

Current City	0
Python (out of 3)	0
R Programming (out of 3)	0
R Programming (out of 3)	0

```
In [48]: # We take a subset of the variables which are important for model building
data=df.drop(['Current City', 'Unnamed: 10', 'Other skills', 'Degree', 'Stream', 'Current Year Of Graduation', 'Performance_PG',
              'Performance_UG', 'Performance_12', 'Performance_10'],axis=1)
```

```
In [53]: # Created a subset of the dataset for Data Scientist Position
datax=data.drop(["PHP (out of 3)", "HTML (out of 3)", "CSS (out of 3)", "JavaScript (out of 3)",
                "AJAX (out of 3)", "Bootstrap (out of 3)", "MongoDB (out of 3)", "Node.js (out of 3)", "ReactJS (out of 3)"],
                axis=1)
```

	Python (out of 3)	R Programming (out of 3)	Deep Learning (out of 3)	MySQL (out of 3)
Application_ID				
ML0001	0	2	0	0

```
In [55]: # Created a subset for Web Development Position
data=data.drop(["Python (out of 3)", "R Programming (out of 3)",
               "Deep Learning (out of 3)"],axis=1)
```

Application_ID	ML0001	2	0	2	3	2	0	2	0	0	0
ML0001	2	0	2	3	2	0	2	0	0	0	0

```
In [57]: # Taking the total score of the different skills which a Data Scientist should possess
datax['total']=datax['Python (out of 3)']+datax['R Programming (out of 3)']+datax['Deep Learning (out of 3)']+datax['MySQL (out of 3)']
```

	Python (out of 3)	R Programming (out of 3)	Deep Learning (out of 3)	MySQL (out of 3)	total	grandtotal
Application_ID						
ML0001	0	2	0	0	2	12

```
In [61]: # build a model with total & garndtotal
# As arrays are lighter in weight hence it will work faster than the dataframe
X=datax.values[:,4:5]]
```

Number of Clusters	Sum of Squared Errors (SSE)
1	1450
2	1000
3	750
4	600
5	500
6	450
7	400
8	350
9	300
10	250

```
In [65]: Y_pred
# the clusters no are allocated on the dataset
# since we have selected the no of clusters = 3, hence 0,1,2
```

```
Out[65]: array([[0, 1, 2, 1, 0, 0, 0, 0, 1, 1, 1, 0, 1, 1, 1, 0, 1, 0, 0, 0, 1, 1, 0,
0, 1, 1, 0, 0, 0, 1, 1, 2, 1, 0, 0, 2, 2, 0, 1, 1, 1, 0, 1, 0, 1, 1,
0, 1, 0, 1, 0, 0, 0, 0, 1, 0, 1, 2, 0, 1, 0, 1, 2, 0, 1, 1, 1, 0,
1, 1, 0, 0, 1, 1, 1, 0, 0, 0, 2, 0, 1, 0, 1, 1, 1, 0, 1, 0, 1, 1,
0, 1, 2, 0, 0, 2, 1, 1, 1, 0, 1, 2, 1, 1, 2, 1, 0, 1, 0, 1, 0,
0, 1, 1, 0, 1, 0, 1, 1, 0, 1, 1, 1, 0, 0, 2, 0, 0, 1, 0, 2, 0, 0,
1, 0, 1, 1, 1, 1, 1, 2, 0, 0, 1, 1, 2, 1, 0, 1, 1, 0, 0, 1, 1, 0,
2, 0, 2, 0, 2, 1, 1, 0, 0, 0, 0, 1, 1, 0, 0, 0, 1, 2, 0, 0, 0, 2,
1, 2, 2, 2, 1, 1, 2, 0, 1, 0, 1, 1, 0, 1, 1, 1, 0, 1, 0, 0, 2,
0, 0, 0, 1, 0, 2, 1, 1, 1, 0, 0, 0, 0, 0, 0, 2, 2, 0, 1, 2, 1,
1, 0, 2, 1, 0, 1, 0, 0, 0, 0, 1, 1, 1, 0, 1, 0, 0, 0, 0, 1, 0,
0, 1, 1, 0, 0, 0, 1, 0, 0, 1, 2, 0, 0, 0, 1, 0, 0, 0, 1, 0,
0, 0, 1, 1, 0, 0, 0, 0, 0, 0, 1, 0, 0, 0, 1, 1, 0, 0, 0, 1, 0,
1, 0, 0, 0, 0, 1, 1, 0, 1, 0, 2, 0, 0, 0, 1, 0, 1, 0, 0, 0, 0, 1,
0, 1, 1, 0, 0, 1, 0, 1, 0, 1, 0, 1, 0, 1, 1, 0, 1, 0, 0, 1, 0,
1, 0, 0, 0, 1, 0, 0, 1, 0, 0, 1, 0, 0, 0, 2, 0, 0, 0, 0, 2, 0, 0,
1, 0, 1, 0, 1, 1, 1, 1, 0, 0, 0, 1, 2, 0, 0, 0, 1, 0, 0, 0, 0, 0,
2, 0, 0, 1, 0, 1, 0, 0, 0, 0, 0, 0, 1, 1, 0, 0, 1, 0, 1, 1]])
```

```
In [66]: datax['clusters']=Y_pred
# we are appending the clusters column to the dataset
```

```
In [67]: datax
```

