

## Q2 Regression analysis:

In [11]:

```
#Importing Libraries
import pandas as pd
import numpy as np
from sklearn import linear_model
import statsmodels.api as sm
import matplotlib.pyplot as plt
import math
```

### Loading the data and looking at the first 15 rows

In [12]:

```
df = pd.read_csv("Comp1_IE529.csv", header = None)
df.head(15)
```

Out[12]:

|    | 0     | 1    |
|----|-------|------|
| 0  | 37.5  | 6.4  |
| 1  | 51.5  | 10.2 |
| 2  | 61.3  | 12.4 |
| 3  | 61.3  | 13.0 |
| 4  | 63.6  | 13.2 |
| 5  | 66.1  | 13.0 |
| 6  | 70.0  | 12.7 |
| 7  | 92.7  | 13.9 |
| 8  | 90.5  | 15.5 |
| 9  | 90.5  | 15.8 |
| 10 | 94.8  | 15.8 |
| 11 | 97.0  | 16.8 |
| 12 | 97.0  | 17.1 |
| 13 | 97.0  | 17.8 |
| 14 | 102.0 | 14.8 |

In [13]:

```
X = np.array(df)
X = np.insert(X, 0, 1, axis=1)  #Inserting a column of 1s so that we can model a simple linear regression
```

## First fit a simple linear regression model to the data

$$\hat{y} = \hat{\beta}_0 + \hat{\beta}_1 x$$

We will use statsmodel.api library in python and call the OLS function which stands for ordinary least squares to fit a simple linear regression model to the data with minimum sum of squared errors

In [14]:

```
model = sm.OLS(X[:,2],X[:, :2]).fit()
model.summary()
```

Out[14]:

OLS Regression Results

|                          |                  |                            |          |
|--------------------------|------------------|----------------------------|----------|
| <b>Dep. Variable:</b>    | y                | <b>R-squared:</b>          | 0.831    |
| <b>Model:</b>            | OLS              | <b>Adj. R-squared:</b>     | 0.828    |
| <b>Method:</b>           | Least Squares    | <b>F-statistic:</b>        | 294.9    |
| <b>Date:</b>             | Thu, 09 Nov 2017 | <b>Prob (F-statistic):</b> | 7.76e-25 |
| <b>Time:</b>             | 11:08:52         | <b>Log-Likelihood:</b>     | -101.07  |
| <b>No. Observations:</b> | 62               | <b>AIC:</b>                | 206.1    |
| <b>Df Residuals:</b>     | 60               | <b>BIC:</b>                | 210.4    |
| <b>Df Model:</b>         | 1                |                            |          |
| <b>Covariance Type:</b>  | nonrobust        |                            |          |

|              | coef   | std err | t      | P> t  | [0.025 | 0.975] |
|--------------|--------|---------|--------|-------|--------|--------|
| <b>const</b> | 5.6796 | 0.580   | 9.784  | 0.000 | 4.518  | 6.841  |
| <b>x1</b>    | 0.0995 | 0.006   | 17.174 | 0.000 | 0.088  | 0.111  |

|                       |        |                          |       |
|-----------------------|--------|--------------------------|-------|
| <b>Omnibus:</b>       | 1.396  | <b>Durbin-Watson:</b>    | 1.071 |
| <b>Prob(Omnibus):</b> | 0.497  | <b>Jarque-Bera (JB):</b> | 1.300 |
| <b>Skew:</b>          | -0.219 | <b>Prob(JB):</b>         | 0.522 |
| <b>Kurtosis:</b>      | 2.443  | <b>Cond. No.</b>         | 365.  |

The coefficients obtained are

$$\hat{\beta}_0 = 5.6796 \quad \hat{\beta}_1 = 0.0995$$

Determine the sum of squared error for this model

In [15]:

