# Web Mining (CSE3024)
# Lab Assignment 8

<u>Name:</u> **Kritika Mishra**

<u>Registration Number:</u> **16BCI0041**

<u>Slot:</u> L15+L16

<u>Date:</u> 3rd October 2018

<u>Question:</u>

**Implement a k-means algorithm with sklearn to partition observations in a dataset into a specific number of clusters in order to aid in analysis of the data.**

- **Use Sklearn Toolkit / Package to perform the process**
- **Import kmeans and PCA through the sklearn library**
- **Devise an elbow curve to select the optimal number of clusters (k)**
- **Generate and visualise a k-means clustering algorithms**

**Note : Dataset in CSV can be generated or downloaded from the internet. Please specify the source of the dataset in the documentation steps of this program.**

<mark>Dataset:</mark>

http://www.michaeljgrogan.com/datasets/  ->sample_stocks.csv

<mark>Code:</mark>

```python
import pandas
import pylab as pl
from sklearn.cluster import KMeans
from sklearn.decomposition import PCA
variables = pandas.read_csv('sample_stocks.csv')
Y = variables[['returns']]
X = variables[['dividendyield']]
X_norm = (X - X.mean()) / (X.max() - X.min())
Y_norm = (Y - Y.mean()) / (Y.max() - Y.min())
Nc = range(1, 20)
kmeans = [KMeans(n_clusters=i) for i in Nc]
kmeans
score = [kmeans[i].fit(Y).score(Y) for i in range(len(kmeans))]
score
```

```python
pl.plot(Nc,score)
pl.xlabel('Number of Clusters')
pl.ylabel('Score')
pl.title('Elbow Curve')
pl.show()
pca = PCA(n_components=1).fit(Y)
pca_d = pca.transform(Y)
pca_c = pca.transform(X)
kmeans=KMeans(n_clusters=3)
kmeansoutput=kmeans.fit(Y)
kmeansoutput
pl.figure('3 Cluster K-Means')
pl.scatter(pca_c[:, 0], pca_d[:, 0], c=kmeansoutput.labels_)
pl.xlabel('Dividend Yield')
pl.ylabel('Returns')
pl.title('3 Cluster K-Means')
pl.show()
```

## Output:

3 Cluster K-Means



Elbow Curve