Out[67]:

	Python (out of 3)	R Programming (out of 3)	Deep Learning (out of 3)	MySQL (out of 3)	total	grandtotal	clusters
Application_ID							
ML0001	0	2	0	0	2	12	0
ML0002	2	0	0	2	4	12	1
ML0003	3	0	1	2	6	12	2
ML0004	2	0	2	0	4	12	1
ML0005	2	0	0	0	2	12	0
...
ML0388	2	1	0	0	3	12	1
ML0389	2	0	0	0	2	12	0
ML0390	1	0	0	0	1	12	0
ML0391	2	2	0	0	4	12	1
ML0392	2	3	0	0	5	12	1

392 rows x 7 columns

```
In [68]: # Data Visualization
import seaborn as sns
sns.lmplot( data=datax, x='total', y='grandtotal',
            fit_reg=False, # No regression line
            hue='clusters',palette='Set1')
```



```
In [69]: # Changing the numerical variable to categorical Variable
# in this case the cluster
datax['clusters']=datax.clusters.map({0:'Not Applicable',1:'Maybe a Data Scientist',2:'Data Scientist'})
```

```
In [70]: # Verifying the dataset
datax
```

Out[70]:

	Python (out of 3)	R Programming (out of 3)	Deep Learning (out of 3)	MySQL (out of 3)	total	grandtotal	clusters
Application_ID							
ML0001	0	2	0	0	2	12	Not Applicable
ML0002	2	0	0	2	4	12	Maybe a Data Scientist
ML0003	3	0	1	2	6	12	Data Scientist
ML0004	2	0	2	0	4	12	Maybe a Data Scientist
ML0005	2	0	0	0	2	12	Not Applicable
...
ML0388	2	1	0	0	3	12	Maybe a Data Scientist
ML0389	2	0	0	0	2	12	Not Applicable
ML0390	1	0	0	0	1	12	Not Applicable
ML0391	2	2	0	0	4	12	Maybe a Data Scientist
ML0392	2	3	0	0	5	12	Maybe a Data Scientist

392 rows x 7 columns

```
In [25]: # Subset data for Web Development Position
datay.head()
```

Out[25]:

