

# MATH 484 Project

We are trying to predict the future times and results of sprints and field events at the Olympics.

```
In [1]: import matplotlib.pyplot as plt
import pandas as pd
import scipy.stats as stats
import numpy as np
import statsmodels.api as sm
import statsmodels.nonparametric as snp
import statsmodels.stats as sms
import scipy.linalg as lg
import plotly

track_results = pd.read_csv('results.csv')
```

```
In [2]: # Creating separate dictionaries for each event.
tf_results = dict(tuple(track_results.groupby("Event")))
```

I am going to create dataframes and plots of the sprint events from the Summer Olympics that we have results from.

```
In [3]: print(track_results.Event.unique())

['10000M Men' '100M Men' '110M Hurdles Men' '1500M Men' '200M Men'
'20Km Race Walk Men' '3000M Steeplechase Men' '400M Hurdles Men'
'400M Men' '4X100M Relay Men' '4X400M Relay Men' '5000M Men'
'50Km Race Walk Men' '800M Men' 'Decathlon Men' 'Discus Throw Men'
'Hammer Throw Men' 'High Jump Men' 'Javelin Throw Men' 'Long Jump Men'
'Marathon Men' 'Pole Vault Men' 'Shot Put Men' 'Triple Jump Men'
'10000M Women' '100M Hurdles Women' '100M Women' '1500M Women'
'200M Women' '20Km Race Walk Women' '3000M Steeplechase Women'
'400M Hurdles Women' '400M Women' '4X100M Relay Women'
'4X400M Relay Women' '5000M Women' '800M Women' 'Discus Throw Women'
'Hammer Throw Women' 'Heptathlon Women' 'High Jump Women'
'Javelin Throw Women' 'Long Jump Women' 'Marathon Women'
'Pole Vault Women' 'Shot Put Women' 'Triple Jump Women']
```

We will create results scatter plots for field events as well overtime at the Olympics.

## Simple Linear Regression Models for 100 Meter, 200 Meter, Long Jump, Shot Put, 1500 M Run

We start by running a simple linear regression for the 100 M Men and Women's Events.

```
In [4]: Men_100M = tf_results['100M Men']  
Men_100M.head()
```

Out[4]:

	Gender	Event	Location	Year	Medal	Name	Nationality	Result
69	M	100M Men	Rio	2016	G	Usain BOLT	JAM	9.81
70	M	100M Men	Rio	2016	S	Justin GATLIN	USA	9.89
71	M	100M Men	Rio	2016	B	Andre DE GRASSE	CAN	9.91
72	M	100M Men	Beijing	2008	G	Usain BOLT	JAM	9.69
73	M	100M Men	Beijing	2008	S	Richard THOMPSON	TTO	9.89

```
In [5]: X = tf_results['100M Men'].Year.astype(float)
y = tf_results['100M Men'].Result.astype(float)
X = sm.add_constant(X)
model_100M_Men = sm.OLS(y, X).fit()
print(model_100M_Men.params)
model_100M_Men.summary()
```

```
const    37.670494
Year     -0.013918
dtype: float64
```

Out[5]: OLS Regression Results

<b>Dep. Variable:</b>	Result	<b>R-squared:</b>	0.736
<b>Model:</b>	OLS	<b>Adj. R-squared:</b>	0.732
<b>Method:</b>	Least Squares	<b>F-statistic:</b>	217.1
<b>Date:</b>	Wed, 28 Nov 2018	<b>Prob (F-statistic):</b>	3.03e-24
<b>Time:</b>	19:48:50	<b>Log-Likelihood:</b>	-19.177
<b>No. Observations:</b>	80	<b>AIC:</b>	42.35
<b>Df Residuals:</b>	78	<b>BIC:</b>	47.12
<b>Df Model:</b>	1		
<b>Covariance Type:</b>	nonrobust		

	coef	std err	t	P> t	[0.025	0.975]
<b>const</b>	37.6705	1.850	20.366	0.000	33.988	41.353
<b>Year</b>	-0.0139	0.001	-14.735	0.000	-0.016	-0.012

<b>Omnibus:</b>	56.123	<b>Durbin-Watson:</b>	0.350
<b>Prob(Omnibus):</b>	0.000	<b>Jarque-Bera (JB):</b>	239.787
<b>Skew:</b>	2.221	<b>Prob(JB):</b>	8.53e-53
<b>Kurtosis:</b>	10.226	<b>Cond. No.</b>	1.04e+05

Warnings:

- [1] Standard Errors assume that the covariance matrix of the errors is correctly specified.
- [2] The condition number is large, 1.04e+05. This might indicate that there are strong multicollinearity or other numerical problems.

```
In [6]: X = tf_results['100M Women'].Year.astype(float)
y = tf_results['100M Women'].Result.astype(float)
X = sm.add_constant(X)
model_100M_Women = sm.OLS(y, X).fit()
print(model_100M_Women.params)
model_100M_Women.summary()
```

```
const    39.259190
Year     -0.014167
dtype: float64
```

Out[6]: OLS Regression Results

<b>Dep. Variable:</b>	Result	<b>R-squared:</b>	0.763
<b>Model:</b>	OLS	<b>Adj. R-squared:</b>	0.758
<b>Method:</b>	Least Squares	<b>F-statistic:</b>	179.9
<b>Date:</b>	Wed, 28 Nov 2018	<b>Prob (F-statistic):</b>	3.93e-19
<b>Time:</b>	19:48:51	<b>Log-Likelihood:</b>	10.628
<b>No. Observations:</b>	58	<b>AIC:</b>	-17.26
<b>Df Residuals:</b>	56	<b>BIC:</b>	-13.14
<b>Df Model:</b>	1		
<b>Covariance Type:</b>	nonrobust		

	coef	std err	t	P> t	[0.025	0.975]
<b>const</b>	39.2592	2.087	18.811	0.000	35.078	43.440
<b>Year</b>	-0.0142	0.001	-13.413	0.000	-0.016	-0.012

<b>Omnibus:</b>	1.636	<b>Durbin-Watson:</b>	0.818
<b>Prob(Omnibus):</b>	0.441	<b>Jarque-Bera (JB):</b>	0.926
<b>Skew:</b>	0.254	<b>Prob(JB):</b>	0.629
<b>Kurtosis:</b>	3.353	<b>Cond. No.</b>	1.53e+05

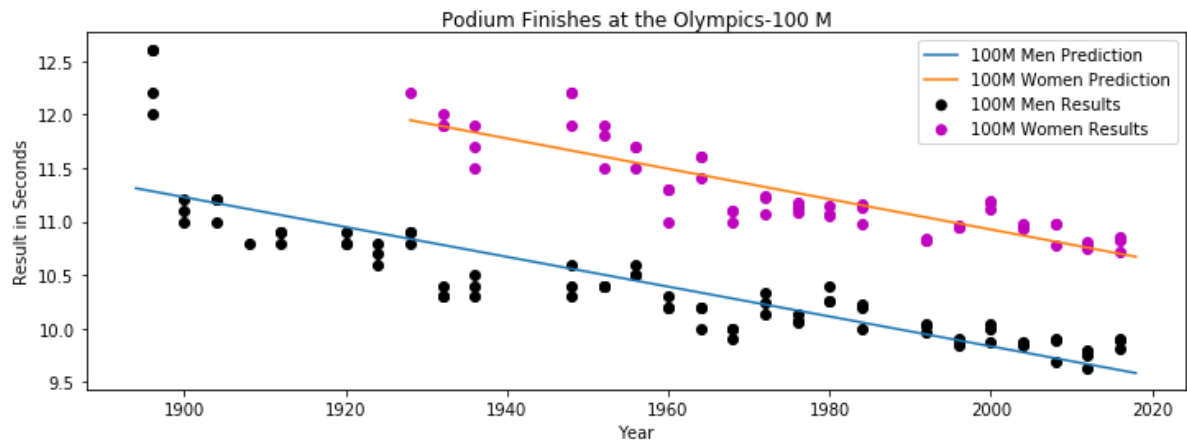
Warnings:

- [1] Standard Errors assume that the covariance matrix of the errors is correctly specified.
- [2] The condition number is large, 1.53e+05. This might indicate that there are strong multicollinearity or other numerical problems.

```

In [7]: pd.to_numeric(tf_results['100M Men'].Result)
pd.to_numeric(tf_results['100M Women'].Result)
time_m = np.linspace(1894, 2018, num = 2000)
time_w = np.linspace(1928, 2018, num = 2000)
Men_pred = 37.670494 - 0.013918 * time_m
Women_pred = 39.259190 - 0.014167 * time_w
plt.figure(figsize=(12, 4))
plt.scatter(tf_results['100M Men'].Year, tf_results['100M Men'].Result, color = 'k')
plt.plot(time_m, Men_pred)
plt.scatter(tf_results['100M Women'].Year, tf_results['100M Women'].Result, color = 'm')
plt.plot(time_w, Women_pred)
plt.xlabel('Year')
plt.ylabel('Result in Seconds')
plt.title('Podium Finishes at the Olympics-100 M')
plt.legend(['100M Men Prediction', '100M Women Prediction', '100M Men Results', '100M Women Results'])
plt.show()

```



200 M Men's and Women's Events now:

```
In [8]: X = tf_results['200M Men'].Year.astype(float)
y = tf_results['200M Men'].Result.astype(float)
X = sm.add_constant(X)
model_200M_Men = sm.OLS(y, X).fit()
print(model_200M_Men.params)
model_200M_Men.summary()
```

```
const      67.599404
Year      -0.023895
dtype: float64
```

Out[8]: OLS Regression Results

<b>Dep. Variable:</b>	Result	<b>R-squared:</b>	0.866
<b>Model:</b>	OLS	<b>Adj. R-squared:</b>	0.864
<b>Method:</b>	Least Squares	<b>F-statistic:</b>	471.6
<b>Date:</b>	Wed, 28 Nov 2018	<b>Prob (F-statistic):</b>	1.40e-33
<b>Time:</b>	19:48:51	<b>Log-Likelihood:</b>	-21.997
<b>No. Observations:</b>	75	<b>AIC:</b>	47.99
<b>Df Residuals:</b>	73	<b>BIC:</b>	52.63
<b>Df Model:</b>	1		
<b>Covariance Type:</b>	nonrobust		

	coef	std err	t	P> t	[0.025	0.975]
<b>const</b>	67.5994	2.159	31.310	0.000	63.297	71.902
<b>Year</b>	-0.0239	0.001	-21.716	0.000	-0.026	-0.022

<b>Omnibus:</b>	0.287	<b>Durbin-Watson:</b>	1.042
<b>Prob(Omnibus):</b>	0.866	<b>Jarque-Bera (JB):</b>	0.470
<b>Skew:</b>	0.059	<b>Prob(JB):</b>	0.791
<b>Kurtosis:</b>	2.631	<b>Cond. No.</b>	1.12e+05

Warnings:

- [1] Standard Errors assume that the covariance matrix of the errors is correctly specified.
- [2] The condition number is large, 1.12e+05. This might indicate that there are strong multicollinearity or other numerical problems.

```
In [9]: X = tf_results['200M Women'].Year.astype(float)
y = tf_results['200M Women'].Result.astype(float)
X = sm.add_constant(X)
model_200M_Women = sm.OLS(y, X).fit()
print(model_200M_Women.params)
model_200M_Women.summary()
```

```
const    94.100684
Year     -0.036020
dtype: float64
```

Out[9]: OLS Regression Results

<b>Dep. Variable:</b>	Result	<b>R-squared:</b>	0.684
<b>Model:</b>	OLS	<b>Adj. R-squared:</b>	0.678
<b>Method:</b>	Least Squares	<b>F-statistic:</b>	106.2
<b>Date:</b>	Wed, 28 Nov 2018	<b>Prob (F-statistic):</b>	7.39e-14
<b>Time:</b>	19:48:51	<b>Log-Likelihood:</b>	-39.140
<b>No. Observations:</b>	51	<b>AIC:</b>	82.28
<b>Df Residuals:</b>	49	<b>BIC:</b>	86.14
<b>Df Model:</b>	1		
<b>Covariance Type:</b>	nonrobust		

	coef	std err	t	P> t	[0.025	0.975]
<b>const</b>	94.1007	6.928	13.583	0.000	80.178	108.023
<b>Year</b>	-0.0360	0.003	-10.303	0.000	-0.043	-0.029

<b>Omnibus:</b>	2.396	<b>Durbin-Watson:</b>	0.507
<b>Prob(Omnibus):</b>	0.302	<b>Jarque-Bera (JB):</b>	2.287
<b>Skew:</b>	0.489	<b>Prob(JB):</b>	0.319
<b>Kurtosis:</b>	2.651	<b>Cond. No.</b>	1.84e+05

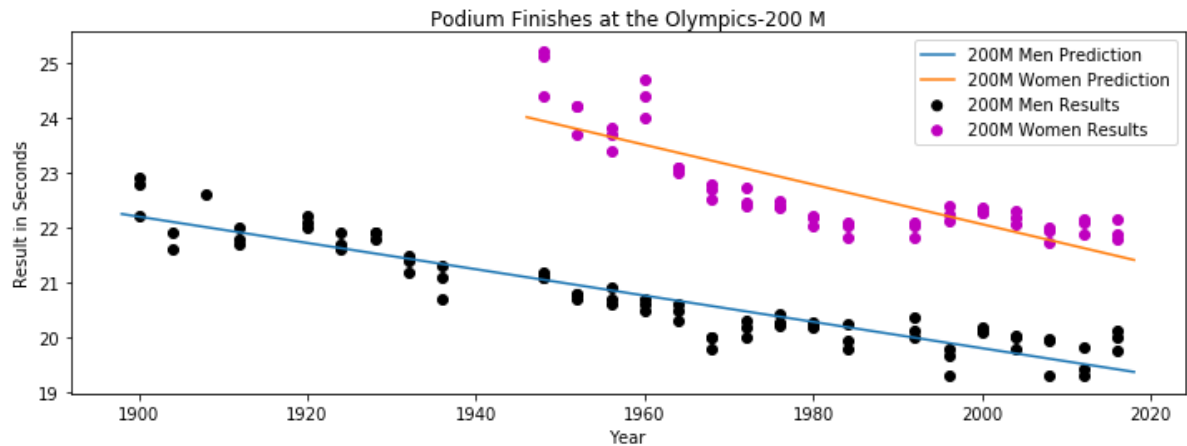
Warnings:

- [1] Standard Errors assume that the covariance matrix of the errors is correctly specified.
- [2] The condition number is large, 1.84e+05. This might indicate that there are strong multicollinearity or other numerical problems.

```

In [10]: time_m = np.linspace(1898, 2018, num = 2000)
time_w = np.linspace(1946, 2018, num = 2000)
Men_pred = 67.599404 - 0.023895 * time_m
Women_pred = 94.100684 - 0.036020 * time_w
plt.figure(figsize=(12, 4))
plt.scatter(tf_results['200M Men'].Year, tf_results['200M Men'].Result, color = 'k')
plt.plot(time_m, Men_pred)
plt.scatter(tf_results['200M Women'].Year, tf_results['200M Women'].Result, color = 'm')
plt.plot(time_w, Women_pred)
plt.xlabel('Year')
plt.ylabel('Result in Seconds')
plt.title('Podium Finishes at the Olympics-200 M')
plt.legend(['200M Men Prediction', '200M Women Prediction', '200M Men Results', '200M Women Results'])
plt.show()

```





```
In [11]: X = tf_results['Long Jump Men'].Year.astype(float)
y = tf_results['Long Jump Men'].Result.astype(float)
X = sm.add_constant(X)
LJ_Men = sm.OLS(y, X).fit()
print(LJ_Men.params)
LJ_Men.summary()
```

```
const    -18.587442
Year       0.013485
dtype: float64
```

Out[11]: OLS Regression Results

<b>Dep. Variable:</b>	Result	<b>R-squared:</b>	0.792
<b>Model:</b>	OLS	<b>Adj. R-squared:</b>	0.788
<b>Method:</b>	Least Squares	<b>F-statistic:</b>	231.9
<b>Date:</b>	Wed, 28 Nov 2018	<b>Prob (F-statistic):</b>	1.87e-22
<b>Time:</b>	19:48:52	<b>Log-Likelihood:</b>	-8.8863
<b>No. Observations:</b>	63	<b>AIC:</b>	21.77
<b>Df Residuals:</b>	61	<b>BIC:</b>	26.06
<b>Df Model:</b>	1		
<b>Covariance Type:</b>	nonrobust		

	coef	std err	t	P> t	[0.025	0.975]
<b>const</b>	-18.5874	1.733	-10.723	0.000	-22.054	-15.121
<b>Year</b>	0.0135	0.001	15.229	0.000	0.012	0.015

<b>Omnibus:</b>	10.035	<b>Durbin-Watson:</b>	0.888
<b>Prob(Omnibus):</b>	0.007	<b>Jarque-Bera (JB):</b>	9.869
<b>Skew:</b>	-0.823	<b>Prob(JB):</b>	0.00719
<b>Kurtosis:</b>	4.024	<b>Cond. No.</b>	9.51e+04