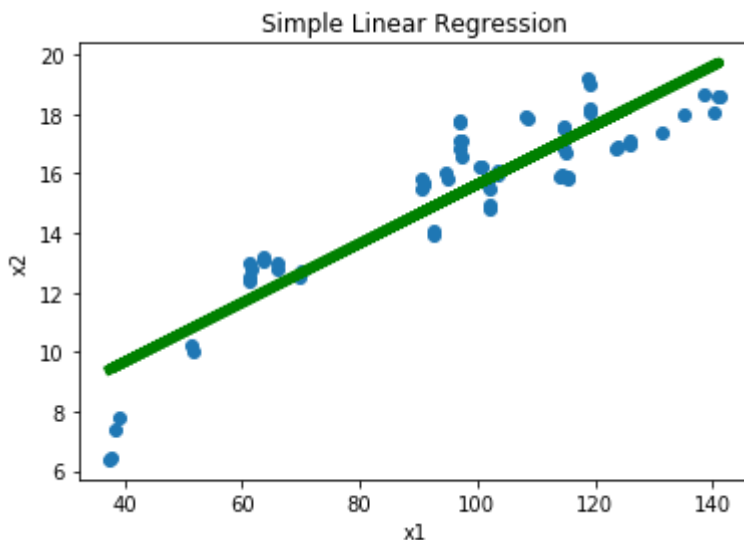
```
ssr = []
ssr.append(model.ssr)
print('\033[1m' + 'Sum of squared residuals= ' + '\033[0m', model.ssr)
```

Sum of squared residuals= 94.5794789741

## Plot the fitted model overlaying the scatter plot of our original dataset

In [16]:

```
x=model.params[0] + model.params[1] * X[:,1]
plt.scatter(X[:,1], X[:,2])
plt.plot(X[:,1], x, color = 'g', linewidth=5.0)
plt.title("Simple Linear Regression")
plt.xlabel("x1")
plt.ylabel("x2")
plt.show()
```



## Fit a second order polynomial linear regression model to the data

$$\hat{y} = \hat{\beta}_0 + \hat{\beta}_1 x + \hat{\beta}_2 x^2$$

In [17]:

```
df['a'] = df[0]**2 #Adding a square term to the dataframe
#Next three lines to move output variable to the end of dataframe
cols = list(df.columns.values) #Make a list of all of the columns in the df
cols.pop(cols.index(1)) #Remove 1 from list
df = df[cols+[1]] #Create new dataframe with columns in the order you want
X = np.array(df) #Converting to numpy array to use further
X = np.insert(X, 0, 1, axis=1) #Inserting a column of 1s for constant term
```

In [18]:

```
model = sm.OLS(X[:,3],X[:,3]).fit()    #Fitting the model
model.summary()
```

Out[18]:

OLS Regression Results

|                          |                  |                            |          |
|--------------------------|------------------|----------------------------|----------|
| <b>Dep. Variable:</b>    | y                | <b>R-squared:</b>          | 0.909    |
| <b>Model:</b>            | OLS              | <b>Adj. R-squared:</b>     | 0.906    |
| <b>Method:</b>           | Least Squares    | <b>F-statistic:</b>        | 293.6    |
| <b>Date:</b>             | Thu, 09 Nov 2017 | <b>Prob (F-statistic):</b> | 2.15e-31 |
| <b>Time:</b>             | 11:09:11         | <b>Log-Likelihood:</b>     | -81.969  |
| <b>No. Observations:</b> | 62               | <b>AIC:</b>                | 169.9    |
| <b>Df Residuals:</b>     | 59               | <b>BIC:</b>                | 176.3    |
| <b>Df Model:</b>         | 2                |                            |          |
| <b>Covariance Type:</b>  | nonrobust        |                            |          |

|              | coef    | std err | t      | P> t  | [0.025 | 0.975] |
|--------------|---------|---------|--------|-------|--------|--------|
| <b>const</b> | -1.5254 | 1.104   | -1.382 | 0.172 | -3.734 | 0.683  |
| <b>x1</b>    | 0.2791  | 0.026   | 10.858 | 0.000 | 0.228  | 0.331  |
| <b>x2</b>    | -0.0010 | 0.000   | -7.088 | 0.000 | -0.001 | -0.001 |

|                       |        |                          |          |
|-----------------------|--------|--------------------------|----------|
| <b>Omnibus:</b>       | 2.658  | <b>Durbin-Watson:</b>    | 1.398    |
| <b>Prob(Omnibus):</b> | 0.265  | <b>Jarque-Bera (JB):</b> | 1.608    |
| <b>Skew:</b>          | -0.104 | <b>Prob(JB):</b>         | 0.448    |
| <b>Kurtosis:</b>      | 2.239  | <b>Cond. No.</b>         | 1.05e+05 |

