



```
In [ ]: # Here are the packages we need

import pandas as pd
import numpy as np
import statsmodels.formula.api as smf
import matplotlib.pyplot as plt
import seaborn as sns

In [ ]: # This command imports our data, which is a Log of games played in 2018 downloaded from Retrosheet
#(you can find the data here: https://www.retrosheet.org/)

MLB = pd.read_excel('../Data/Week 1/Retrosheet MLB game log 2018.xlsx')

In [ ]: # For the Pythagorean Expectation we need only runs scored and conceded. Of course, we also need the names of the teams.
# and the date will also be useful. We put these into a new dataframe (df) which we call MLB18.
# The variable names are rather lengthy, so to make life easier we can rename columns to give them short names.
# If we want to see what the data looks like, we can just type the name of the df.

MLB18 = MLB[['VisitingTeam', 'HomeTeam', 'VisitorRunsScored', 'HomeRunsScore', 'Date']]
MLB18 = MLB18.rename(columns={'VisitorRunsScored': 'VisR', 'HomeRunsScore': 'HomR'})
MLB18

In [ ]: # We will need to know who won the game - which we can tell by who scored the more runs, the home team or the visiting teams
#(there are no ties in baseball)
# The variable 'hwin' is defined here as equaling 1 if the home team scored more runs, and zero otherwise.
# The variable 'awin' is defined in a similar way for the away team.
# we also create a 'count' variable = 1 for each row.

MLB18['hwin'] = np.where(MLB18['HomR'] > MLB18['VisR'], 1, 0)
MLB18['awin'] = np.where(MLB18['HomR'] < MLB18['VisR'], 1, 0)
MLB18['count'] = 1
MLB18

In [ ]: # Since our data refers to games, for each game there are two teams, but what we want is a List of runs scored and conceded
# by each team and its win percentage.
# To create this we are going to define two dfs, one for home teams and one for away teams, which we can then merge to get
# the stats for the entire season.
# Here we define a df for home teams. The command is called ".groupby" and we will use this often. We group by home team
# to obtain the sum of wins and runs (scored and conceded) and also the counter variable to show how many games were played
# (in MLB the teams do not necessarily play the same number of games in the regular season)
# Finally we rename the columns.

MLBhome = MLB18.groupby('HomeTeam')['hwin', 'HomR', 'VisR', 'count'].sum().reset_index()
MLBhome = MLBhome.rename(columns={'HomeTeam': 'team', 'VisR': 'VisRh', 'HomR': 'HomRh', 'count': 'Gh'})
MLBhome
```

Self test - 1 Solution

```
In [ ]: # Now we create a similar df for teams playing as visitors - To write this code all you need to do is to copy and paste
# the previous cell and then change any reference to the home team into a reference to the visiting team.

MLBaway = MLB18.groupby('VisitingTeam')['awin', 'HomR', 'VisR', 'count'].sum()
MLBaway = MLBaway.rename(columns={'VisitingTeam': 'team', 'VisR': 'VisRa', 'HomR': 'HomRa', 'count': 'Ga'})
MLBaway

In [ ]: # We now merge MLBhome and MLBaway so that we have a List of all the clubs with home and away records for the 2018 season
# We will be using pd.merge frequently during the course to combine dfs
# Note that we've called this new df "MLB18", which is name we had already used for earlier df. By doing this we are simply
# overwriting the old MLB18 - which is fine in this case since we don't need the data in the old MLB18 any more.
# If we did want to retain the data in the old MLB18 df, we should have given this new df a different name.

MLB18 = pd.merge(MLBhome, MLBaway, on='team')
MLB18
```

Self test - 2 Solution

```
In [ ]: MLBhome

In [ ]: MLBaway

In [ ]: MLB18 = pd.merge(MLBhome, MLBaway, on='team')
MLB18

In [ ]: # Now we create the total wins, games, played, runs scored and run conceded by summing the totals as home team and away team

MLB18['W'] = MLB18['hwin'] + MLB18['awin']
MLB18['G'] = MLB18['Gh'] + MLB18['Ga']
MLB18['R'] = MLB18['HomRh'] + MLB18['VisRa']
MLB18['RA'] = MLB18['VisRh'] + MLB18['HomRa']
MLB18

In [ ]: # The Last step in preparing the data is to define win percentage and the Pythagorean Expectation.

MLB18['wpc'] = MLB18['W'] / MLB18['G']
MLB18['pyth'] = MLB18['R']**2 / (MLB18['R']**2 + MLB18['RA']**2)
MLB18
```

Self Test - 3 Solution

```
In [ ]: # Having prepared the data, we are now ready to examine it. First, we generate and xy plot use the Seaborn package.
# This illustrates nicely the close correlation between win percentage and the Pythagorean Expectation.

sns.relplot(x="pyth", y="wpc", data = MLB18)
```

Self Test - 4 Solution

```
In [ ]: ▶ # Finally we generate a regression.  
        pyth_lm = smf.ols(formula = 'wpc ~ W', data=MLB18).fit()  
        pyth_lm.summary()
```

```
In [ ]: ▶
```