Warnings:

- [1] Standard Errors assume that the covariance matrix of the errors is correctly specified.
- [2] The condition number is large, 9.51e+04. This might indicate that there are strong multicollinearity or other numerical problems.

```
In [12]: X = tf_results['Long Jump Women'].Year.astype(float)
y = tf_results['Long Jump Women'].Result.astype(float)
X = sm.add_constant(X)
LJ_Women = sm.OLS(y, X).fit()
print(LJ_Women.params)
LJ_Women.summary()
```

```
const    -21.278026
Year       0.014137
dtype: float64
```

Out[12]: OLS Regression Results

<b>Dep. Variable:</b>	Result	<b>R-squared:</b>	0.734
<b>Model:</b>	OLS	<b>Adj. R-squared:</b>	0.726
<b>Method:</b>	Least Squares	<b>F-statistic:</b>	90.97
<b>Date:</b>	Wed, 28 Nov 2018	<b>Prob (F-statistic):</b>	5.22e-11
<b>Time:</b>	19:48:52	<b>Log-Likelihood:</b>	11.192
<b>No. Observations:</b>	35	<b>AIC:</b>	-18.38
<b>Df Residuals:</b>	33	<b>BIC:</b>	-15.27
<b>Df Model:</b>	1		
<b>Covariance Type:</b>	nonrobust		

	coef	std err	t	P> t	[0.025	0.975]
<b>const</b>	-21.2780	2.947	-7.221	0.000	-27.273	-15.283
<b>Year</b>	0.0141	0.001	9.538	0.000	0.011	0.017

<b>Omnibus:</b>	0.431	<b>Durbin-Watson:</b>	0.745
<b>Prob(Omnibus):</b>	0.806	<b>Jarque-Bera (JB):</b>	0.537
<b>Skew:</b>	0.228	<b>Prob(JB):</b>	0.765
<b>Kurtosis:</b>	2.601	<b>Cond. No.</b>	1.91e+05

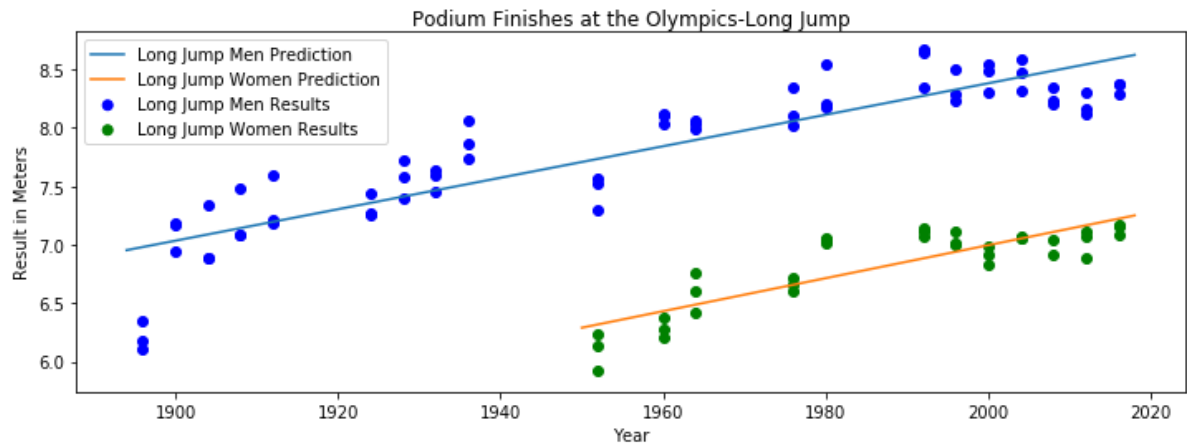
Warnings:

- [1] Standard Errors assume that the covariance matrix of the errors is correctly specified.
- [2] The condition number is large, 1.91e+05. This might indicate that there are strong multicollinearity or other numerical problems.

```

In [13]: time_m = np.linspace(1894, 2018, num = 2000)
time_w = np.linspace(1950, 2018, num = 2000)
Men_pred = -18.587442 + 0.013485 * time_m
Women_pred = -21.278026 + 0.014137 * time_w
plt.figure(figsize=(12, 4))
plt.scatter(tf_results['Long Jump Men'].Year, tf_results['Long Jump Men'].Result, color = 'b')
plt.plot(time_m, Men_pred)
plt.scatter(tf_results['Long Jump Women'].Year, tf_results['Long Jump Women'].Result, color = 'g')
plt.plot(time_w, Women_pred)
plt.xlabel('Year')
plt.ylabel('Result in Meters')
plt.title('Podium Finishes at the Olympics-Long Jump')
plt.legend(['Long Jump Men Prediction', 'Long Jump Women Prediction', 'Long Jump Men Results', 'Long Jump Women Results'])
plt.show()

```



```
In [14]: X = tf_results['Shot Put Men'].Year.astype(float)
y = tf_results['Shot Put Men'].Result.astype(float)
X = sm.add_constant(X)
SP_Men = sm.OLS(y, X).fit()
print(SP_Men.params)
SP_Men.summary()
```

```
const    -142.645216
Year       0.082039
dtype: float64
```

Out[14]: OLS Regression Results

<b>Dep. Variable:</b>	Result	<b>R-squared:</b>	0.930
<b>Model:</b>	OLS	<b>Adj. R-squared:</b>	0.929
<b>Method:</b>	Least Squares	<b>F-statistic:</b>	774.3
<b>Date:</b>	Wed, 28 Nov 2018	<b>Prob (F-statistic):</b>	3.06e-35
<b>Time:</b>	19:48:52	<b>Log-Likelihood:</b>	-80.094
<b>No. Observations:</b>	60	<b>AIC:</b>	164.2
<b>Df Residuals:</b>	58	<b>BIC:</b>	168.4
<b>Df Model:</b>	1		
<b>Covariance Type:</b>	nonrobust		

	coef	std err	t	P> t	[0.025	0.975]
<b>const</b>	-142.6452	5.777	-24.692	0.000	-154.209	-131.081
<b>Year</b>	0.0820	0.003	27.826	0.000	0.076	0.088

<b>Omnibus:</b>	0.027	<b>Durbin-Watson:</b>	0.658
<b>Prob(Omnibus):</b>	0.986	<b>Jarque-Bera (JB):</b>	0.152
<b>Skew:</b>	-0.045	<b>Prob(JB):</b>	0.927
<b>Kurtosis:</b>	2.770	<b>Cond. No.</b>	9.38e+04

Warnings:

- [1] Standard Errors assume that the covariance matrix of the errors is correctly specified.
- [2] The condition number is large, 9.38e+04. This might indicate that there are strong multicollinearity or other numerical problems.

```
In [15]: X = tf_results['Shot Put Women'].Year.astype(float)
y = tf_results['Shot Put Women'].Result.astype(float)
X = sm.add_constant(X)
SP_Women = sm.OLS(y, X).fit()
print(SP_Women.params)
SP_Women.summary()
```

```
const    -103.642325
Year         0.062040
dtype: float64
```

Out[15]: OLS Regression Results

<b>Dep. Variable:</b>	Result	<b>R-squared:</b>	0.411
<b>Model:</b>	OLS	<b>Adj. R-squared:</b>	0.396
<b>Method:</b>	Least Squares	<b>F-statistic:</b>	27.24
<b>Date:</b>	Wed, 28 Nov 2018	<b>Prob (F-statistic):</b>	6.25e-06
<b>Time:</b>	19:48:52	<b>Log-Likelihood:</b>	-72.319
<b>No. Observations:</b>	41	<b>AIC:</b>	148.6
<b>Df Residuals:</b>	39	<b>BIC:</b>	152.1
<b>Df Model:</b>	1		
<b>Covariance Type:</b>	nonrobust		

	coef	std err	t	P> t	[0.025	0.975]
<b>const</b>	-103.6423	23.576	-4.396	0.000	-151.329	-55.955
<b>Year</b>	0.0620	0.012	5.220	0.000	0.038	0.086

<b>Omnibus:</b>	0.304	<b>Durbin-Watson:</b>	0.593
<b>Prob(Omnibus):</b>	0.859	<b>Jarque-Bera (JB):</b>	0.445
<b>Skew:</b>	0.173	<b>Prob(JB):</b>	0.800
<b>Kurtosis:</b>	2.625	<b>Cond. No.</b>	2.07e+05

