


```
In [34]: # We use the same function as before but modify it to take into account years
# We use the same function as defined before.
kicker_list = data['Kicker'].unique()
```

```
def data_per_kicker(kickers):
    accuracies = []
    for i in kickers:
        # this is the dataframe for each kicker
        data_kicker = data.loc[(data['Kicker'] == i)]
        # for each kicker we only look at the shots they take that year
        years_played = data_kicker['Year'].unique()
        year = []
        for year in years_played:
            # this is the dataframe for each kicker for each year
            data_kicker_yearly = data_kicker.loc[(data['Year'] == year)]
            success = data_kicker_yearly['Success'].value_counts()
            if len(success) != 1:
                accuracy = success[1]/(success[1] + success[0])
            else:
                accuracy.append({'Year': year, 'Kicker': i, 'Shots_taken': (success[1] + success[0]), 'Success': success[1]})
    return accuracies

results = data_per_kicker(kicker_list)
frame_years = pd.DataFrame(results)
frame_years
```

	Year	Kicker	Shots_taken	Success	Accuracy
0	2005	Akers	22	16	0.727273
1	2006	Akers	27	22	0.814815
2	2007	Akers	32	24	0.750000
3	2008	Akers	50	42	0.840000
4	2009	Akers	37	32	0.864865
...
390	2015	Franks	16	13	0.812500
391	2015	Hockers	14	10	0.714286
392	2015	Hopkins	29	26	0.896552
393	2015	Lambo	32	26	0.812500
394	2015	Myers	30	26	0.866667

395 rows x 5 columns

Since we can't plot all players. Let's create another table where we identify the players with the highest variance in accuracy (e.g. the least consistent) and those who have the lowest variance over the sample space.

For each kicker i calculate the variance $S^2 = \frac{\sum (x_i - \bar{x})^2}{n-1}$ where:

- x_i = the value of the observation
- \bar{x} = the mean value of all observations
- n = the number of observations

```
In [35]: # we can use the var() to do this for us by specifying the column
kickers_names = frame_years['Kicker'].unique()
results_var = []
for name in kickers_names:
    data = frame_years.loc[frame_years['Kicker'] == name]
    var = data.var(numeric_only=True)
    results_var.append({'Kicker': name, 'Variance': var['Accuracy']})

results_var = pd.DataFrame(results_var)
results_var = results_var.sort_values(ascending=False, by='Variance').dropna()
# these are the 5 highest variance players
results_var.head()
```

	Kicker	Variance
58	Henery	0.106169
44	Prater	0.043287
40	Gramatica	0.023802
48	Hauschka	0.021864
25	Novak	0.021560

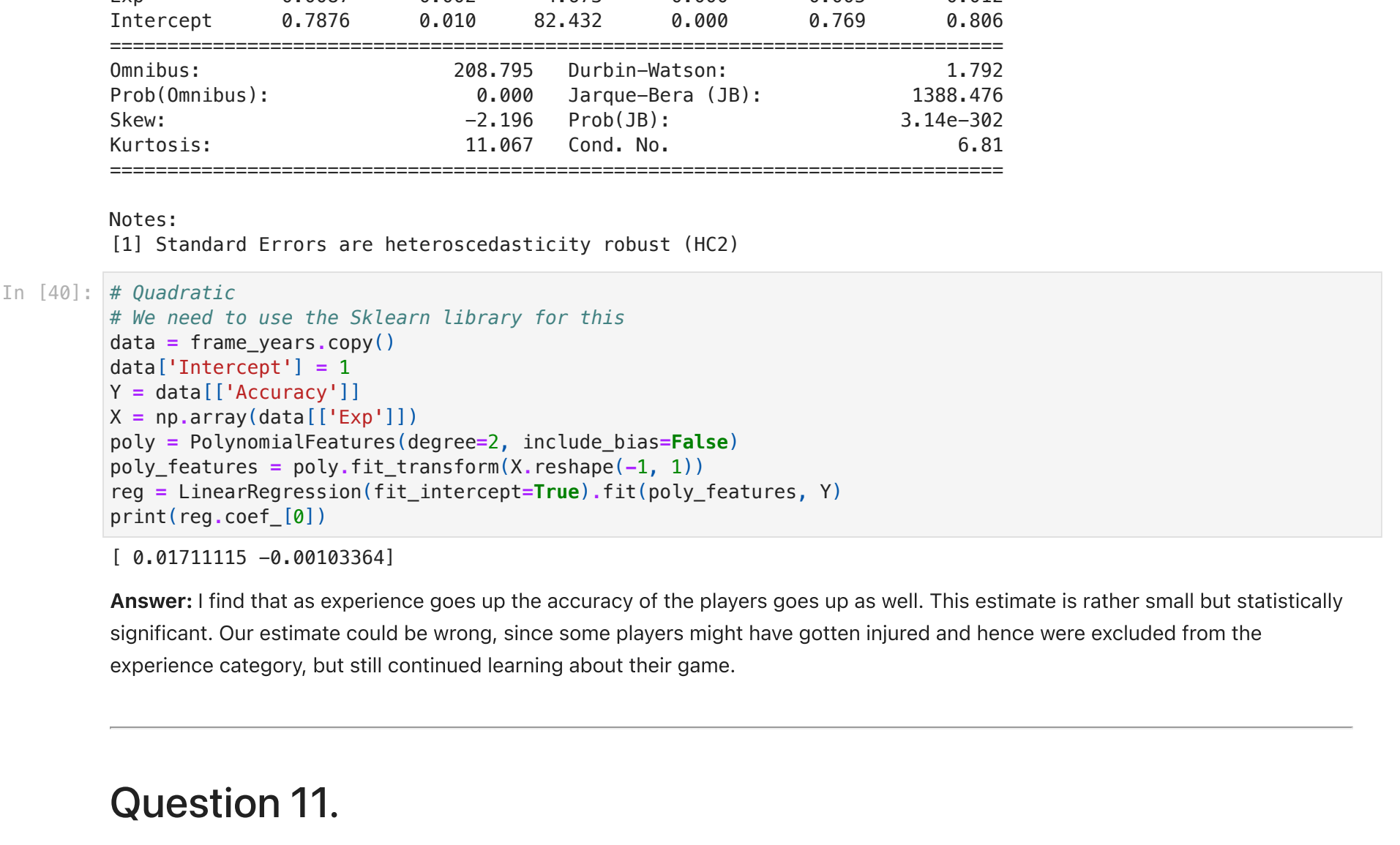
In [36]: # these have the lowest variance
results_var.tail()