**The coefficients obtained are**

$\hat{\beta}_0 = -1.5254$   $\hat{\beta}_1 = 0.2791$   $\hat{\beta}_2 = -0.0010$

**The sum of squared error is for this model**

In [19]:

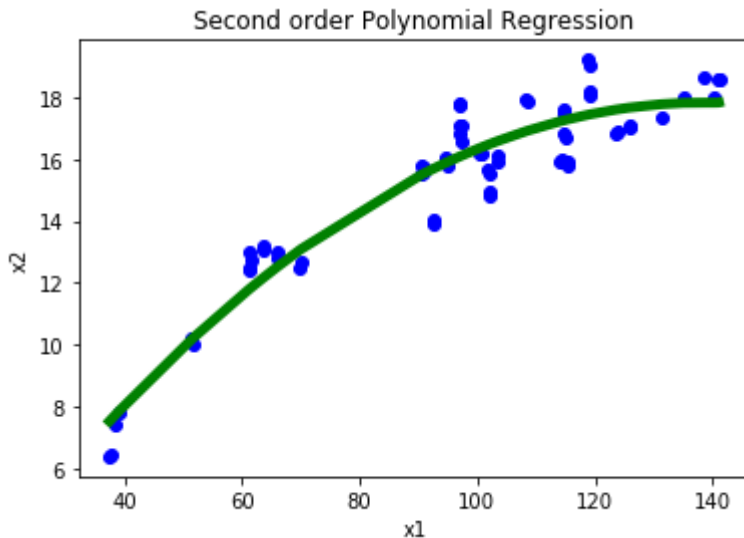
```
ssr.append(model.ssr)
print('Sum of squared residuals= ' + str(model.ssr))
```

Sum of squared residuals= 51.0810573979

**Plot the fitted model overlaying the scatter plot of our original dataset**

In [21]:

```
x=model.params[0] + model.params[1] * X[:,1] + model.params[2] * X[:,2]
plt.scatter(X[:,1], X[:,3], color = 'b')
plt.plot(np.sort(X[:,1]), np.sort(x), color = 'g', linewidth=5.0)
plt.title("Second order Polynomial Regression")
plt.xlabel("x1")
plt.ylabel("x2")
plt.show()
```



The second order polynomial regression model fits the data points much better and causes the sum of squared errors to fall considerably as well. From the plot and ssr, it seems like the second order fits the data better than the simple linear regression model

We will further examine higher order polynomials before coming to a conclusion about our preferred model

**Fit a third order polynomial linear regression model to the data**

$$\hat{y} = \hat{\beta}_0 + \hat{\beta}_1 x + \hat{\beta}_2 x^2 + \hat{\beta}_3 x^3$$

In [22]:

```
df['b'] = df[0]**3 #Adding a square term to the dataframe
#Next three lines to move output variable to end of dataframe
cols = list(df.columns.values) #Make a list of all of the columns in the df
cols.pop(cols.index(1)) #Remove 1 from list
df = df[cols+[1]] #Create new dataframe with columns in the order you want
X = np.array(df) #Converting to numpy array to use further
X = np.insert(X, 0, 1, axis=1) #Inserting a column of 1s for constant term
```

In [23]:

```
model = sm.OLS(X[:,4],X[:,4]).fit()
model.summary()
```

Out[23]:

## OLS Regression Results

|                          |                  |                            |          |
|--------------------------|------------------|----------------------------|----------|
| <b>Dep. Variable:</b>    | y                | <b>R-squared:</b>          | 0.920    |
| <b>Model:</b>            | OLS              | <b>Adj. R-squared:</b>     | 0.916    |
| <b>Method:</b>           | Least Squares    | <b>F-statistic:</b>        | 223.4    |
| <b>Date:</b>             | Thu, 09 Nov 2017 | <b>Prob (F-statistic):</b> | 8.10e-32 |
| <b>Time:</b>             | 11:13:03         | <b>Log-Likelihood:</b>     | -77.744  |
| <b>No. Observations:</b> | 62               | <b>AIC:</b>                | 163.5    |
| <b>Df Residuals:</b>     | 58               | <b>BIC:</b>                | 172.0    |
| <b>Df Model:</b>         | 3                |                            |          |
| <b>Covariance Type:</b>  | nonrobust        |                            |          |