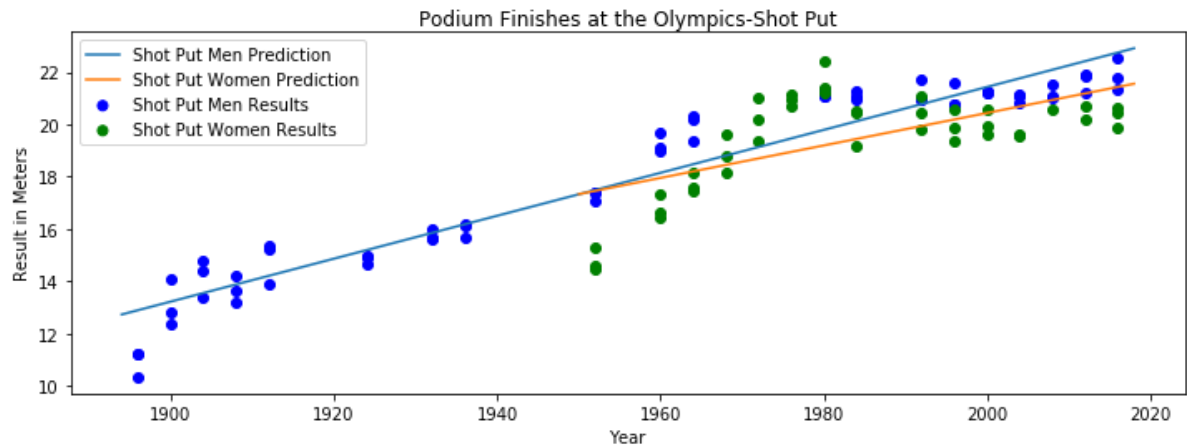
Warnings:

- [1] Standard Errors assume that the covariance matrix of the errors is correctly specified.
- [2] The condition number is large, 2.07e+05. This might indicate that there are strong multicollinearity or other numerical problems.

```

In [16]: time_m = np.linspace(1894, 2018, num = 2000)
time_w = np.linspace(1950, 2018, num = 2000)
Men_pred = -142.645216 + 0.082039 * time_m
Women_pred = -103.642325 + 0.062040 * time_w
plt.figure(figsize=(12, 4))
plt.scatter(tf_results['Shot Put Men'].Year, tf_results['Shot Put Men'].Result,
            color = 'b')
plt.plot(time_m, Men_pred)
plt.scatter(tf_results['Shot Put Women'].Year, tf_results['Shot Put Women'].Result,
            color = 'g')
plt.plot(time_w, Women_pred)
plt.xlabel('Year')
plt.ylabel('Result in Meters')
plt.title('Podium Finishes at the Olympics-Shot Put')
plt.legend(['Shot Put Men Prediction', 'Shot Put Women Prediction', 'Shot Put Men Results', 'Shot Put Women Results'])
plt.show()

```



```
In [17]: #Men's 1500M
men_1500m = pd.read_csv('Men 1500M.csv')
X = men_1500m['Year'].astype(float)
y = men_1500m['seconds'].astype(float)
X = sm.add_constant(X)
Men_1500M = sm.OLS(y, X).fit()
print(Men_1500M.params)
Men_1500M.summary()
```

```
const      846.904980
Year       -0.316465
dtype: float64
```

Out[17]: OLS Regression Results

<b>Dep. Variable:</b>	seconds	<b>R-squared:</b>	0.682
<b>Model:</b>	OLS	<b>Adj. R-squared:</b>	0.677
<b>Method:</b>	Least Squares	<b>F-statistic:</b>	164.8
<b>Date:</b>	Wed, 28 Nov 2018	<b>Prob (F-statistic):</b>	8.04e-21
<b>Time:</b>	19:48:53	<b>Log-Likelihood:</b>	-275.71
<b>No. Observations:</b>	79	<b>AIC:</b>	555.4
<b>Df Residuals:</b>	77	<b>BIC:</b>	560.2
<b>Df Model:</b>	1		
<b>Covariance Type:</b>	nonrobust		

	coef	std err	t	P> t	[0.025	0.975]
<b>const</b>	846.9050	48.285	17.540	0.000	750.756	943.053
<b>Year</b>	-0.3165	0.025	-12.837	0.000	-0.366	-0.267

<b>Omnibus:</b>	48.286	<b>Durbin-Watson:</b>	0.292
<b>Prob(Omnibus):</b>	0.000	<b>Jarque-Bera (JB):</b>	132.075
<b>Skew:</b>	2.138	<b>Prob(JB):</b>	2.09e-29
<b>Kurtosis:</b>	7.673	<b>Cond. No.</b>	1.05e+05

Warnings:

- [1] Standard Errors assume that the covariance matrix of the errors is correctly specified.
- [2] The condition number is large, 1.05e+05. This might indicate that there are strong multicollinearity or other numerical problems.

```
In [18]: #Women's 1500M
women_1500m = pd.read_csv('Women 1500M.csv')
women_1500m.head(3)
X = women_1500m['Year'].astype(float)
y = women_1500m['seconds'].astype(float)
X = sm.add_constant(X)
Women_1500M = sm.OLS(y, X).fit()
print(Women_1500M.params)
Women_1500M.summary()
```

```
const    40.274876
Year      0.101555
dtype: float64
```

Out[18]: OLS Regression Results

<b>Dep. Variable:</b>	seconds	<b>R-squared:</b>	0.110
<b>Model:</b>	OLS	<b>Adj. R-squared:</b>	0.081
<b>Method:</b>	Least Squares	<b>F-statistic:</b>	3.722
<b>Date:</b>	Wed, 28 Nov 2018	<b>Prob (F-statistic):</b>	0.0632
<b>Time:</b>	19:48:53	<b>Log-Likelihood:</b>	-90.462
<b>No. Observations:</b>	32	<b>AIC:</b>	184.9
<b>Df Residuals:</b>	30	<b>BIC:</b>	187.9
<b>Df Model:</b>	1		
<b>Covariance Type:</b>	nonrobust		

	coef	std err	t	P> t	[0.025	0.975]
<b>const</b>	40.2749	104.971	0.384	0.704	-174.104	254.653
<b>Year</b>	0.1016	0.053	1.929	0.063	-0.006	0.209

<b>Omnibus:</b>	8.608	<b>Durbin-Watson:</b>	0.832
<b>Prob(Omnibus):</b>	0.014	<b>Jarque-Bera (JB):</b>	2.354
<b>Skew:</b>	-0.147	<b>Prob(JB):</b>	0.308
<b>Kurtosis:</b>	1.704	<b>Cond. No.</b>	2.80e+05

Warnings:

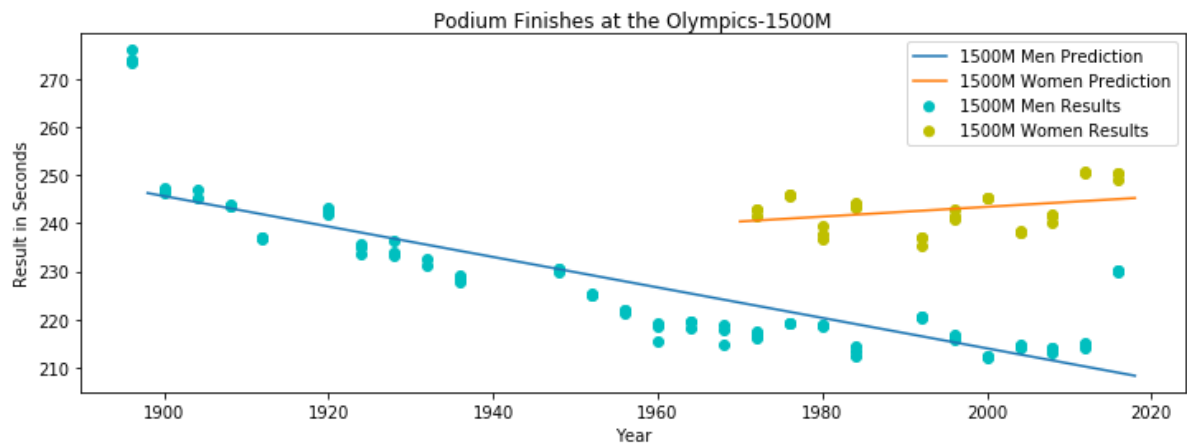
- [1] Standard Errors assume that the covariance matrix of the errors is correctly specified.
- [2] The condition number is large, 2.8e+05. This might indicate that there are strong multicollinearity or other numerical problems.



```

In [19]: time_m = np.linspace(1898, 2018, num = 2000)
time_w = np.linspace(1970, 2018, num = 2000)
Men_pred = 846.904980 - 0.316465 * time_m
Women_pred = 40.274876 + 0.101555 * time_w
plt.figure(figsize=(12, 4))
plt.scatter(men_1500m['Year'], men_1500m['seconds'], color = 'c')
plt.plot(time_m, Men_pred)
plt.scatter(women_1500m['Year'], women_1500m['seconds'], color = 'y')
plt.plot(time_w, Women_pred)
plt.xlabel('Year')
plt.ylabel('Result in Seconds')
plt.title('Podium Finishes at the Olympics-1500M')
plt.legend(['1500M Men Prediction', '1500M Women Prediction', '1500M Men Results', '1500M Women Results'])
plt.show()