	PHP (out of 3)	MySQL (out of 3)	HTML (out of 3)	CSS (out of 3)	JavaScript (out of 3)	AJAX (out of 3)	Bootstrap (out of 3)	MongoDB (out of 3)	Node.js (out of 3)	ReactJS (out of 3)
Application_ID										
ML0001	2	0	2	3	2	0	2	0	0	0
ML0002	2	2	2	2	2	0	0	0	0	0
ML0003	2	2	2	0	2	0	0	0	0	0
ML0004	1	0	2	0	0	0	0	0	0	0
ML0005	2	0	2	1	1	0	0	2	2	2

```
In [26]: datay.columns
```

```
Out[26]: Index(['PHP (out of 3)', 'MySQL (out of 3)', 'HTML (out of 3)',
              'CSS (out of 3)', 'JavaScript (out of 3)', 'AJAX (out of 3)',
              'Bootstrap (out of 3)', 'MongoDB (out of 3)', 'Node.js (out of 3)',
              'ReactJS (out of 3)'],
              dtype='object')
```

```
In [27]: datay['total']=datay['PHP (out of 3)']+datay['MySQL (out of 3)']+datay['HTML (out of 3)']+datay['CSS (out of 3)']+datay['JavaScript (out of 3)']+datay['AJAX (out of 3)']+datay['Bootstrap (out of 3)']+datay['MongoDB (out of 3)']+datay['Node.js (out of 3)']+datay['ReactJS (out of 3)']
```

```
In [28]: datay['grandtotal']= 30
```

```
In [29]: datay
```

Out[29]:

	PHP (out of 3)	MySQL (out of 3)	HTML (out of 3)	CSS (out of 3)	JavaScript (out of 3)	AJAX (out of 3)	Bootstrap (out of 3)	MongoDB (out of 3)	Node.js (out of 3)	ReactJS (out of 3)	total	grandtotal
Application_ID												
ML0001	2	0	2	3	2	0	2	0	0	0	11	30
ML0002	2	2	2	2	2	0	0	0	0	0	10	30
ML0003	2	2	2	0	2	0	0	0	0	0	8	30
ML0004	1	0	2	0	0	0	0	0	0	0	3	30
ML0005	2	0	2	1	1	0	0	2	2	2	12	30
...
ML0388	0	0	2	0	0	0	0	0	0	0	2	30
ML0389	2	0	2	2	1	0	0	0	0	0	7	30
ML0390	0	0	2	2	1	0	0	0	0	0	5	30
ML0391	0	0	0	0	0	0	0	0	0	0	0	30
ML0392	2	0	2	2	3	0	0	0	0	0	9	30

