

Prop Bet - LeBron James

Ken Marciel

7/22/2021

Packages

```
1 # Packages installed in a terminal session:
2 ## pip install --user matplotlib
3 ## pip install --user numpy
4 ## pip install --user pandas
5 ## pip install --user scipy
6 ## pip install --user seaborn
7 ## pip install --user sklearn
8 ## pip install --user statsmodels
9
10 # Packages loaded
11 import matplotlib.pyplot as plt
12 import numpy as np
13 import pandas as pd
14 import scipy.stats as stats
15 import seaborn as sns
16 from sklearn import feature_selection
17 from sklearn import linear_model
18 from sklearn.linear_model import LinearRegression
19 from sklearn.preprocessing import PolynomialFeatures
20 from statsmodels.graphics.tsaplots import plot_pacf
21 from statsmodels.graphics.tsaplots import plot_acf
22 from statsmodels.tsa.statespace.sarimax import SARIMAX
23
```

1. Data Collection

Data was collected from the website, Basketball-Reference.com, by Sports Reference. The data set consists of all the games played by LeBron James in the most recent season. This includes the 2020-2021 regular season and the six games played by in the 2021 playoffs, downloaded as two Excel files, respectively. Using Excel, the files were converted to CSV files. These were then read into Python. The subsets of all games that LeBron James played in were concatenated to exclude the games in which he was inactive. The index of the new file was renumbered from 0 to 50, for the 51 games that LeBron James played in for 2020-2021 regular season and playoffs.