```



## Diagonostics of the Simple Linear Models

We now will look at the residuals to try to determine if we can make the assumption that the error is normally distributed with mean 0 and variance  $\sigma^2$ , and run some diagnostic tests to determine whether the error is constant.

```

In [20]: #Residuals for 100M
res_100_m = tf_results['100M Men'].Result.astype(float) - \
(37.670494 - 0.013918 * tf_results['100M Men'].Year.astype(float))
res_100_w = tf_results['100M Women'].Result.astype(float) - \
(39.259190 - 0.014167 * tf_results['100M Women'].Year.astype(float))

#Residuals for 200M
res_200_m = tf_results['200M Men'].Result.astype(float) - \
(67.599404 - 0.023895 * tf_results['200M Men'].Year.astype(float))
res_200_w = tf_results['200M Women'].Result.astype(float) - \
(94.100684 - 0.036020 * tf_results['200M Women'].Year.astype(float))

#Residuals for Long Jump
res_lj_m = tf_results['Long Jump Men'].Result.astype(float) - \
(-18.587442 + 0.013485 * tf_results['Long Jump Men'].Year.astype(float))
res_lj_w = tf_results['Long Jump Women'].Result.astype(float) - \
(-21.278026 + 0.014137 * tf_results['Long Jump Women'].Year.astype(float))

#Residuals for Shot Put
res_sp_m = tf_results['Shot Put Men'].Result.astype(float) - \
(-142.645216 + 0.082039 * tf_results['Shot Put Men'].Year.astype(float))
res_sp_w = tf_results['Shot Put Women'].Result.astype(float) - \
(-103.642325 + 0.062040 * tf_results['Shot Put Women'].Year.astype(float))

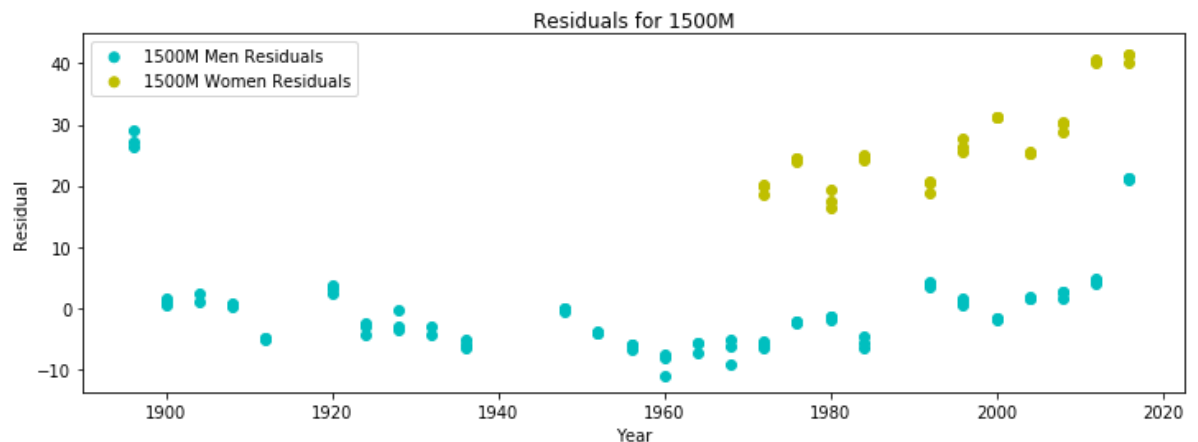
#Residuals for 1500M
res_1500_m = men_1500m['seconds'] - (846.904980 - 0.316465 * men_1500m['Year'])
res_1500_w = women_1500m['seconds'] - (846.904980 - 0.316465 * women_1500m['Year'])

```

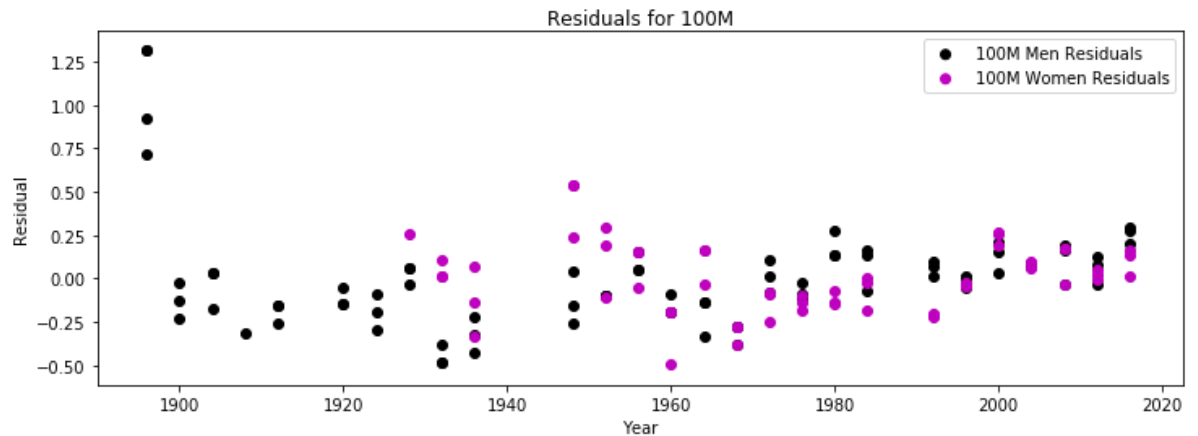
```

In [21]: plt.figure(figsize=(12, 4))
plt.scatter(men_1500m['Year'], res_1500_m, color = 'c')
plt.scatter(women_1500m['Year'], res_1500_w, color = 'y')
plt.xlabel('Year')
plt.ylabel('Residual')
plt.title('Residuals for 1500M')
plt.legend(['1500M Men Residuals', '1500M Women Residuals'])
plt.show()

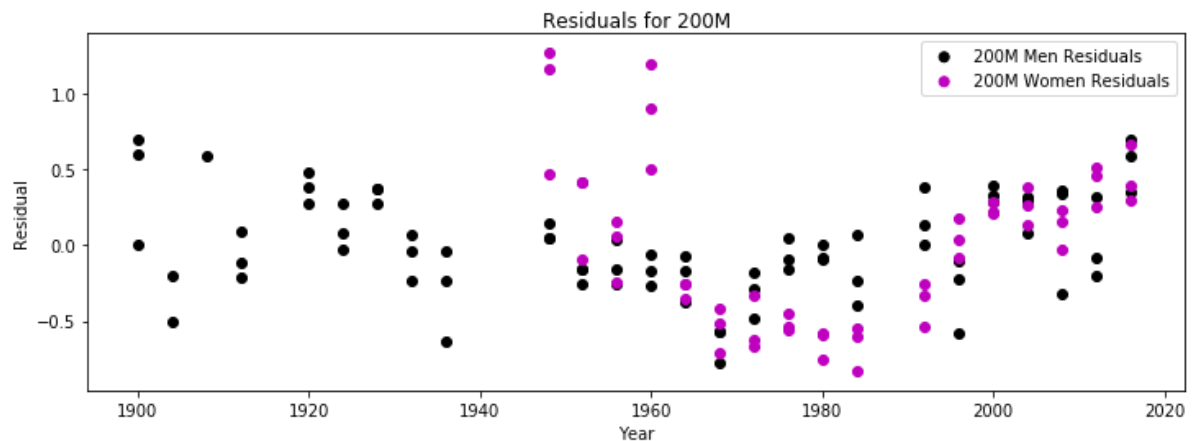
```



