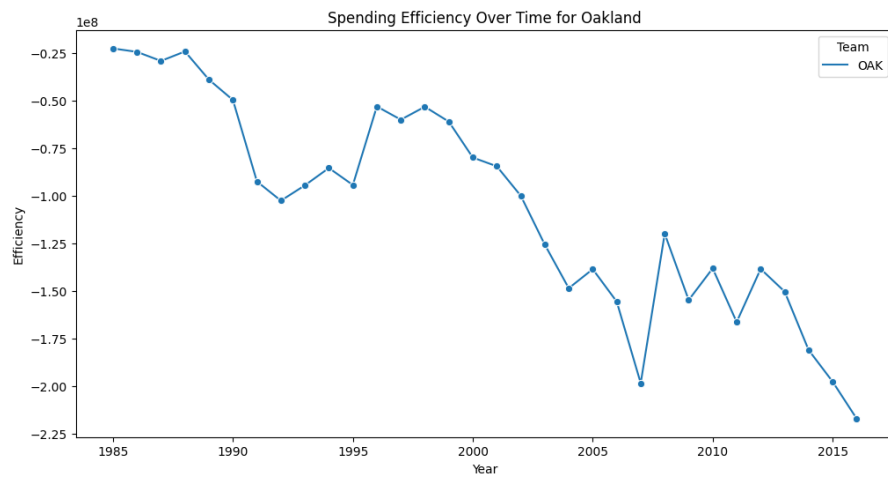


From the above bar plot, the data on the +ve the y axis indicate that the more money the team spent, the more wins they secured. Similarly the data on -ve y axis indicate that the teams spending more money did not result in more wins. MIA, TBA and MIL were particularly good with paying for their wins.

```
1  data['expected_win_pct'] = 50 + 2.5 * data['total_payroll']
2  data['efficiency'] = data['winning_percentage'] - data['expected_win_pct']
3
4  selected_team = ['OAK']
5  selected_team_data = data[data['teamID'].isin(selected_team)]
6
7  plt.figure(figsize=(12, 6))
8  sns.lineplot(x='yearID', y='efficiency', hue='teamID', data=selected_team_data, marker='o')
9
10 plt.xlabel('Year')
11 plt.ylabel('Efficiency')
12 plt.title('Spending Efficiency Over Time for Oakland')
13 plt.legend(title='Team')
14 plt.show()
```

Oakland's spending efficieny decreased non-linearly over the years.

Part 3:

Problem 5:

```
1 data['avg_payroll'] = data.groupby('yearID')['total_payroll'].transform('mean')
2 data['std_payroll'] = data.groupby('yearID')['total_payroll'].transform('std')
3 data['standardized_payroll'] = (data['total_payroll'] - data['avg_payroll']) / data['std_payroll']
4 data.head()
5
```

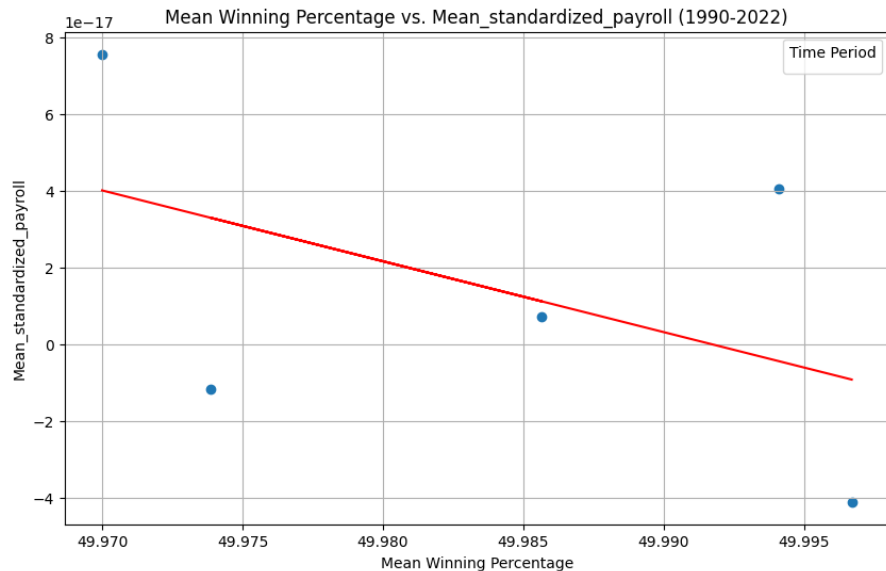| | teamID | yearID | franchID | W | G | total_payroll | winning_percentage | time_period | ex |
|---|---|---|---|---|---|---|---|---|---|
| 0 | ATL | 1985 | ATL | 66 | 162 | 14807000.0 | 40.740741 | 1985-1991 | |
| 1 | BAL | 1985 | BAL | 83 | 161 | 11560712.0 | 51.552795 | 1985-1991 | |
| 2 | BOS | 1985 | BOS | 81 | 163 | 10897560.0 | 49.693252 | 1985-1991 | |
| 3 | CAL | 1985 | ANA | 90 | 162 | 14427894.0 | 55.555556 | 1985-1991 | |
| 4 | CHA | 1985 | CHW | 85 | 163 | 9846178.0 | 52.147239 | 1985-1991 | |