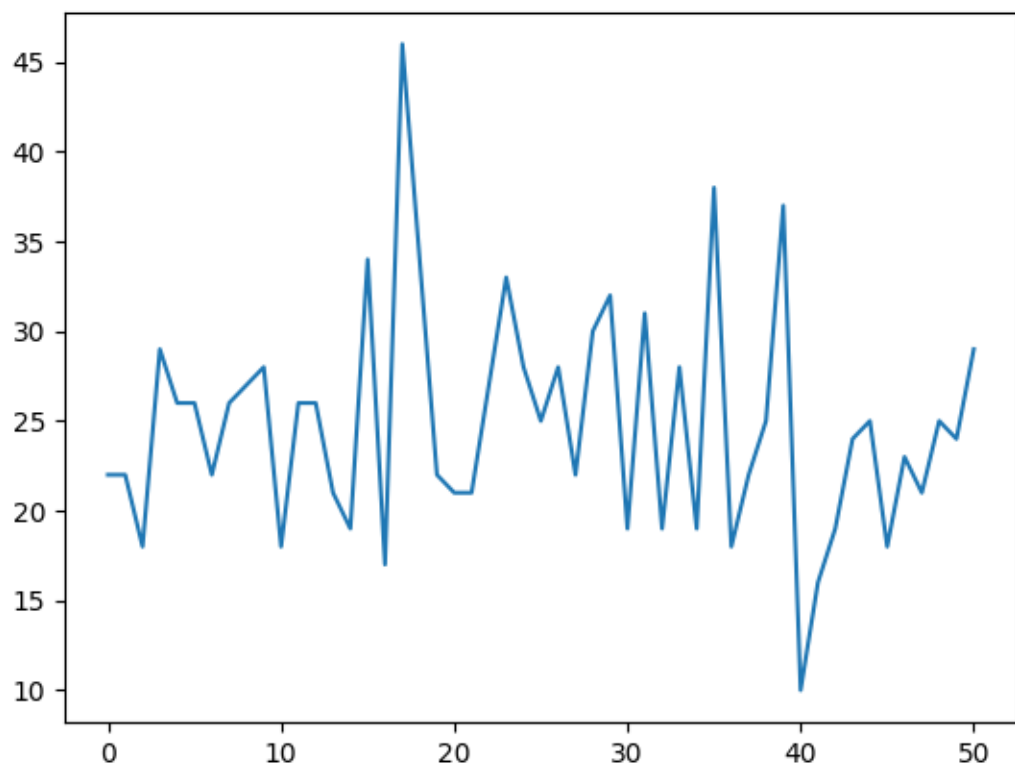
```

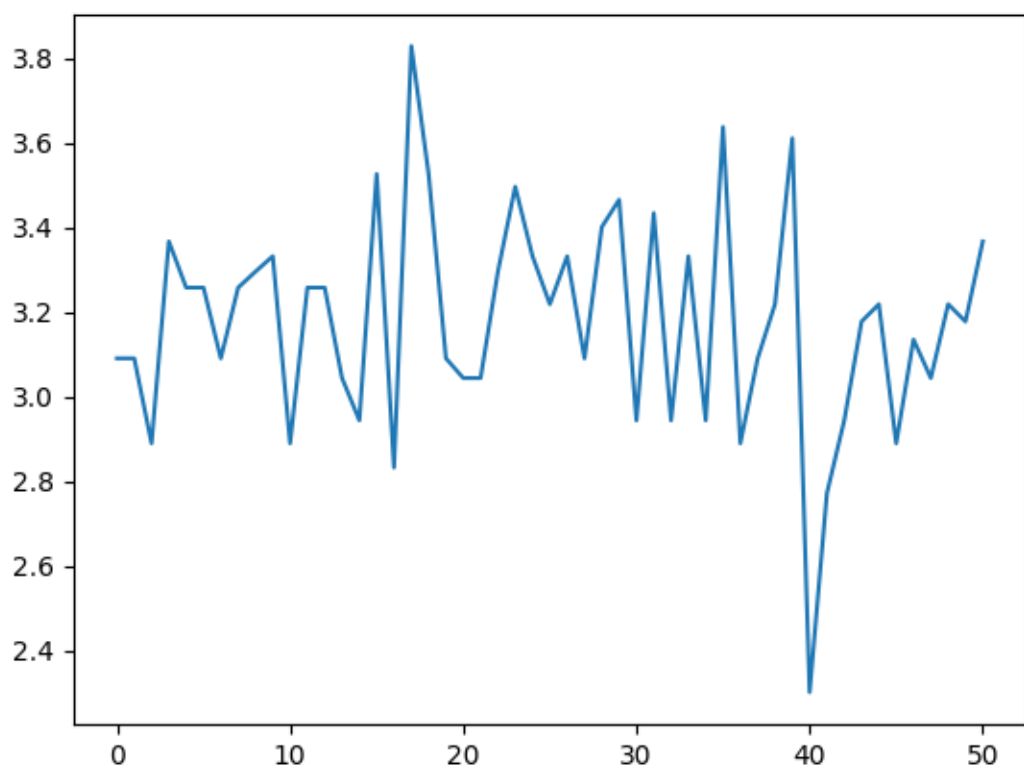
24 # LeBron James: 2020-21 Regular Season
25 # Source: https://www.basketball-reference.com/players/j/jamesle01/gamelog/2021
26 data1 = 'C:/Users/keoka/Desktop/python_work/lbj_2020_21_regular_season.csv'
27 df1 = pd.read_csv(data1)
28 df1.head()
29 len(df1) # 72 regular season games
30
31 # LeBron James: 2021 Playoffs
32 # Source: https://www.basketball-reference.com/players/j/jamesle01/gamelog/2021
33 data2 = 'C:/Users/keoka/Desktop/python_work/lbj_2021_playoffs.csv'
34 df2 = pd.read_csv(data2)
35 df2.head()