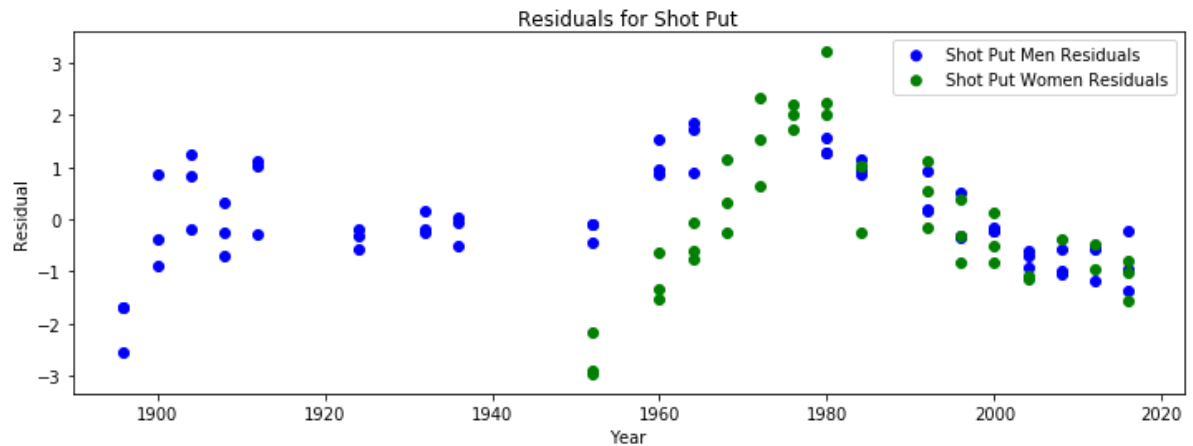
```
In [22]: plt.figure(figsize=(12, 4))
plt.scatter(tf_results['100M Men'].Year, res_100_m, color = 'k')
plt.scatter(tf_results['100M Women'].Year, res_100_w, color = 'm')
plt.xlabel('Year')
plt.ylabel('Residual')
plt.title('Residuals for 100M')
plt.legend(['100M Men Residuals', '100M Women Residuals'])
plt.show()
```



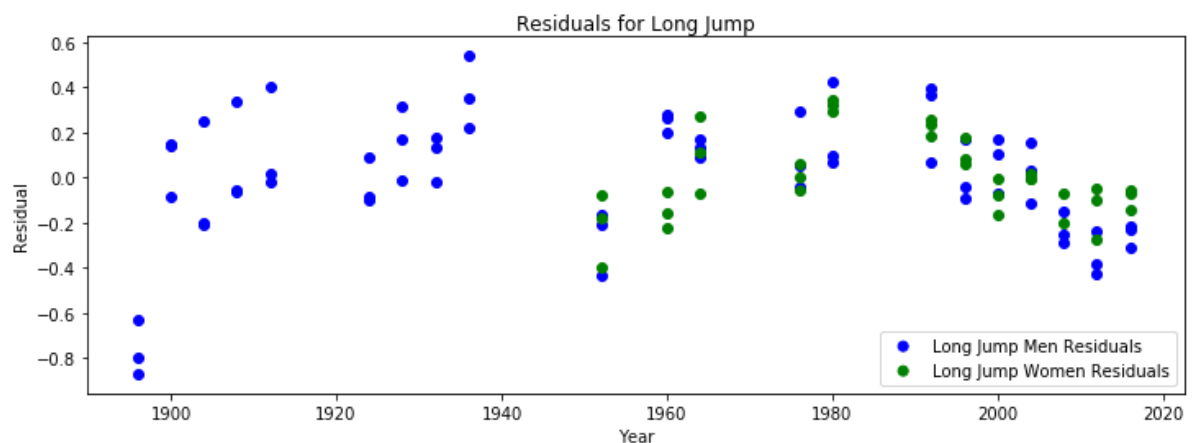
```
In [23]: plt.figure(figsize=(12, 4))
plt.scatter(tf_results['200M Men'].Year, res_200_m, color = 'k')
plt.scatter(tf_results['200M Women'].Year, res_200_w, color = 'm')
plt.xlabel('Year')
plt.ylabel('Residual')
plt.title('Residuals for 200M')
plt.legend(['200M Men Residuals', '200M Women Residuals'])
plt.show()
```



```
In [24]: plt.figure(figsize=(12, 4))
plt.scatter(tf_results['Shot Put Men'].Year, res_sp_m, color = 'b')
plt.scatter(tf_results['Shot Put Women'].Year, res_sp_w, color = 'g')
plt.xlabel('Year')
plt.ylabel('Residual')
plt.title('Residuals for Shot Put')
plt.legend(['Shot Put Men Residuals', 'Shot Put Women Residuals'])
plt.show()
```



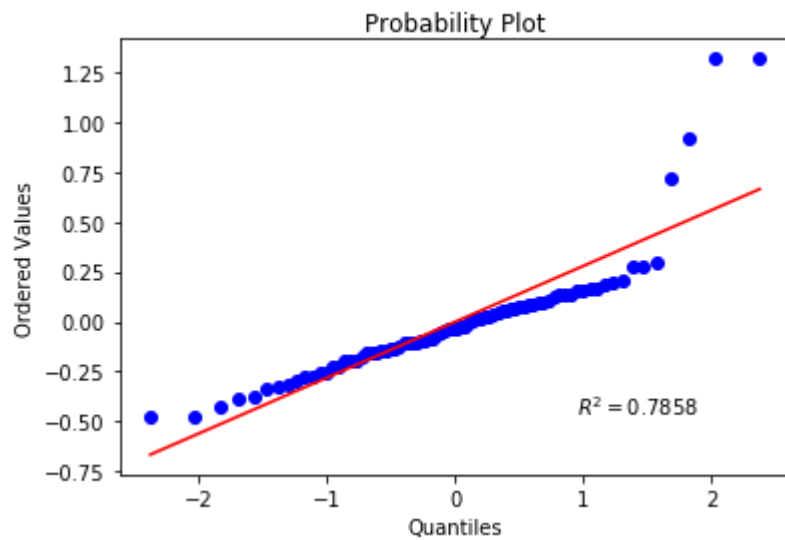
```
In [25]: plt.figure(figsize=(12, 4))
plt.scatter(tf_results['Long Jump Men'].Year, res_lj_m, color = 'b')
plt.scatter(tf_results['Long Jump Women'].Year, res_lj_w, color = 'g')
plt.xlabel('Year')
plt.ylabel('Residual')
plt.title('Residuals for Long Jump')
plt.legend(['Long Jump Men Residuals', 'Long Jump Women Residuals'])
plt.show()
```



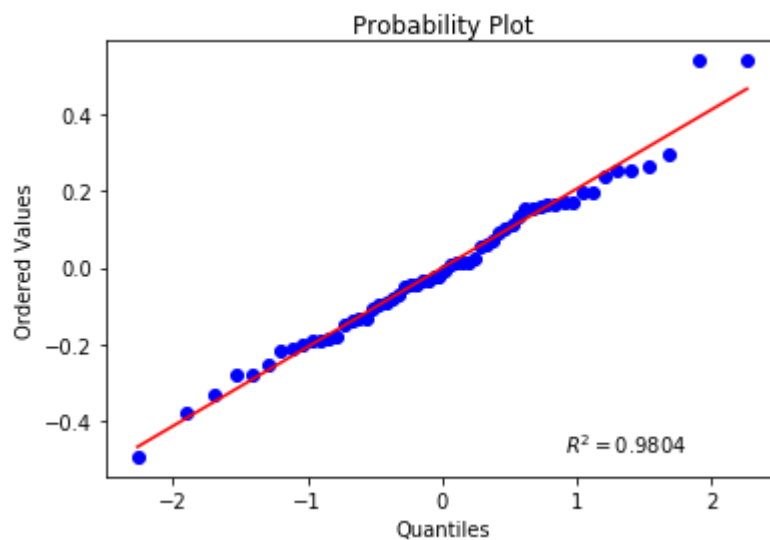
We are somewhat suspicious of whether or not the error variance is constant or normally distributed. We also believe that we may need higher order terms or a Box-Cox transformation to fit our regression model better than it currently fits the data.