392 rows x 12 columns

```
In [30]: X=datay.values[:,[10,11]]
```

```
In [31]: from sklearn.cluster import KMeans
wsse = []
for i in range(1, 10):
    kmeans = KMeans(n_clusters = i, random_state = 10)
    kmeans.fit(X)
    wsse.append(kmeans.inertia_)
plt.plot(range(1, 10), wsse)
plt.title('The Elbow Method')
plt.xlabel('Number of clusters')
plt.ylabel('WSSE')
plt.show()
```



```
In [32]: kmeans=KMeans(n_clusters=3,random_state=10)
Y_pred=kmeans.fit_predict(X)
```

```
In [33]: Y_pred
```

```
Out[33]: array([[1, 1, 0, 2, 1, 0, 2, 2, 0, 0, 2, 2, 0, 2, 2, 1, 0, 0, 2, 1, 1, 2,
0, 2, 0, 2, 0, 2, 0, 0, 0, 2, 0, 1, 2, 0, 0, 1, 2, 2, 1, 0, 2,
2, 2, 1, 2, 0, 0, 2, 1, 0, 2, 0, 2, 2, 0, 2, 1, 2, 2, 2, 2, 2, 1,
1, 2, 2, 2, 0, 1, 2, 0, 2, 0, 2, 0, 2, 0, 2, 0, 2, 1, 2, 2, 2, 2,
1, 1, 2, 2, 0, 1, 2, 0, 2, 0, 1, 0, 0, 0, 2, 0, 2, 2, 1, 0, 2,
0, 0, 1, 2, 2, 2, 0, 2, 0, 2, 1, 0, 0, 0, 1, 2, 2, 2, 2, 1, 2, 0,
2, 2, 1, 0, 0, 2, 1, 1, 0, 2, 0, 2, 0, 2, 0, 0, 0, 0, 0, 0, 2, 0,
1, 2, 0, 0, 1, 2, 2, 1, 0, 2, 2, 2, 1, 2, 0, 0, 2, 1, 0, 0, 0, 0,
2, 0, 1, 2, 0, 0, 1, 2, 2, 1, 0, 2, 2, 2, 1, 2, 0, 0, 2, 1, 0, 0,
0, 0, 2, 0, 1, 2, 0, 0, 1, 2, 2, 1, 0, 2, 2, 1, 2, 0, 0, 2, 1, 0,
1, 1, 0, 2, 1, 0, 2, 2, 0, 0, 2, 2, 0, 2, 1, 0, 0, 0, 2, 1, 1, 2,
0, 2, 0, 2, 0, 2, 0, 0, 0, 0, 2, 0, 1, 2, 0, 0, 1, 2, 2, 1, 0, 2,
2, 2, 1, 2, 0, 0, 2, 1, 0, 2, 0, 2, 2, 0, 2, 1, 2, 2, 2, 2, 2, 1,
1, 2, 2, 2, 0, 1, 2, 0, 2, 0, 2, 0, 2, 0, 2, 1, 2, 2, 2, 2, 2, 2,
1, 1, 2, 2, 0, 0, 1, 2, 0, 2, 0, 1, 0, 0, 0, 2, 0, 2, 2, 1, 0, 2,
0, 0, 1, 2, 2, 2, 0, 2, 0, 2, 1, 0, 0, 0, 1, 2, 2, 2, 2, 1, 2, 0,
2, 2, 1, 0, 0, 2, 1, 1, 0, 2, 0, 2, 0, 2, 0, 2, 0, 0, 0, 0, 2, 0,
1, 2, 0, 0, 1, 2, 2, 1, 0, 2, 2, 2, 1, 2, 0, 0, 2, 1]])
```

```
In [34]: datax['clusters']=Y_pred
```

```
In [35]: datay
```

Out[35]:

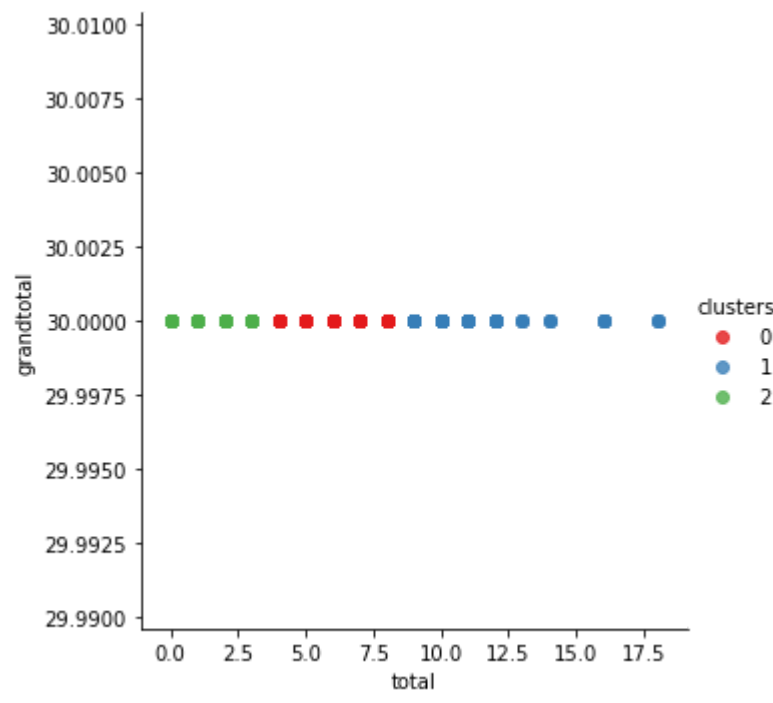
	PHP (out of 3)	MySQL (out of 3)	HTML (out of 3)	CSS (out of 3)	JavaScript (out of 3)	AJAX (out of 3)	Bootstrap (out of 3)	MongoDB (out of 3)	Node.js (out of 3)	ReactJS (out of 3)	total	grandtotal	clusters
Application_ID													
ML0001	2	0	2	3	2	0	2	0	0	0	11	30	1
ML0002	2	2	2	2	2	0	0	0	0	0	10	30	1
ML0003	2	2	2	0	2	0	0	0	0	0	8	30	0
ML0004	1	0	2	0	0	0	0	0	0	0	3	30	2
ML0005	2	0	2	1	1	0	0	2	2	2	12	30	1
...
ML0388	0	0	2	0	0	0	0	0	0	0	2	30	2
ML0389	2	0	2	2	1	0	0	0	0	0	7	30	0
ML0390	0	0	2	2	1	0	0	0	0	0	5	30	0
ML0391	0	0	0	0	0	0	0	0	0	0	0	30	2
ML0392	2	0	2	2	3	0	0	0	0	0	9	30	1