|              | coef      | std err  | t      | P> t  | [0.025   | 0.975]   |
|--------------|-----------|----------|--------|-------|----------|----------|
| <b>const</b> | -9.4058   | 2.901    | -3.243 | 0.002 | -15.212  | -3.599   |
| <b>x1</b>    | 0.5894    | 0.109    | 5.391  | 0.000 | 0.371    | 0.808    |
| <b>x2</b>    | -0.0047   | 0.001    | -3.687 | 0.001 | -0.007   | -0.002   |
| <b>x3</b>    | 1.348e-05 | 4.63e-06 | 2.910  | 0.005 | 4.21e-06 | 2.27e-05 |

|                       |       |                          |          |
|-----------------------|-------|--------------------------|----------|
| <b>Omnibus:</b>       | 0.477 | <b>Durbin-Watson:</b>    | 1.542    |
| <b>Prob(Omnibus):</b> | 0.788 | <b>Jarque-Bera (JB):</b> | 0.377    |
| <b>Skew:</b>          | 0.186 | <b>Prob(JB):</b>         | 0.828    |
| <b>Kurtosis:</b>      | 2.914 | <b>Cond. No.</b>         | 3.47e+07 |

The coefficients obtained are  $\begin{aligned}$

$\hat{\beta}_0 = -9.4058$   $\hat{\beta}_1 = 0.5894$   $\hat{\beta}_2 = -0.0047$   $\hat{\beta}_3 = 1.348e^{-05}$   $\end{aligned}$

Lets now see what the sum of squared error is for this model

In [25]:

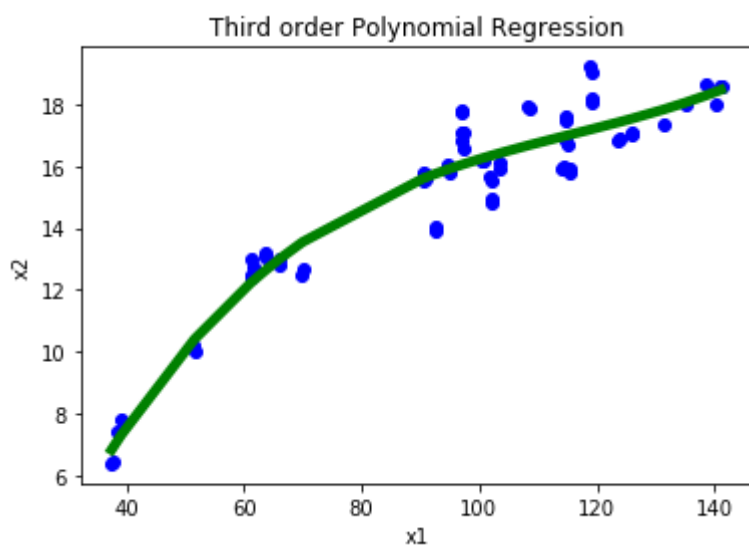
```
ssr.append(model.ssr)
print('\033[1m'+ 'Sum of squared residuals= ' + '\033[0m',model.ssr)
```

Sum of squared residuals= 44.5729109596

## Plot the fitted model overlaying the scatter plot of the original dataset

In [26]:

```
x=model.params[0] + model.params[1] * X[:,1] + model.params[2] * X[:,2] + model.params[3] * X[:,3]
plt.scatter(X[:,1], X[:,4], color = 'b')
plt.plot(np.sort(X[:,1]), np.sort(x), color = 'g', linewidth=5.0)
plt.title("Third order Polynomial Regression")
plt.xlabel("x1")
plt.ylabel("x2")
plt.show()
```



**The third order polynomial fits better but the sum of squared residuals has not reduced much. One needs to keep in mind the simplicity of the model before making any decision**