## Normal Probability Plots

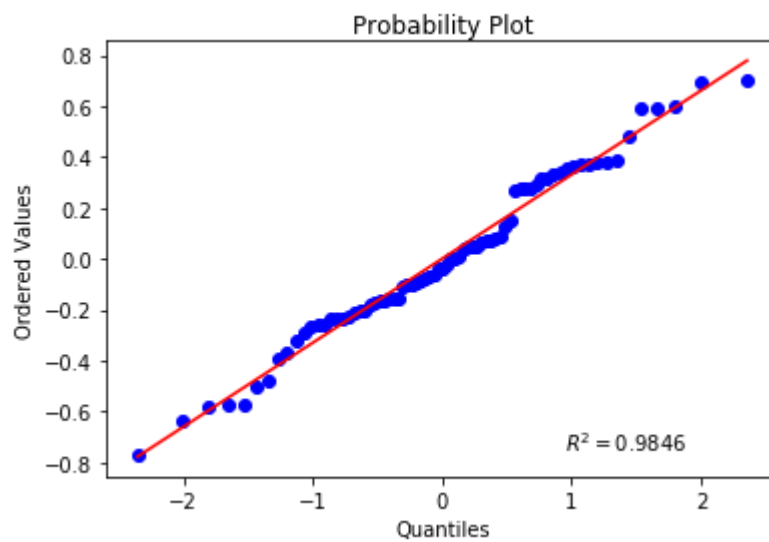
```
In [26]: #100M Men  
res_prob_100m_m = stats.probplot(res_100_m, plot= plt)  
plt.show()
```



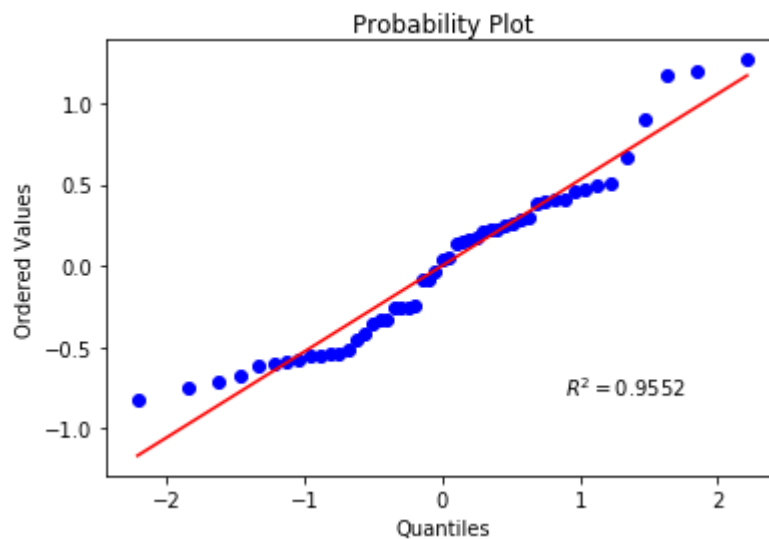
```
In [27]: #100M Women  
res_prob_100m_w = stats.probplot(res_100_w, plot= plt)  
plt.show()
```



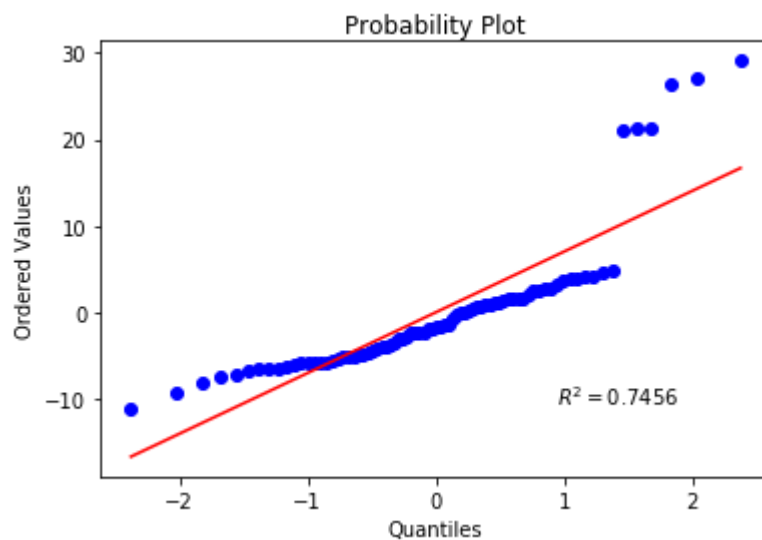
```
In [28]: #200M Men  
res_prob_200m_m = stats.probplot(res_200_m, plot= plt)  
plt.show()
```



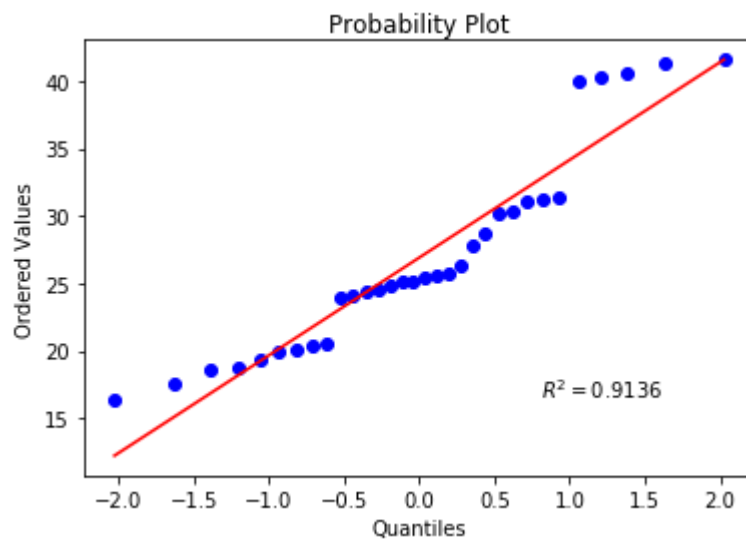
```
In [29]: #200M Women  
res_prob_200m_w = stats.probplot(res_200_w, plot= plt)  
plt.show()
```



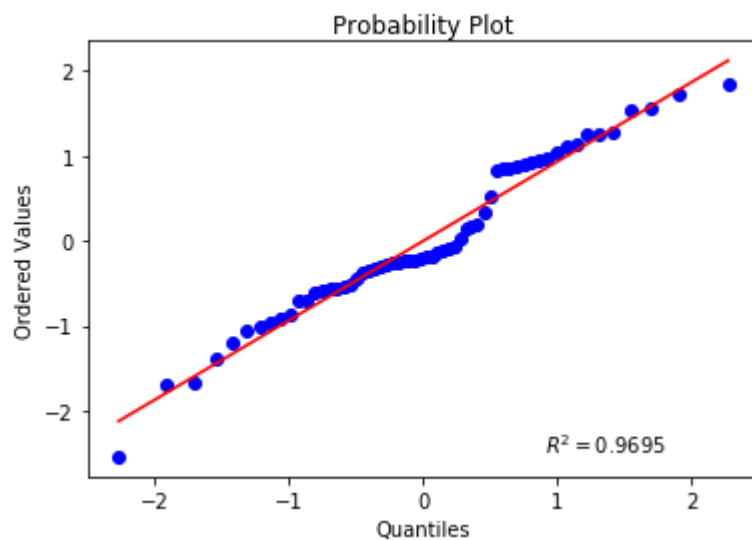
```
In [30]: #1500M Men  
res_prob_1500m_m = stats.probplot(res_1500_m, plot= plt)  
plt.show()
```



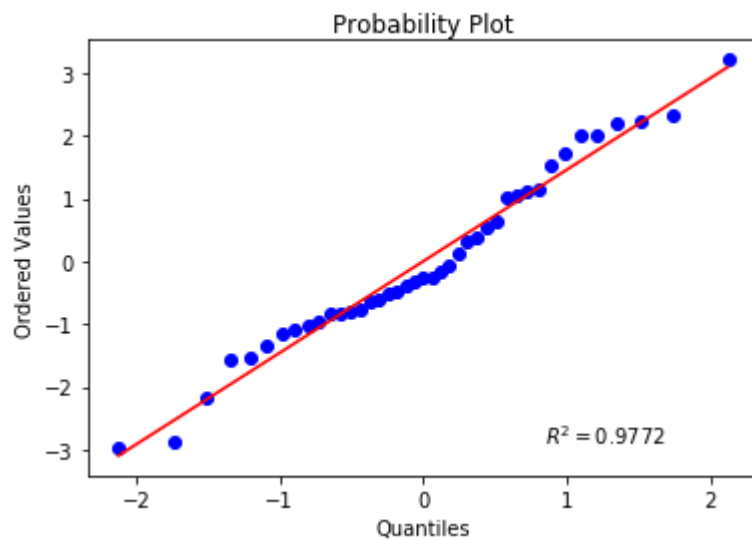
```
In [31]: #1500M Women  
res_prob_1500m_w = stats.probplot(res_1500_w, plot= plt)  
plt.show()
```



```
In [32]: #Men's Shot Put  
res_prob_sp_m = stats.probplot(res_sp_m, plot= plt)  
plt.show()
```

