

Name : Sumon Singh

Roll No. 16

Paper : Statistical Method

Import Libraries

```
In [1]: import pandas as pd
from sklearn.cluster import KMeans
from sklearn import metrics
from sklearn.model_selection import train_test_split
from sklearn.linear_model import LinearRegression
import numpy as np
from sklearn.decomposition import PCA
from pandas_profiling import ProfileReport
```

Modeling On Unigram Dataset

Read the unigram dataset

```
In [2]: df_uni = pd.read_excel('~/home/sumon/Documents/Datasets/DS-sem1.xlsx',sheet_name='uni_euclid_mat_intersection')

In [3]: df_uni

Out[3]:
```

	Sr #	raga	closest	farthest	Unnamed: 4	euclidean distance	AbhaireeTodi	Abhogi	Adana	Adi Bharivi	...
0	1	AbhaireeTodi	YogYjoti	Plasi	NaN	AbhaireeTodi	0.000000	152.272782	257.699437	171.452617	...
1	2	Abhogi	DevMukhani	Bhageshari	NaN	Abhogi	152.272782	0.000000	250.868491	67.461030	...
2	3	Adana	KalyanBasanti	Bangal Bhaivar	NaN	Adana	257.699437	250.868491	0.000000	229.780330	...
3	4	Adi Bharivi	Rudher,Gandhar	Jaywanti	NaN	Adi Bharivi	171.452617	67.461030	229.780330	0.000000	...
4	5	Aheer Bhariv	Hamir	Bhageshari	NaN	Aheer Bhariv	230.822443	204.151904	104.584894	249.242853	...
...
324	325	Yamani Bilawal	Vegveghini	Bhageshari	NaN	Yamani Bilawal	526.480769	214.932342	606.164994	238.470124	...
326	326	YamanKalyan	Pahadi	Bangal Bhaivar	NaN	YamanKalyan	538.697905	268.233108	658.758661	406.006900	...
325	327	YashPriya	Rsik_priya	SerswatSarang	NaN	YashPriya	424.915075	135.173962	393.010178	278.688867	...
327	328	YashRanjani	Anandkila	Jaywanti	NaN	YashRanjani	426.909827	322.225077	221.621299	599.581003	...
328	329	YogYjoti	ChanderRekha	NagSwarabali	NaN	YogYjoti	227.525823	155.730537	385.112970	285.115766	...

329 rows × 335 columns

Data cleaning of unigram dataset

```
In [4]: df_uni.columns

Out[4]: Index(['Sr #', 'raga', 'closest', 'farthest', 'Unnamed: 4', 'euclidean distance', 'AbhaireeTodi', 'Abhogi', 'Adana', 'Adi Bharivi',
      ...,
      'Vidyapati', 'Vijay', 'VijayKokila', 'Vinayak', 'Yaman', 'Yamani Bilawal', 'YamanKalyan', 'YashPriya', 'YashRanjani',
      'YogYjoti'],
      dtype='object', length=335)

In [5]: df_uni.shape

Out[5]: (329, 335)

Drop the unnecessary columns

In [6]: df_uni.drop(df_uni.columns[0:5],axis=1,inplace=True)

In [7]: df_uni.dtypes

Out[7]: euclidean distance      object
AbhaireeTodi              float64
Abhogi                    float64
Adana                     float64
Adi Bharivi                float64
Yamani Bilawal            float64
YamanKalyan              float64
YashPriya                 float64
YashRanjani               float64
YogYjoti                  float64
Length: 330, dtype: object
```

Clustering on unigram dataset

Select only numeric data for clustering

```
In [8]: X=df_uni.iloc[:,1:]

K-means cluster model

In [9]: model_uni=KMeans(n_clusters=10,random_state=0)

In [10]: model_uni.fit(X)

Out[10]: KMeans(n_clusters=10, random_state=0)

In [11]: p_uni=model_uni.predict(X)

In [12]: df_uni["target"]=p_uni

Quality of cluster

In [13]: metrics.silhouette_score(X,model_uni.labels_)

Out[13]: 0.118143413849179

In [14]: result_pred_uni=DataFrame({'X':df_uni.iloc[:,0], 'Y':p_uni})

In [15]: result_pred_uni.sort_values('Y')

Out[15]:
```

	X	Y
72	Devgiri	0
21	Barba	0
225	Noopur	0
142	KalyanKaisri	0
145	KamalManoharee	0
...
5	Aheerfalti	9
175	LachanTodi	9
50	Chanderika Bhairvi	9
283	Sawanikalyan	9
113	Hembanli	9