36 len(df2) # 6 playoff games
37
38 # LeBron James: 2020-21 Combined Regular Season and Playoffs
39 df3 = pd.concat([df1, df2])
40 df3.head()
41 len(df3) # 78 total games for 2020-21
42
43 # Remove missing values (for inactivity) from data series
44 a = df3[0:35 + 1] # first 36 games played
45 b = df3[37:41 + 1] # next 5 games played
46 c = df3[62:63 + 1] # next 2 games played
47 d = df3[70:77 + 1] # last 8 games played
48 df4 = pd.concat([a,b,c,d])
49 df4 = df4.reset_index() # renumber from 0 to 50
50 len(df4) # played in 51 total games for 2020-21
51

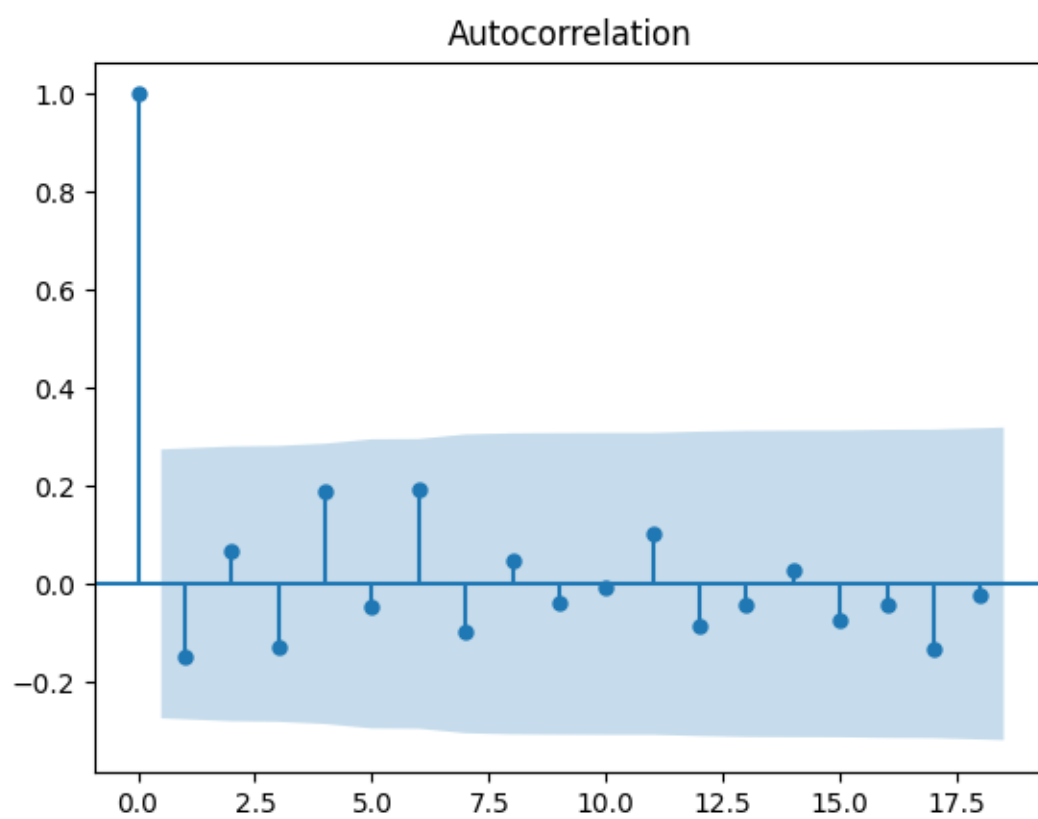
```

