

```
In [224]: #Importing modules
import pandas as pd
import numpy as np
import math
import datetime
import matplotlib.pyplot as plt
import seaborn as sns
```

```
In [225]: ## Loading dataframes to analyze and cleaning datasets
sal = pd.read_csv("Salary.csv")
stats = pd.read_csv("TeamWins.csv")
stats.drop(stats.iloc[:, 4], inplace = True, axis = 1)
stats.dropna(inplace = True)

print(sal)
print(stats)

      Year      Team  Payroll
0  2016-17  Celtics      93.5
1  2015-16  Celtics      77.1
2  2014-15  Celtics      62.2
3  2013-14  Celtics      70.0
4  2012-13  Celtics      62.6
..      ...      ...      ...
869      NaN      NaN      NaN
870      NaN      NaN      NaN
871      NaN      NaN      NaN
872      NaN      NaN      NaN
873      NaN      NaN      NaN

[874 rows x 3 columns]

      Year      Team Record  Winning Percentage
0  2016-17  Celtics  25-15      0.625
1  2015-16  Celtics  48-34      0.585
2  2014-15  Celtics  40-42      0.488
3  2013-14  Celtics  25-57      0.305
4  2012-13  Celtics  41-40      0.506
..      ...      ...      ...
503  2004-05  Wizards  45-37      0.549
504  2003-04  Wizards  25-57      0.305
505  2002-03  Wizards  37-45      0.451
506  2001-02  Wizards  37-45      0.451
507  2000-01  Wizards  19-63      0.232

[508 rows x 4 columns]

In the previous cell, we have just gotten our data loaded into a dataframe using pandas. From here, we can go on to interpret the date to see if the winning percentage has a correlation with the millions in pay roll a team is spending. In this cell below, we are comparing the data based on a shared index to see if we are missing any values.
```

```
In [226]: compare = datacompy.Compare(sal, stats, on_index = True)
print(compare.report())
output = pd.concat((sal, stats), axis=1)
print(output)

DataCompy Comparison
-----

DataFrame Summary
-----

      DataFrame  Columns  Rows
0      df1          3      874
1      df2          4      508

Column Summary
-----

Number of columns in common: 2
Number of columns in df1 but not in df2: 1
Number of columns in df2 but not in df1: 2

Row Summary
-----

Matched on: index
Any duplicates on match values: No
Absolute Tolerance: 0
Relative Tolerance: 0
Number of rows in common: 508
Number of rows in df1 but not in df2: 366
Number of rows in df2 but not in df1: 0

Number of rows with some compared columns unequal: 0
Number of rows with all compared columns equal: 508

Column Comparison
-----

Number of columns compared with some values unequal: 0
Number of columns compared with all values equal: 2
Total number of values which compare unequal: 0

Sample Rows Only in df1 (First 10 Columns)
-----

      year team payroll
681  NaN  NaN      NaN
873  NaN  NaN      NaN
718  NaN  NaN      NaN
562  NaN  NaN      NaN
793  NaN  NaN      NaN
677  NaN  NaN      NaN
570  NaN  NaN      NaN
752  NaN  NaN      NaN
564  NaN  NaN      NaN
711  NaN  NaN      NaN

      year      team      payroll      year      team record  winning percentage
0  2016-17  Celtics      93.5  2016-17  Celtics  25-15      0.625
1  2015-16  Celtics      77.1  2015-16  Celtics  48-34      0.585
2  2014-15  Celtics      62.2  2014-15  Celtics  40-42      0.488
3  2013-14  Celtics      70.0  2013-14  Celtics  25-57      0.305
4  2012-13  Celtics      62.6  2012-13  Celtics  41-40      0.506
..      ...      ...      ...      ...      ...      ...      ...
503  2004-05  Wizards      60.5  45-37      0.549
504  2003-04  Wizards      52.3  25-57      0.305
505  2002-03  Wizards      53.4  37-45      0.451
506  2001-02  Wizards      51.4  37-45      0.451
507  2000-01  Wizards      50.7  19-63      0.232

[508 rows x 5 columns]

In the previous cell, we merged together the two dataframes and deleted any duplicate entries for the column headers. In the next cell, I will compute some statistics based off of winning percentage and payroll for each team and the overall league.