329 rows × 2 columns

Regression on unigram dataset

Extract top 20 best correlated features with the cluster

```
In [17]: df_uni.corr().iloc[:,~1].abs().sort_values(ascending=False).iloc[0:20].index

Out[17]: target      1.000000
Girija      0.455178
KamaRanjani 0.361103
SerswatRanjani 0.359402
NagSwarabali 0.340339
Varati      0.339283
ChandarNandan 0.321111
DevkipurTiya 0.315450
Dhaney Dhaivat 0.306547
JaunaPuri    0.297888
Neetmat      0.285860
Chhanika     0.275447
Adana        0.274022
Manvi        0.271846
KalaBati     0.271613
MangalGujari 0.271463
Abhogi       0.267413
Sohni        0.262053
KaliyanBasant 0.260845
SgurnRanjani 0.259410
Name: target, dtype: float64

In [18]: fea_col_uni=df_uni.corr().iloc[:,~1].abs().sort_values(ascending=False).iloc[1:21].index

In [19]: fea_col_uni

Out[19]: Index(['Girija', 'KamaRanjani', 'SerswatRanjani', 'NagSwarabali', 'varati',
      'Neetmat', 'Chhanika', 'Adana', 'Manvi', 'KalaBati', 'MangalGujari',
      'Abhogi', 'Sohni', 'KaliyanBasant', 'SgurnRanjani', 'Rsik_priya'],
      dtype='object')
```

Create new dataset of top 20 correlated features with the target

```
In [20]: df_uni_new=df_uni[fea_col_uni]

In [21]: df_uni_new

Out[21]:
```

	Girija	KamaRanjani	SerswatRanjani	NagSwarabali	varati	ChandarNandan	DevkipurTiya	Dhaney Dhaivat	JaunaPuri	Neetmat
0	545.947000	386.139871	162.769776	486.569625	424.123803	352.681159	250.700117	273.867627	186.316934	340.386630
1	345.421482	230.282435	93.749667	381.241393	285.220667	305.638676	100.970293	151.383619	197.909070	249.094360
2	533.702164	425.779286	341.764539	647.787734	534.815856	425.806294	349.455290	329.913625	161.830158	470.059577
3	410.369346	138.032005	244.965304	406.135445	422.407386	286.052443	192.296646	194.480503	158.154987	281.524422
4	297.307189	213.475994	262.994297	610.286818	417.435025	426.075111	307.866086	155.521703	184.591635	331.561435
...
324	516.111422	249.819355	372.478645	438.421031	353.980225	690.703577	429.633565	299.604740	531.968044	474.968310
325	573.869323	284.416947	300.932219	607.080720	395.845930	414.627529	356.164288	278.154653	619.099346	432.468877
326	456.404426	310.665415	163.719272	287.416075	336.740256	348.162318	211.928874	235.538906	343.250637	279.089590
327	334.232451	315.412428	500.236944	337.615166	435.873835	333.151617	509.374126	291.982876	297.889241	446.882210
328	599.976666	323.677038	131.045794	603.688661	394.124346	480.960497	270.992620	209.475536	351.239235	235.792280

329 rows × 20 columns

Spitting dataset in training and testing

```
In [22]: X_train,X_test,Y_train,Y_test=train_test_split(df_uni_new,p_uni,test_size=0.3,random_state=0)

Dimension reduction using PCA

In [23]: pca=PCA(n_components=2)

In [24]: X_train=pca.fit_transform(X_train)

In [25]: X_test=pca.transform(X_test)

Regression Model

In [26]: lm_model_uni = LinearRegression().fit(X_train,Y_train)

In [27]: lm_model_pred_uni=lm_model_uni.predict(X_test)

In [28]: lm_model_pred_uni=np.round(lm_model_pred_uni)

Accuracy of regression on model

In [29]: metrics.mean_absolute_error(Y_test,lm_model_pred_uni)

Out[29]: 1.6666666666666667

In [30]: metrics.mean_squared_error(Y_test,lm_model_pred_uni)

Out[30]: 4.717171717171717

In [31]: metrics.r2_score(Y_test,lm_model_pred_uni)

Out[31]: 0.2665154287386724
```

Modeling On bigram Dataset

Read the bigram dataset

```
In [32]: df_bi = pd.read_excel('~/home/sumon/Documents/Datasets/DS-sem1.xlsx',sheet_name='bi_euclid_mat_intersection')

In [33]: df_bi

Out[33]:
```