392 rows x 13 columns

```
In [36]: import seaborn as sns

sns.lmplot(data=datay, x='total', y='grandtotal',
          fit_reg=False, # No regression line
          hue='clusters',palette="Set1")
```

Out[36]: <seaborn.axisgrid.FacetGrid at 0x11806afae08>



```
In [37]: datay['clusters']=datay.clusters.map({0:'Not Applicable',1:'May be in Web Development',2:'Web Development'})
```

```
In [38]: datay.head()
```

	PHP (out of 3)	MySQL (out of 3)	HTML (out of 3)	CSS (out of 3)	JavaScript (out of 3)	AJAX (out of 3)	Bootstrap (out of 3)	MongoDB (out of 3)	Node.js (out of 3)	ReactJS (out of 3)	total	grandtotal	clusters
Application_ID													
ML0001	2	0	2	3	2	0	2	0	0	0	11	30	May be in Web Development
ML0002	2	2	2	2	2	0	0	0	0	0	10	30	May be in Web Development
ML0003	2	2	2	0	2	0	0	0	0	0	8	30	Not Applicable
ML0004	1	0	2	0	0	0	0	0	0	0	3	30	Web Development
ML0005	2	0	2	1	1	0	0	2	2	2	12	30	May be in Web Development

```
In [41]: # Merging the Data Scientist & Web Development DataFrame
jd = pd.merge(datay, datay, on = 'Application_ID', how = 'inner', indicator = False)
jd.head()
```

	Python (out of 3)	R Programming (out of 3)	Deep Learning (out of 3)	MySQL (out of 3)_x	total_x	grandtotal_x	clusters_x	PHP (out of 3)	MySQL (out of 3)_y	HTML (out of 3)	CSS (out of 3)	JavaScript (out of 3)	AJAX (out of 3)	Bootstrap (out of 3)	MongoDB (out of 3)	Node.js (out of 3)	ReactJS (out of 3)	total_y	grandtotal_y	clusters_y
Application_ID																				
ML0001	0	2	0	0	2	12	Not Applicable	2	0	2	3	2	0	2	0	0	0	11	30	May be in Web Development
ML0002	2	0	0	2	4	12	Maybe a Data Scientist	2	2	2	2	2	0	0	0	0	0	10	30	May be in Web Development
ML0003	3	0	1	2	6	12	Data Scientist	2	2	2	0	2	0	0	0	0	0	8	30	Not Applicable
ML0004	2	0	2	0	4	12	Maybe a Data Scientist	1	0	2	0	0	0	0	0	0	0	3	30	Web Development
ML0005	2	0	0	0	2	12	Not Applicable	2	0	2	1	1	0	0	0	2	2	12	30	May be in Web Development

```
In [71]: # Saving df to a file
df.to_csv(r"F:\DSP\Class\Python\Basics\Wavoi\JobRole.csv",
          index = True, header = True)
```