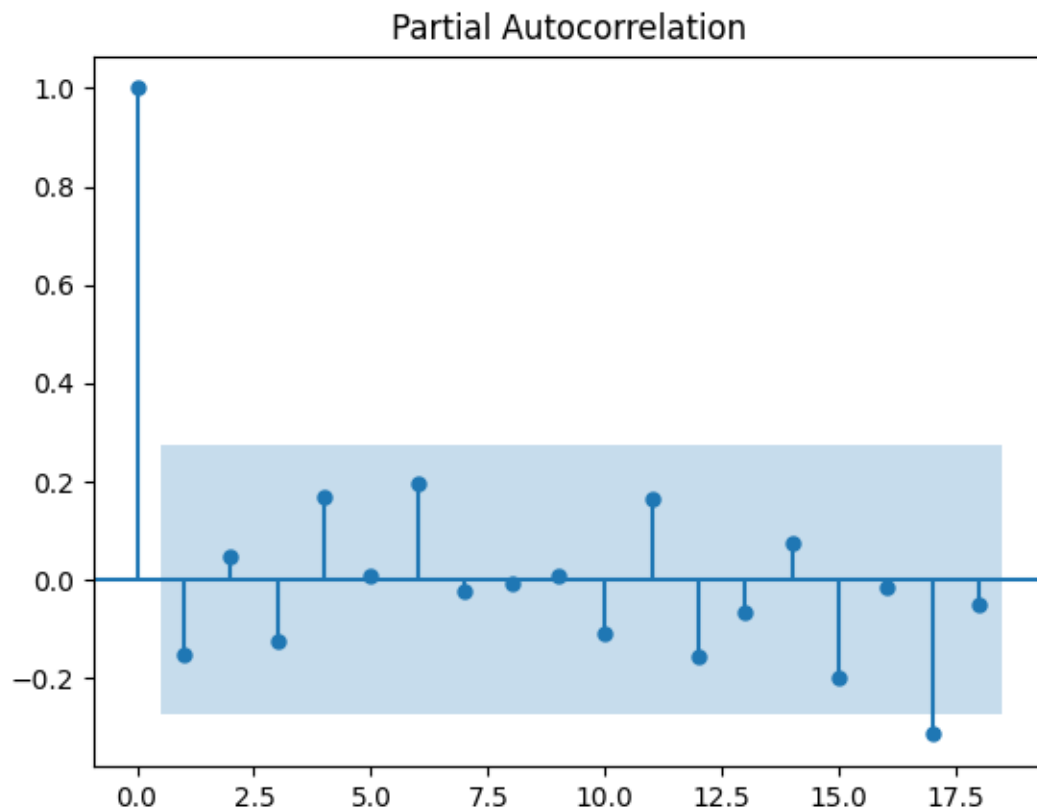
2. Data Transformation

A plot of the time series data displays no noticeable trend, so the series appears stationary. The plot suggests the presence of heteroscedasticity (time-varying variance) and possible autocorrelation. After applying a log transformation to the data, the plot shows stabilized variance. The correlograms (ACF and PACF plots) show a lack of evidence for significant autocorrelation (within the confidence band of $\pm 2/\sqrt{n}$).









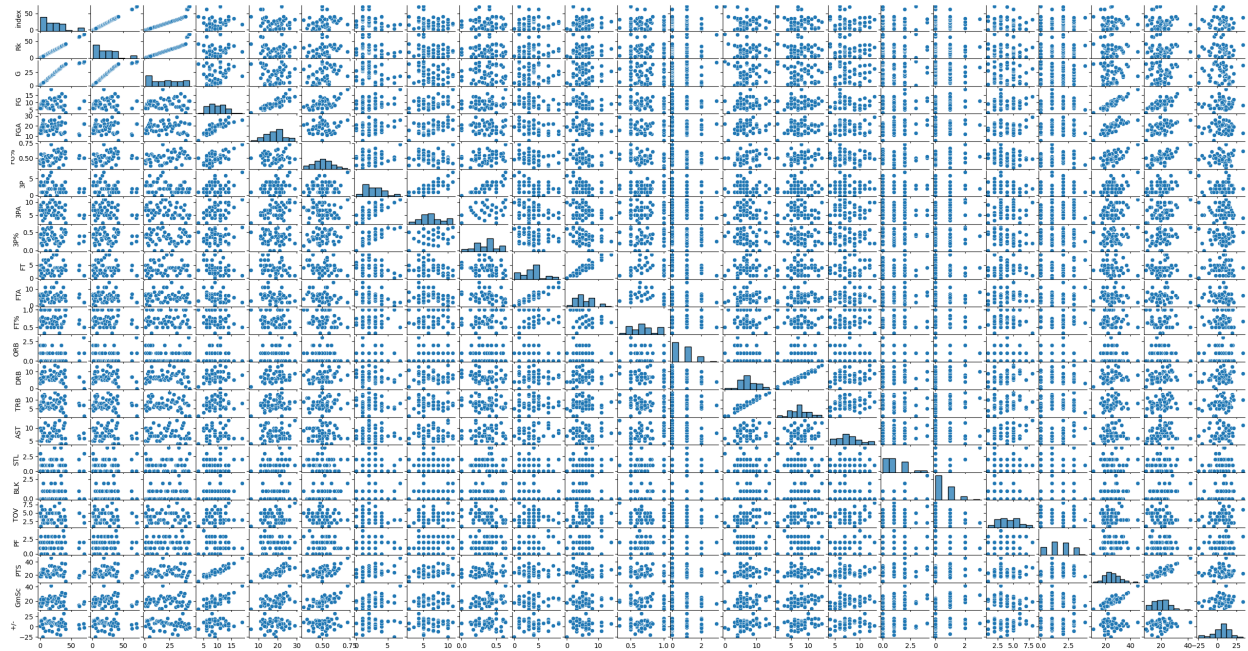
```

52 # Plot the time series data
53 plt.plot(df4['PTS'])
54 plt.show()
55 ## There is no noticeable trend, so the series appears
56 ## stationary. The plot suggests the presence of
57 ## heteroscedasticity and possible autocorrelation.
58
59 # Log transformations to stabilize variance
60 df5 = np.log(df4['PTS'])
61
62 # Plots of transformed data
63 plt.plot(df5) # heteroscedasticity has been reduced
64 plot_acf(df5) # lack of evidence for autocorrelation
65 plot_pacf(df5) # lack of evidence for autocorrelation
66 plt.show()
67