**Similarly for fourth order polynomial**

In [30]:

```
df['c'] = df[0]**4
cols = list(df.columns.values) #Make a list of all of the columns in the df
cols.pop(cols.index(1)) #Remove 1 from list
df = df[cols+[1]] #Create new dataframe with columns in the order you want
X = np.array(df)
X = np.insert(X, 0, 1, axis=1)
model = sm.OLS(X[:,5],X[:,5]).fit()
model.summary()
```

Out[30]:

#### OLS Regression Results

|                          |                  |                            |          |
|--------------------------|------------------|----------------------------|----------|
| <b>Dep. Variable:</b>    | y                | <b>R-squared:</b>          | 0.921    |
| <b>Model:</b>            | OLS              | <b>Adj. R-squared:</b>     | 0.915    |
| <b>Method:</b>           | Least Squares    | <b>F-statistic:</b>        | 165.1    |
| <b>Date:</b>             | Thu, 09 Nov 2017 | <b>Prob (F-statistic):</b> | 1.23e-30 |
| <b>Time:</b>             | 11:38:22         | <b>Log-Likelihood:</b>     | -77.669  |
| <b>No. Observations:</b> | 62               | <b>AIC:</b>                | 165.3    |
| <b>Df Residuals:</b>     | 57               | <b>BIC:</b>                | 176.0    |
| <b>Df Model:</b>         | 4                |                            |          |
| <b>Covariance Type:</b>  | nonrobust        |                            |          |

|              | coef       | std err  | t      | P> t  | [0.025    | 0.975]   |
|--------------|------------|----------|--------|-------|-----------|----------|
| <b>const</b> | -12.2319   | 8.145    | -1.502 | 0.139 | -28.542   | 4.078    |
| <b>x1</b>    | 0.7447     | 0.432    | 1.724  | 0.090 | -0.120    | 1.610    |
| <b>x2</b>    | -0.0076    | 0.008    | -0.948 | 0.347 | -0.024    | 0.008    |
| <b>x3</b>    | 3.684e-05  | 6.3e-05  | 0.585  | 0.561 | -8.94e-05 | 0.000    |
| <b>x4</b>    | -6.556e-08 | 1.76e-07 | -0.372 | 0.711 | -4.19e-07 | 2.88e-07 |

|                       |       |                          |          |
|-----------------------|-------|--------------------------|----------|
| <b>Omnibus:</b>       | 0.448 | <b>Durbin-Watson:</b>    | 1.535    |
| <b>Prob(Omnibus):</b> | 0.799 | <b>Jarque-Bera (JB):</b> | 0.435    |
| <b>Skew:</b>          | 0.190 | <b>Prob(JB):</b>         | 0.804    |
| <b>Kurtosis:</b>      | 2.842 | <b>Cond. No.</b>         | 1.19e+10 |

The coefficients obtained are  $\begin{aligned}$

$\hat{\beta}_0 = -12.2319$   $\hat{\beta}_1 = 0.7447$   $\hat{\beta}_2 = -0.0076$   $\hat{\beta}_3 = 3.684e^{-05}$   $\hat{\beta}_4 = -6.556e^{-08}$   $\end{aligned}$

Lets see what the sum of squared error is for this model



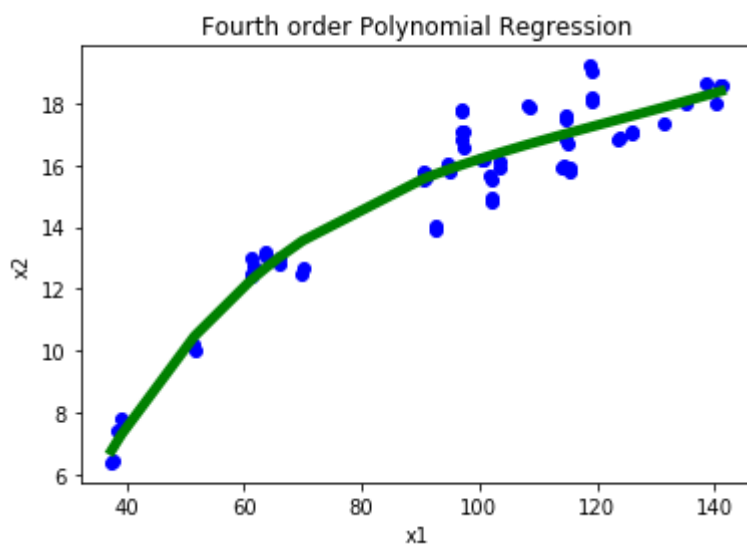
In [31]:

```
ssr.append(model.ssr)
print('\033[1m' + 'Sum of squared residuals= ' + '\033[0m', model.ssr)
```