Problem 6:

```
1  data['time_period'] = pd.cut(data['yearID'], bins=5)
2  labels = data['time_period'].unique()
3
4  xlabels = []
5
6  for label in labels:
7    xlabels.append('{}-{}'.format(round(label.left), round(label.right)))
8  data['time_period'] = pd.cut(data['yearID'], bins=5, labels=xlabels)
9  labels = data['time_period'].unique()
10
11 mean_winning_percentage = data.groupby('time_period')['winning_percentage'].mean()
12 mean_standardized_payroll = data.groupby('time_period')['standardized_payroll'].mean()
13
14
15
16 plt.figure(figsize=(10, 6))
17 plt.scatter(list(mean_winning_percentage), list(mean_standardized_payroll))
18
19 plt.title('Mean Winning Percentage vs. Mean_standardized_payroll (1990-2022)')
20 plt.ylabel('Mean_standardized_payroll')
21 plt.xlabel('Mean Winning Percentage')
22 plt.legend(title='Time Period')
23 plt.grid(True)
24
25 # x = np.array(range(len(mean_standardized_payroll)))
26 x = mean_winning_percentage
27 y = mean_standardized_payroll.values
28 fit = np.polyfit(x, y, deg=1)
29 plt.plot(x, fit[0] * x + fit[1], color='red')
30 plt.show()
31 data.head()
```

```
WARNING:matplotlib.legend:No artists with labels found to put in legend.  Note that arti
```



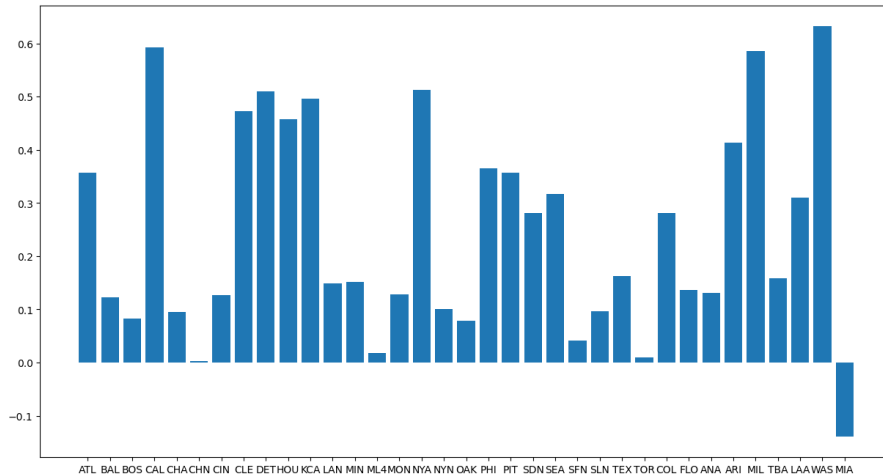| | teamID | yearID | franchID | W | G | total_payroll | winning_percentage | time_period | ex |
|---|--------|--------|----------|-----|-----|---------------|--------------------|-------------|-----|
| 0 | ATL | 1985 | ATL | 66 | 162 | 14807000.0 | 40.740741 | 1985-1991 | |
| 1 | BAL | 1985 | BAL | 83 | 161 | 11560712.0 | 51.552795 | 1985-1991 | |
| 2 | BOS | 1985 | BOS | 81 | 163 | 10897560.0 | 49.693252 | 1985-1991 | |
| 3 | CAL | 1985 | ANA | 90 | 162 | 14427894.0 | 55.555556 | 1985-1991 | |
| 4 | CHA | 1985 | CHW | 85 | 163 | 9846178.0 | 52.147239 | 1985-1991 | |