```

3. Model Features

Of the 30 columns in the data file, the 22 metrics recorded for each game were the quantitative features under consideration to formulate a predictive model. A scatterplot matrix reveals that Points (PTS) exhibit the strongest correlation with Field Goals (FG) and Game Score (GmSc), respectively.



```

68 # Scatterplot matrix
69 sns.pairplot(df4)
70 plt.show()
71 ## PTS exhibits strongest correlation with FG and GmSc
72 df4['PTS'].corr(df4['FG']) # 0.9387
73 df4['PTS'].corr(df4['GmSc']) # 0.8608
74 df4['FG'].corr(df4['GmSc']) # 0.7739
75

```

4. Modeling Process

Four models were fit parsimoniously. The first model regresses PTS onto FG. The second model regresses PTS onto GmSc. The third model regresses PTS onto FG and GmSc. The fourth model adds an FG*GmSc interaction term to the third model.


```

76 # Linear model 1: Field Goals
77 x = df4[['FG']]
78 y = df4['PTS']
79 reg1 = linear_model.LinearRegression()
80 reg1.fit(x,y)
81
82 # Linear model 2: Game Score
83 x = df4[['GmSc']]
84 y = df4['PTS']
85 reg2 = linear_model.LinearRegression()
86 reg2.fit(x,y)
87
88 # Linear model 3: Field Goals + Game Score
89 x = df4[['FG', 'GmSc']]
90 y = df4['PTS']
91 reg3 = linear_model.LinearRegression()
92 reg3.fit(x,y)
93
94 # Linear model 4: FG + GmSc + FG*GmSc
95 x = df4[['FG', 'GmSc']]
96 interaction = PolynomialFeatures(degree=1)
97 x = interaction.fit_transform(x)
98 y = df4['PTS']
99 reg4 = linear_model.LinearRegression()
100 reg4.fit(x,y)
101

```

5. Model Testing and Evaluation

The third regression model was selected because it had the highest coefficient of determination (R^2) of 0.9262. This suggests that about 93 percent of the variation in the data is explained by the model. The fourth regression model had the same coefficient of determination, suggesting that the interaction between Field Goals and Game Score is not significant. A one-step-ahead forecast was employed, using the values of FG and GmSc from the last playoff game of 2021, to predict that LeBron James will score 28 points in the second game of the next preseason, on October 6, 2021, between the L.A. Lakers and the Phoenix Suns, in Phoenix.

```

102 # Coefficient of determination
103 reg1.score(x,y) # 0.8812
104 reg2.score(x,y) # 0.7410
105 reg3.score(x,y) # 0.9262
106 reg4.score(x,y) # 0.9262
107
108 # Selected model
109 intercept = reg3.intercept_ # 4.955
110 coef = reg3.coef_ # [1.395, 0.3394]
111
112 # Prediction using data from final game of 2021 playoffs
113 intercept + coef[0]*11 + coef[1]*22.8 # 28 points
114
115 # Same prediction using different code
116 dim1 = np.array([[11],[22.8]])
117 dim2 = np.reshape(dim1, (1,1))
118 reg3.predict(dim2)

```