**Sum of squared residuals= 44.4651118142**

In [32]:

```
x=model.params[0] + model.params[1] * X[:,1] + model.params[2] * X[:,2] + model.params[3] * X[:,3] + model.params[4] * X[:,4]
plt.scatter(X[:,1], X[:,5], color = 'b')
plt.plot(np.sort(X[:,1]), np.sort(x), color = 'g', linewidth=5.0)
plt.title("Fourth order Polynomial Regression")
plt.xlabel("x1")
plt.ylabel("x2")
plt.show()
```

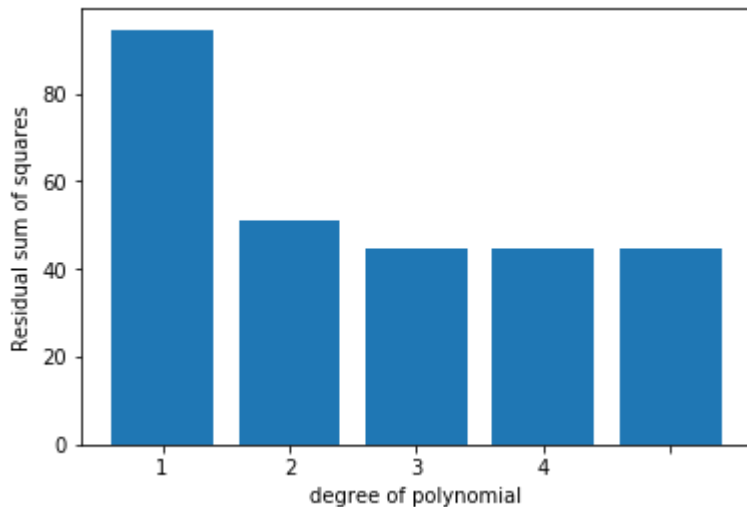


**The ssr for the fourth order polynomial has barely even changed.**

**Observing the RSS of different order models**

In [33]:

```
plt.bar(np.arange(len(ssr)),ssr)
plt.xticks(np.arange(len(ssr)),[i for i in range(1,5)])
plt.xlabel("degree of polynomial")
plt.ylabel("Residual sum of squares")
plt.show()
```



**We can infer from the above plots that the 2nd degree of polynomial would be the best fit as increasing complexity further does not reduce the sum of squared errors to a large extent.**

## Logistic Regression

In [53]:

```
lm = linear_model.LogisticRegression(C=1e10)
```

In [54]:

```
X = np.array(df)
Y = X[:,4]
X = X[:,0].reshape((62,1))
model = lm.fit(X,Y)
```

```
-----
-
ValueError                                Traceback (most recent call las
t)
<ipython-input-54-73ff87589d59> in <module>()
      2 Y = X[:,4]
      3 X = X[:,0].reshape((62,1))
----> 4 model = lm.fit(X,Y)

C:\Users\karan\Anaconda3\lib\site-packages\sklearn\linear_model\logistic.p
y in fit(self, X, y, sample_weight)
    1172         X, y = check_X_y(X, y, accept_sparse='csr', dtype=np.float
64,
    1173                                     order="C")
-> 1174         check_classification_targets(y)
    1175         self.classes_ = np.unique(y)
    1176         n_samples, n_features = X.shape

C:\Users\karan\Anaconda3\lib\site-packages\sklearn\utils\multiclass.py in
check_classification_targets(y)
    170     if y_type not in ['binary', 'multiclass', 'multiclass-multiout
put',
    171                     'multilabel-indicator', 'multilabel-sequences']:
--> 172         raise ValueError("Unknown label type: %r" % y_type)
    173
    174
```

ValueError: Unknown label type: 'continuous'

**Model does not fit the data as logistic regression is used when output variable lies between 0 and 1. But here the output variable has values above 1**

In [ ]: