

Nikhil Panse - IIT Guwahati

Importing text file and splitting it to create data variables

In [154]:

```
1import pandas as pd
2data = [line.strip() for line in open("train_set.txt", "r").readlines()]
3data = [{"LineNumber": line.split("~")[0], "Number": line.split("~")[1].split("IN")
4data = pd.DataFrame(data)
5dfctest = [line.strip() for line in open("public_test_set.txt", "r").readlines()]
6dfctest = [{"LineNumber": line.split("~")[0], "Number": line.split("~")[1].split("IN")
7dfctest = pd.DataFrame(dfctest)
8dp = [line.strip() for line in open("private_test_set.txt", "r").readlines()]
9dp = [{"LineNumber": line.split("~")[0], "Number": line.split("~")[1].split("INFC")
10dp = pd.DataFrame(dp)
```

In [155]:

1	dfctest					
5	L5	"GET /v2/2019-09-10 meta_data.json HTTP/1.1" ...	62161	req-y5tj-ufwp-34nme6-6do12x	gs.api.openapi.compute.wsgi	
6	L6	HTTP exception thrown: No entry found for any...	13691	req-x52i0-95ewzi-tnxuz	gs.openstack.wsgi.server	
7	L7	"GET /v2/2019-09-12 meta_data.json HTTP/1.1" ...	25112	req-joh69-ydjr7v-rhezju-a7nq-rytxe8	gs.openapi.wsgi.server	
8	L8	Instance spawned correctly	94733	req-vz2nx-5eva8-78jxq	gs.api.openapi.compute.wsgi	

Importing data as text file and separating by new lines.

Using Term Frequency- Inverse Document Frequency to understand word frequencies

In [180]:

```
1 with open('train_set.txt') as f:
2     file = [line.rstrip('\n') for line in f]
3
4 from sklearn.feature_extraction.text import TfidfVectorizer
5 vectorizer = TfidfVectorizer(analyzer="word", max_features=125)
6 X = vectorizer.fit_transform(file)
7 from scipy.sparse import csr_matrix
8 X = csr_matrix.todense(X)
9
```

K Means clustering

In [190]:

```
1 from sklearn.cluster import KMeans
2 kmeans = KMeans(n_clusters=18, random_state=0).fit(X)
```

Public dataset

In [191]:

```
1 with open('public_test_set.txt') as f:
2     test = [line.rstrip('\n') for line in f]
3 vectorizer = TfidfVectorizer(analyzer="word", max_features=125)
4 Xt = vectorizer.fit_transform(test)
5 Xt = csr_matrix.todense(Xt)
6
```

Clustering public dataset

In [192]:

```
1 d = {"LineNumber":dfstest["LineNumber"].values, "Class":kmeans.predict(Xt)}
2 df = pd.DataFrame(d)
3 df.to_csv("public_test.csv", index=False)
4
```

Private Dataset

In [193]:

```
1 with open('private_test_set.txt') as f:
2     testp = [line.rstrip('\n') for line in f]
3 vectorizer = TfidfVectorizer(analyzer="word", max_features=125)
4 Xp = vectorizer.fit_transform(testp)
5 Xp = csr_matrix.todense(Xp)
6 p = {"LineNumber":dp["LineNumber"].values, "Class":kmeans.predict(Xp)}
7
8 dp = pd.DataFrame(p)
9 dp.to_csv("private_test.csv", index=False)
```

In []:

```
1
```