	raga	closest	farthest	Unnamed: 3	euclidean distance	AbhaireeTodi	Abhogi	Adana	Adi Bharivi	Aheer Bhariv
0	AbhaireeTodi	YogYjoti	Plasi	NaN	AbhaireeTodi	0.000000	123.511133	122.073748	41.827748	75.903688
1	Adana	KalyanBasanti	Bangal Bhaivar	NaN	Adana	123.511133	99.030298	0.000000	138.996361	2.829427
2	Adi Bharivi	Rudher,Gandhar	Jaywanti	NaN	Adi Bharivi	123.511133	43.817805	138.949631	0.000000	24.041831
3	Aheer Bhariv	Hamir	Bhageshari	NaN	Aheer Bhariv	70.590368	41.82521	2.829427	24.041831	0.000000
...
324	Yamani Bilawal	Vegveghini	Bhageshari	NaN	Yamani Bilawal	242.950612	83.264638	142.720006	80.199751	148.243044
325	YamanKalyan	Pahadi	Bangal Bhaivar	NaN	YamanKalyan	210.853029	103.013274	132.174128	133.592664	86.029065
326	YashPriya	Rsik_priya	SerswatSarang	NaN	YashPriya	123.656783	78.752778	154.369686	96.010416	56.187187
327	YashRanjani	Anandkila	Jaywanti	NaN	YashRanjani	159.132021	13.000000	37.588892	157.308862	177.910090
328	YogYjoti	ChanderRekha	NagSwarabali	NaN	YogYjoti	117.025638	53.507009	130.610107	68.942505	51.458721

329 rows × 334 columns

Data cleaning of bigram dataset

```
In [34]: df_bi.columns

Out[34]: Index(['raga', 'closest', 'farthest', 'Unnamed: 3', 'euclidean distance',
      'AbhaireeTodi', 'Abhogi', 'Adana', 'Adi Bharivi', 'Aheer Bhariv',
      ...,
      'Vidyapati', 'Vijay', 'VijayKokila', 'Vinayak', 'Yaman', 'Yamani Bilawal', 'YamanKalyan', 'YashPriya', 'YashRanjani',
      'YogYjoti'],
      dtype='object', length=334)

Drop the unnecessary columns

In [35]: df_bi.drop(df_bi.columns[:4],axis=1,inplace=True)

In [36]: df_bi.dtypes

Out[36]: euclidean distance      object
AbhaireeTodi              float64
Abhogi                    float64
Adana                     float64
Adi Bharivi                float64
Yamani Bilawal            float64
YamanKalyan              float64
YashPriya                 float64
YashRanjani               float64
YogYjoti                  float64
Length: 330, dtype: object
```

Clustering on bigram dataset

Select only numeric data for clustering

```
In [37]: X_bi=df_bi.iloc[:,1:]

K-means cluster model

In [38]: model_bi=KMeans(n_clusters=10,random_state=0)

In [39]: model_bi.fit(X_bi)

Out[39]: KMeans(n_clusters=10, random_state=0)

In [40]: p_bi=model_bi.predict(X_bi)

In [41]: df_bi["target"]=p_bi

Quality of cluster

In [42]: metrics.silhouette_score(X_bi,model_bi.labels_)

Out[42]: 0.0998404978314635

In [43]: result_pred_bi=pd.DataFrame({'X':df_bi.iloc[:,0], 'Y':p_bi})

In [44]: result_pred_bi.sort_values('Y')

Out[44]:
```

	X	Y
217	Nand	0
25	Bawani	0
52	ChanderRekha	0
127	Jayant Malhar	0
125	Jay-Yaji Banti	0
...
193	Maligaura	9
178	LalitMalini	9
194	Malini	9
314	Tivneri	9
269	Roomwals	9

329 rows × 2 columns

```
In [45]: model_bi.cluster_centers_

Out[45]: array([[114.40996745,  80.58139664,  98.97289896, ..., 109.48968632,
      41.38217233,  91.69671444],
      [121.68925892,  61.60983111, 122.73233388, ..., 111.61396375,
      128.8085676 , 115.41845699],
      [ 89.04394764,  73.94789384,  83.9214368 , ..., 74.81666919,
      95.34894875,  73.0497658 ],
      [121.68925892, 124.19586952, 20.35961384, ..., 74.41977266,
      171.23181279,  82.95947827]])

Regression on bigram dataset

Extract top 20 best correlated features with the cluster

In [46]: df_bi.corr().iloc[:,~1].abs().sort_values(ascending=False).iloc[0:20].index

Out[46]: target      1.000000
Gope Kamodiji 0.488877
Pahadi Kamodi 0.477454
RaketHans     0.468059
Samairee      0.464335
Karankwralii  0.455289
Hans Manjari  0.444946
Kukubh        0.432078
Kamaiyee      0.423506
GujariTodi    0.412382
Kukubh        0.412265
KanakBasant   0.407818
RasRanjani    0.407316
SamaireeMalhar 0.407052
ManRang       0.399556
Chandrlika   0.388684
Jayant        0.389092
Khoker        0.378607
ChhyavaTodi  0.374847
Sur malhar    0.370452
Name: target, dtype: float64

In [47]: fea_col_bi=df_bi.corr().iloc[:,~1].abs().sort_values(ascending=False).iloc[1:21].index

In [48]: fea_col_bi

Out[48]: Index(['Gope Kamodiji', 'Pahadi Kamodi', 'RaketHans', 'Samairee', 'Karankwralii',
      'Hans Manjari', 'Kekhas', 'Kamaiyee', 'GujariTodi', 'kukubh',
      'KanakBasant', 'RasRanjani', 'SamaireeMalhar', 'ManRang', 'Chandrlika',
      'Jayant', 'Khoker', 'ChhyavaTodi', 'Sur malhar', 'RaeRanjani'],
      dtype='object')
```