Question 3: Unlike in problem 4, the standardized payroll has decreased over the years but the win percentage has remained constant.

```
1 teams = data['teamID'].unique()
2 correlation = []
3 for team_id in teams:
4   team_data = data[data['teamID'] == team_id]
5   cor = team_data['winning_percentage'].corr(team_data['standardized_payroll'])
6   correlation.append([team_id, cor])
7
8 corr = pd.DataFrame(correlation, columns = ['team_id', 'corr'])
9
10 plt.figure(figsize = (15, 8))
11 plt.bar(corr['team_id'], corr['corr'])
12 plt.show()
```
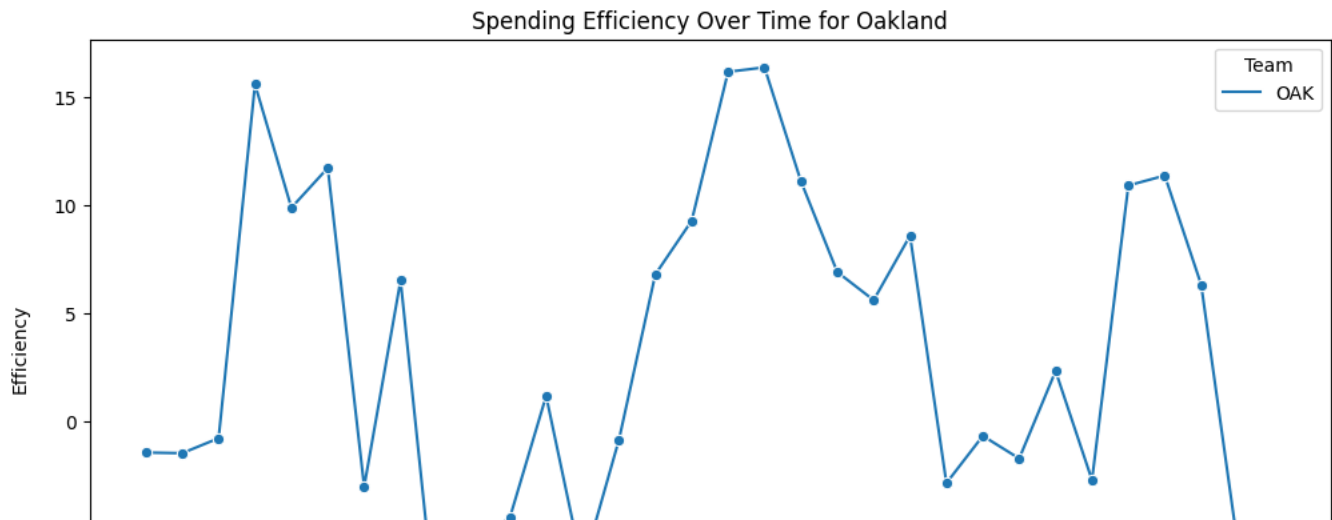