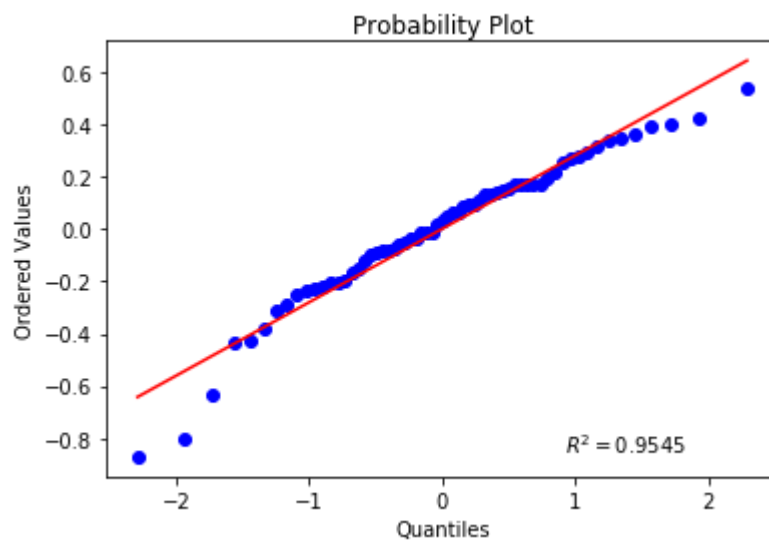


```
In [33]: #Women's Shot Put  
res_prob_sp_w = stats.probplot(res_sp_w, plot= plt)  
plt.show()
```

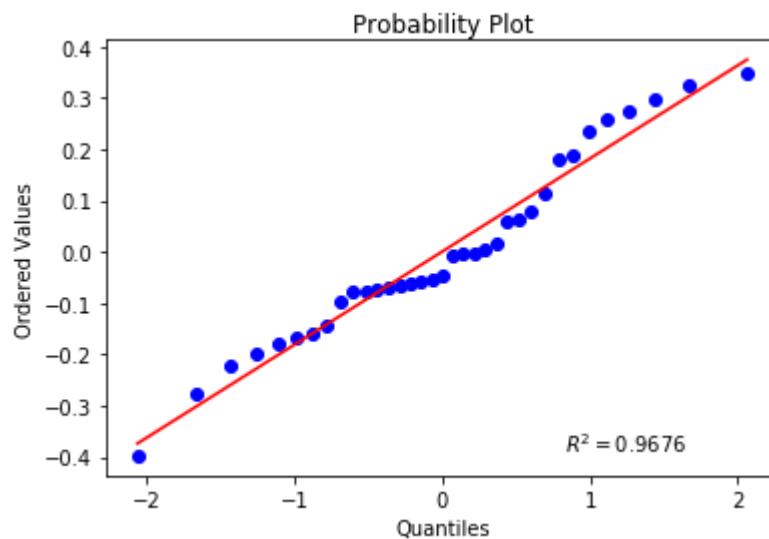




```
In [34]: #Men's Long Jump  
res_prob_lj_m = stats.probplot(res_lj_m, plot= plt)  
plt.show()
```



```
In [35]: #Women's Long Jump  
res_prob_lj_w = stats.probplot(res_lj_w, plot= plt)  
plt.show()
```



It appears that there are some departures in normality in every single normal probability plot in varying degrees of seriousness except for Shot Put, Women's Long Jump, and Men's 200M dash. This again strongly implies that the error variance is not constant for this model and that we may need to utilize some higher order terms to better fit the data.

## Bruesch-Pagan Tests for Constant Error Variance

We will also run a Bruesch-Pagan test to test form departures in constancy of the error terms. Remember, this test is carried out as follows:

We regress the function:  $\log(\sigma_i^2) = \gamma_0 + \gamma_1 x_{i1}$

$$H_0 : \gamma_1 = 0$$

$$H_1 : \gamma_1 \neq 0$$

$$\text{Test statistic: } \frac{SSR^*/2}{(SSE/n)^2}$$

Where  $SSR^*$  is the regression sum of squares for the log regression of the residuals.

Critical value for the rejection region would be  $\chi_{\alpha/2;1}^2$ , and we will test at an  $\alpha = 0.05$  level.

```
In [36]: from statsmodels.stats import diagnostic as dn
bp_100_m = dn.het_breuschpagan(model_100M_Men.resid, model_100M_Men.model.exog)
print("The Breusch-Pagan test yields a p-value of: ", bp_100_m[3], ".")

('The Breusch-Pagan test yields a p-value of: ', 0.00066013566230912351, '.')
```

```
In [37]: bp_100_w = dn.het_breuschpagan(model_100M_Women.resid, model_100M_Women.model.exog)
print("The Breusch-Pagan test yields a p-value of: ", bp_100_w[3], ".")

('The Breusch-Pagan test yields a p-value of: ', 0.021171634781806045, '.')
```

```
In [38]: bp_200_m = dn.het_breuschpagan(model_200M_Men.resid, model_200M_Men.model.exog)
print("The Breusch-Pagan test yields a p-value of: ", bp_200_m[3], ".")

('The Breusch-Pagan test yields a p-value of: ', 0.82010810471300388, '.')
```

```
In [39]: bp_200_w = dn.het_breuschpagan(model_200M_Women.resid, model_200M_Women.model.exog)
print("The Breusch-Pagan test yields a p-value of: ", bp_200_w[3], ".")

('The Breusch-Pagan test yields a p-value of: ', 0.0047673629262334448, '.')
```

```
In [40]: bp_1500_m = dn.het_breuschpagan(Men_1500M.resid, Men_1500M.model.exog)
print("The Breusch-Pagan test yields a p-value of: ", bp_1500_m[3], ".")

('The Breusch-Pagan test yields a p-value of: ', 0.25447848550659419, '.')
```

```
In [41]: bp_1500_w = dn.het_breuschpagan(Women_1500M.resid, Women_1500M.model.exog)
print("The Breusch-Pagan test yields a p-value of: ", bp_1500_w[3],".")

('The Breusch-Pagan test yields a p-value of: ', 0.10160731583559358, '.')
```

```
In [42]: bp_lj_m = dn.het_breuschpagan(LJ_Men.resid, LJ_Men.model.exog)
print("The Breusch-Pagan test yields a p-value of: ", bp_lj_m[3],".")

('The Breusch-Pagan test yields a p-value of: ', 0.04485189184383321, '.')
```

```
In [43]: bp_lj_w = dn.het_breuschpagan(LJ_Women.resid, LJ_Women.model.exog)
print("The Breusch-Pagan test yields a p-value of: ", bp_lj_w[3],".")

('The Breusch-Pagan test yields a p-value of: ', 0.073930238062038287, '.')
```

```
In [44]: bp_sp_m = dn.het_breuschpagan(SP_Men.resid, SP_Men.model.exog)
print("The Breusch-Pagan test yields a p-value of: ", bp_sp_m[3],".")

('The Breusch-Pagan test yields a p-value of: ', 0.25290089457670523, '.')
```

```
In [45]: bp_sp_w = dn.het_breuschpagan(SP_Women.resid, SP_Women.model.exog)
print("The Breusch-Pagan test yields a p-value of: ", bp_sp_w[3],".")

('The Breusch-Pagan test yields a p-value of: ', 0.0072550835390873597, '.')
```

Based on the Bruesch-Pagan tests we have run, it appears that the only track results that we have regressed that do not have constant error variance are the 100M Results, the Women's 200M, the Men's Long Jump, and the Women's Shot Put at an significance level of  $\alpha = 0.05$ .

## Goodness of Fit Tests on First Order Model:

Since we have repeat observations at each independent predictor variable, we should be able to run an F test for lack of fit. Remember, the F-test for lack of fit is carried out as follows:

$$H_0 : E(Y) = \beta_0 + \beta_1 X$$

$$H_1 : E(Y) \neq \beta_0 + \beta_1 X$$

The test statistic here would be  $F^* = \frac{SSLF}{c-2} / \frac{SSPE}{n-c} = \frac{MSLF}{MSPE}$ , where  $c$  is the number of years with replicates and  $n = \sum_{j=1}^c n_j$  and  $n$  is the number of observations.

Our decision rule is:

$F^* \leq F(1 - \alpha; c - 2, n - c)$  then we conclude that our simple linear model is appropriate.

$F^* > F(1 - \alpha; c - 2, n - c)$  then we conclude that our simple linear model is not appropriate for our data.

We will control the error at  $\alpha = 0.05$ .

## ALL THE CODE WILL BE IN R FOR THESE TESTS.

All of the results are on Github. In all cases we conclude that the simple linear models are not appropriate. Clearly, we should use some Box-Cox transformations on our results in order to better predict the results.