```

```
In [229]: teams = output['team'].unique()
print(teams)
```

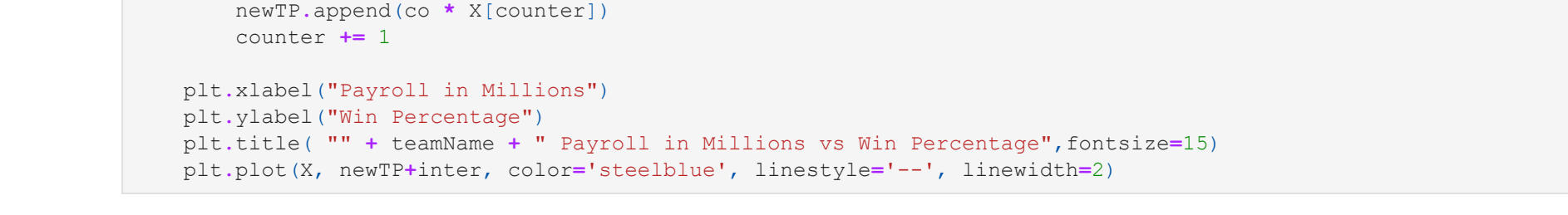
```
['Celtics' 'Hawks' 'Hats' 'Hornets' 'Bulls' 'Cavaliers' 'Mavericks'
 'Nuggets' 'Pistons' 'Warriors' 'Rockets' 'Pacers' 'Clippers' 'Lakers'
 'Grizzlies' 'Heat' 'Bucks' 'Timberwolves' 'Pelicans' 'Knicks' 'Thunder'
 'Magic' '76ers' 'Suns' 'Trail Blazers' 'Kings' 'Spurs' 'Raptors' 'Jazz'
 'Wizards']
```

```
In [230]: teams_map = {}
for i in teams:
    teams_map[i] = output.loc[output['team'] == i]
    team_stat = {}

for key, value in teams_map.items():
    avgpr = value['payroll'].mean()
    avgwp = value['winning percentage'].mean()
    team_stat[key] = (avgpr, avgwp)
print('Team, Average Payroll, Winning Percentage')
totalPay = []
totalWin = []
teamList = []
for key, value in team_stat.items():
    teamList.append(key)
    totalPay.append(value[0])
    totalWin.append(value[1])
    print(key, value)
plt.figure(figsize=(20,20))
plt.scatter(totalPay, totalWin)
for i in range(len(teamList)):
    plt.annotate(teamList[i], (totalPay[i], totalWin[i]))

plt.xlabel("Payroll in Millions")
plt.ylabel("Win Percentage")
plt.title("Payroll in Millions vs Win Percentage", fontsize=15)
AvgPR = sum(totalPay) / len(teamList)
AvgWin = sum(totalWin) / len(teamList)
print(AvgPR, AvgWin)
a, b = np.polyfit(totalPay, totalWin, 1)
newTP = []
for i in range(len(teamList)):
    newTP.append(totalPay[i] * a)
print("Coefficient: " + str(a) + " Intercept = " + str(b))
plt.plot(totalPay, newTP, color='steelblue', linestyle='--', linewidth=2)
```

```
Team, Average Payroll, Winning Percentage
Celtics [65.1058823529412, 0.5418235294117647]
Hawks [67.72352941176469, 0.4718823529411765]
Heat [68.66470588235294, 0.4354117647058824]
Cavaliers [61.06666666666667, 0.4127333333333334]
Bulls [71.6470588235294, 0.49441176470588244]
Hornets [74.5, 0.49905882352941167]
Mavericks [69.8529411764706, 0.6274705882352942]
Mavericks [75.8941176470882, 0.5090588235294118]
Nuggets [75.8941176470882, 0.5090588235294118]
Pistons [61.2470588235294, 0.5132941176470589]
Warriors [76.65294117647058, 0.5016470588235294]
Rockets [76.6411764705882, 0.5656470588235294]
Pacers [76.12352941176471, 0.527164705882353]
Clippers [89.24117647058823, 0.4852352941176471]
Lakers [97.83235294117647, 0.5489411764705883]
Grizzlies [77.2823529411764, 0.4822352941176471]
Heat [86.72352941176472, 0.5446470588235295]
Bucks [75.70588235294119, 0.44376470588235295]
Timberwolves [75.0747058823529, 0.40970588235294114]
Pelicans [73.66470588235293, 0.4725882352941177]
Knicks [89.0235294117647, 0.42188235294117643]
Thunder [82.37058823529412, 0.5424117647058824]
Magic [69.48235294117647, 0.4771176470588236]
76ers [75.08823529411767, 0.5318823529411766]
Suns [81.17352941176471, 0.520588235294118]
Trail Blazers [86.98235294117666, 0.5119999999999999]
Kings [74.17058823529412, 0.4566235294117647]
Spurs [81.17647058823529, 0.42021647058823526]
Raptors [68.05882352941177, 0.471]
Jazz [63.51176470588234, 0.5318823529411766]
Wizards [66.6647058823529, 0.42021647058823526]
Coefficient: 0.0016297329796934065 Intercept = 0.3774097568370732
[<matplotlib.lines.Line2D at 0x1fa8d184acd>]
```



In the previous cell, we plotted the correlation of all the teams payroll in millions over the past two decades in comparison to the win percentage they are able to accomplish. Furthermore, by using simple linear regression through a $a * x + b$, I was able to plot the line of best fit. In the following cell, I will be plotting a regression model for each teams future payroll versus wins.

```
In [231]: from sklearn import linear_model
teamModels = {}
counter = 0
for key, value in teams_map.items():
    teamName = value
    dataframe = value
    X = dataframe['payroll']
    y = dataframe['winning percentage']
    dates = dataframe['year']
    reg = linear_model.LinearRegression()
    reg.fit(dataframe[X], y)
    co = reg.coef_[0]
    inter = reg.intercept_
    print(" " + key + " Coefficient: " + str(co) + " Intercept = " + str(inter))
    plt.figure(figsize=(10,10))
    plt.scatter(X, y)
    newTP = []
    for i in range(len(dates)):
        plt.annotate(dates[counter], (X[counter], y[counter]))
        newTP.append(co * X[counter])
        counter += 1
    plt.xlabel("Payroll in Millions")
    plt.ylabel("Win Percentage")
    plt.title(" " + teamName + " Payroll in Millions vs Win Percentage", fontsize=15)
    plt.plot(X, newTP, color='steelblue', linestyle='--', linewidth=2)
```