Apart from MIA most of the teams have a +ve correlation between money spent and number of wins secured. The teams that are close to the x axis had no effect on wins even when they spent more money.
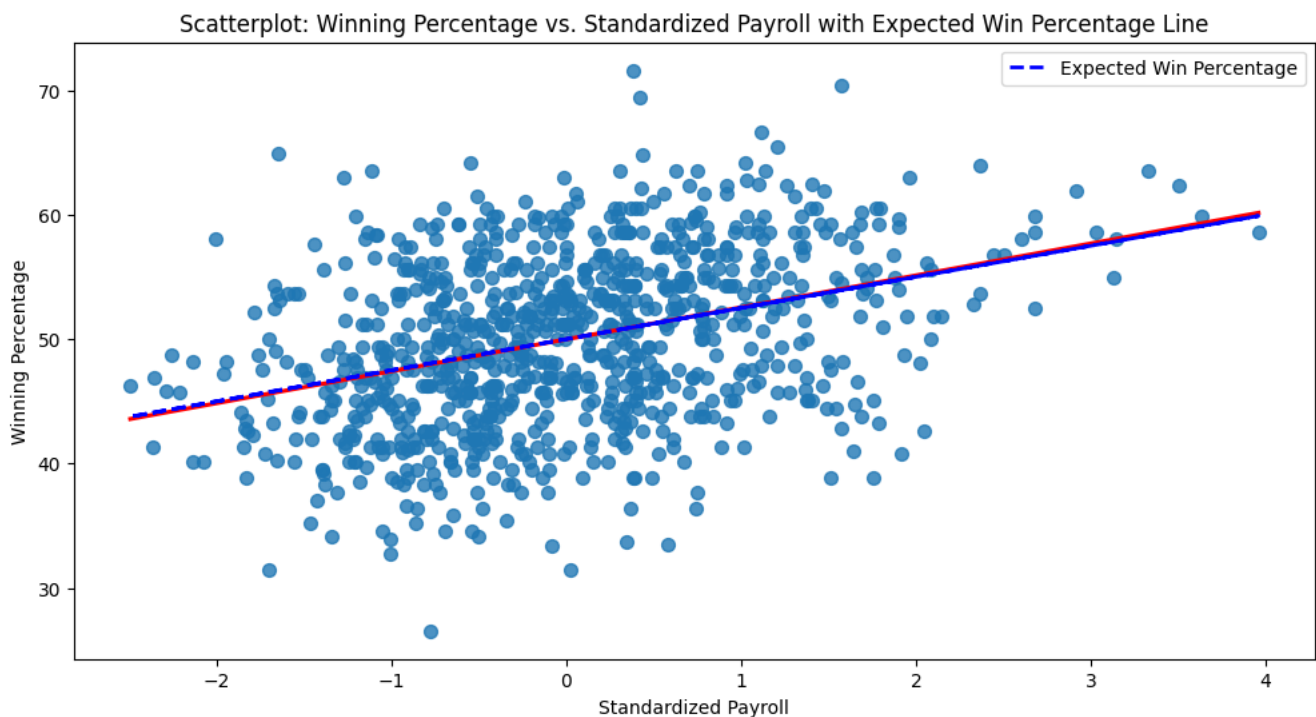
```
1 data['expected_win_pct'] = 50 + 2.5 * data['standardized_payroll']
2 data['efficiency'] = data['winning_percentage'] - data['expected_win_pct']
3
4 selected_team = ['OAK']
5 selected_team_data = data[data['teamID'].isin(selected_team)]
6
7 plt.figure(figsize=(12, 6))
8 sns.lineplot(x='yearID', y='efficiency', hue='teamID', data=selected_team_data, marker='o')
9
10 plt.xlabel('Year')
11 plt.ylabel('Efficiency')
12 plt.title('Spending Efficiency Over Time for Oakland')
13 plt.legend(title='Team')
14 plt.show()
```

## Spending Efficiency Over Time for Oakland



Problem 7:

```
 1 plt.figure(figsize=(12, 6))
 2 sns.regplot(x='standardized_payroll', y='winning_percentage', data=data, scatter_kws={'s': 50}, ci=None, line_kws={'color': 'red'})
 3 expected_win_pct = 50 + 2.5 * data['standardized_payroll']
 4
 5 plt.plot(data['standardized_payroll'], expected_win_pct, color='blue', linestyle='--', linewidth=2, label='Expected Win Percentage')
 6
 7 plt.xlabel('Standardized Payroll')
 8 plt.ylabel('Winning Percentage')
 9 plt.title('Scatterplot: Winning Percentage vs. Standardized Payroll with Expected Win Percentage Line')
10 plt.legend()
11 plt.show()
12
```



Problem 8:

```
 1 data['expected_win_pct'] = 50 + 2.5 * data['standardized_payroll']
 2 data['efficiency'] = data['winning_percentage'] - data['expected_win_pct']
 3
 4 selected_teams = ['OAK', 'NYA', 'BOS', 'ATL', 'TBA']
 5 selected_teams_data = data[data['teamID'].isin(selected_teams)]
 6
 7 plt.figure(figsize=(12, 6))
```

```
 7 plt.figure(figsize=(12, 6))
 8 sns.lineplot(x='yearID', y='efficiency', hue='teamID', data=selected_teams_data, marker='o')
 9
10 plt.xlabel('Year')
11 plt.ylabel('Efficiency')
12 plt.title('Spending Efficiency Over Time for Selected Teams')
13 plt.legend(title='Team')
14 plt.show()
15
```