Create new dataset of top 20 correlated features with the target

```
In [49]: df_bi_new=df_bi[fea_col_bi]

In [50]: df_bi_new

Out[50]:
```

	Gope Kamodiji	Pahadi Kamodi	RaketHans	Samairee	Karankwralii	Hans Manjari	Kolhas	Kamaiyee	GujariTodi	kukubh	KanakBasar
0	146.400137	200.544259	167.017963	147.193750	133.555232	172.365890	164.170843	168.433963	61.530400	183.806964	138.95393
1	105.806427	113.225439	132.238421	138.747973	100.184829	132.657454	104.661359	91.809586	0.000000	103.889899	140.35995
2	175.982554	163.367683	167.017841	178.275068	155.048268	187.054489	198.924899	169.826382	49.819675	154.673801	131.14649
3	96.954535	126.035151	119.653667	151.043040	107.536040	143.289916	118.941612	80.890403	90.872215	131.43722	112.43722
4	62.072538	100.054985	69.318107	43.174066	60.415230	83.078376	67.305275	135.738720	88.475985	123.328829	75.21967
...
324	161.095738	163.767059	194.797331	151.115950	236.282192	212.903734	130.418567	200.117466	25.455844	203.803827	157.80367
325	192.919637	157.200509	194.797331	163.205392	182.477390	240.299813	195.662176	206.891276	31.304952	225.959236	145.25494
326	108.415866	109.538947	129.410201	69.783888	64.845971	150.519102	171.749007	56.787887	54.543734	160.993789	108.21373
327	44.508428	73.579889	67.941151	74.390860	63.521650	121.954910	92.892411	64.288413	233.942301	90.360380	68.11020
328	93.420118	104.591607	135.547040	85.545310	104.527508	143.553385	113.313724	94.942088	60.659964	162.728837	98.62555

329 rows × 20 columns

Spitting dataset in training and testing

```
In [51]: X_train,X_test,Y_train,Y_test=train_test_split(df_bi_new,p_bi,test_size=0.3,random_state=0)

Dimension reduction using PCA

In [52]: pca=PCA(n_components=2)

In [53]: X_train=pca.fit_transform(X_train)

In [54]: X_test=pca.transform(X_test)

Regression Model

In [55]: lm_model_bi = LinearRegression().fit(X_train,Y_train)

In [56]: lm_model_pred_bi=lm_model_bi.predict(X_test)

In [57]: lm_model_pred_bi=np.round(lm_model_pred_bi)

Accuracy of regression model

In [58]: metrics.mean_absolute_error(Y_test,lm_model_pred_bi)

Out[58]: 2.3636363636363638

In [59]: metrics.mean_squared_error(Y_test,lm_model_pred_bi)

Out[59]: 7.838383838383838

In [60]: metrics.r2_score(Y_test,lm_model_pred_bi)

Out[60]: 0.1388665457634913
```

Regression on trigram Dataset

Read the trigram dataset

```
In [61]: df_tri = pd.read_excel('~/home/sumon/Documents/Datasets/DS-sem1.xlsx',sheet_name='tri_euclid_mat_intersection')

In [62]: df_tri

Out[62]:
```

	raga	closest	farthest	Unnamed: 3	euclidean distance	AbhaireeTodi	Abhogi	Adana
0	AbhaireeTodi	vallajani	NagSwarabali	NaN	AbhaireeTodi	0.000000	34.088121	43.335899
1	Abhogi	Aheeralti,Bilashkani,Dhoulambeni,Gaumna,Gumna...	Haivri	NaN	Abhogi	34.088121	0.000000	22.922078
2	Adana	Aheeralti,Bilashkani,Dhoulambeni	Bangal Bhaivar	NaN	Adana	43.335897	32.939338	0.000000
3	Adi Bharivi	Anand Bhariv,BhalYaar,Malini,Schri	Jayant	NaN	Adi Bharivi	48.363209	41.082777	64.845971
4	Aheer Bhariv	Chandrika	Bhageshari	NaN	Aheer Bhariv	17.805751	0.000000	0.000000
...
324	Yamani Bilawal	Gundkali	Chhaya	NaN	Yamani Bilawal	78.835271	17.748239	34.713131
325	YamanKalyan	Bilashkani	NagSwarabali	NaN	YamanKalyan			