	Kicker	Variance
18	Kaeding	0.000680
66	Catanzaro	0.000459
38	Andersen	0.000271
15	Hall	0.000014
71	Santos	0.000000

In [37]: # Let's grab the first 5 kickers and last 5 kickers to plot
kickers = list(results_var['Kicker'].unique())
kickers = kickers[5] + kickers[-5]

fig, axs = plt.subplots(1,1,figsize=(16,9))

```
for x in kickers:
    data = frame_years.loc[frame_years['Kicker'] == x]
    plt.plot(data['Year'], data['Accuracy'], label=f'({x})')
    plt.xlabel('Years')
    plt.ylabel('Accuracy Score')
plt.title('Variation in accuracy score for the 5 most consistent and least consistent players')
plt.legend()
```



Interestingly, we can observe that some players such as Hauschka improved their kick accuracy year by year. This concludes question 9 part 2.

Question 10

10. Some argue that kickers get better with experience, in this dataset do you see evidence to support this conjecture? Try both a linear and quadratic specification. (For simplicity assume that there were no kicks attempted before the beginning of the dataset). What might be wrong with your estimates (besides incomplete data)? Explain!

We can answer this question by looking at the accuracies over time and whether this increases as experience goes up.

```
In [38]: # Let's compute experience in this dataset for each individual
frame_years['Exp'] = 0
names = frame_years['Kicker'].unique()
for i in names:
    frame_years.loc[(frame_years['Kicker'] == i), 'Exp'] = range(0, len(frame_years.loc[frame_years['Kicker'] == i]))
```

```
In [39]: # Linear
data = frame_years.copy()
data['Intercept'] = 1
Y = data['Accuracy']
X = data[['Exp', 'Intercept']]
mod = sm.OLS(Y, X)
res = mod.fit(cov_type='HC2')
print(res.summary())
```

OLS Regression Results

Dep. Variable:	Accuracy	R-squared:	0.053
Model:	OLS	Adj. R-squared:	0.051
Method:	Least Squares	F-statistic:	21.84
Date:	Wed, 01 Feb 2023	Prob (F-statistic):	4.08e-06
Time:	21:59:24	Log-Likelihood:	342.44
No. Observations:	395	AIC:	-680.9
Df Residuals:	393	BIC:	-672.9
Df Model:	1		
Covariance Type:	HC2		

	coef	std err	z	P> z	[0.025	0.975]
Exp	0.0087	0.002	4.673	0.000	0.005	0.012
Intercept	0.7876	0.010	82.432	0.000	0.769	0.806

Prob(Omnibus): 206.795 Durbin-Watson: 1.792
Skewness: -2.196 Prob(JB): 3.14e-302
Kurtosis: 11.867 Cond. No. 6.81

Notes: [1] Standard Errors are heteroscedasticity robust (HC2)

```
In [40]: # Quadratic
# We need to use the Sklearn library for this
data = frame_years.copy()
data['Intercept'] = 1
Y = data['Accuracy']
X = np.array(data[['Exp']])
poly = PolynomialFeatures(degree=2, include_bias=False)
poly_features = poly.fit_transform(X.reshape(-1, 1))
reg = LinearRegression(fit_intercept=True).fit(poly_features, Y)
print(reg.coef_[0])
```

[0.0111115 -0.00183364]

Answer: I find that as experience goes up the accuracy of the players goes up as well. This estimate is rather small but statistically significant. Our estimate could be wrong, since some players might have gotten injured and hence were excluded from the experience category, but still continued learning about their game.

Question 11.

What are the omitted variables you would want in this dataset? List at least 3 and which direction the bias of excluding them could go and why.

I think it would be cool to look at factors such as age and also positional data. I know that there are datasets on Kaggle for example in the (NFL big data bowl competition) where movement(X, Y coordinates) as well as speed at relative times are recorded. I also think it would be interesting to see whether Kickers are better when their contracts are about to expire or whether that does not apply to Kickers. Excluding age could be a positive bias since, we could imagine that younger Kickers are more fit than older ones. On the other hand, we might discover that age is statistically insignificant.

In []: