```
1 from google.colab import files
2 uploaded = files.upload()
```

Choose Files | No file chosen    Upload widget is only available when the cell has been executed in the current browser session. Please rerun this cell to

Saving dummydata.csv to dummydata.csv

```
1 import pandas as pd
2 import numpy as np
3 df = pd.read_csv('dummydata.csv')
4 pd.DataFrame.from_records(df)
5 df.head()
```

| | Unnamed: 0 | age | self_employed | family_history | mh_treatment | interfere | company_size | remote | tech_company | mh_negativ |
|---|---|---|---|---|---|---|---|---|---|---|
| 0 | 1 | 3 | 1 | 1 | 0 | 3 | 1 | 1 | 1 | |
| 1 | 2 | 2 | 0 | 1 | 1 | 3 | 4 | 0 | 1 | |
| 2 | 3 | 2 | 1 | 0 | 0 | 1 | 1 | 1 | 1 | |
| 3 | 4 | 3 | 0 | 0 | 1 | 4 | 3 | 1 | 1 | |
| 4 | 5 | 3 | 0 | 0 | 1 | 1 | 6 | 0 | 0 | |

```
1 df.describe
```

```
<bound method NDFrame.describe of      Unnamed: 0  age  ...  country_Ireland  country_India
0              1    3  ...                0              0
1              2    2  ...                0              0
2              3    2  ...                0              0
3              4    3  ...                0              0
4              5    3  ...                0              0
..           ...  ...  ...              ...            ...
891          892    2  ...                0              0
892          893    2  ...                0              0
893          894    2  ...                0              0
894          895    2  ...                0              0
895          896    3  ...                0              0

[896 rows x 51 columns]>
```

```
1 df.columns
```

```
Index(['Unnamed: 0', 'age', 'self_employed', 'family_history', 'mh_treatment',
       'interfere', 'company_size', 'remote', 'tech_company',
       'mh_negative_consequence_flag', 'ph_negative_consequence_flag',
       'mh_disscuss_coworker', 'mh_disscuss_supervisor',
       'interview_mh_bringup', 'interview_ph_bringup', 'witness_mh_nc',
       'anonymity_protected_Yes', 'anonymity_protected_No',
       'anonymity_protected_Don't know', 'awareness_mh_benefits_Not sure',
       'awareness_mh_benefits_Yes', 'awareness_mh_benefits_No', 'gender_M',
       'gender_F', 'gender_T', 'medical_leave_easy_Very easy',
       'medical_leave_easy_Somewhat difficult',
       'medical_leave_easy_Don't know', 'medical_leave_easy_Very difficult',
       'medical_leave_easy_Somewhat easy', 'mh_benefits_Yes', 'mh_benefits_No',
       'mh_benefits_Don't know', 'mh_discuss_Yes', 'mh_discuss_No',
       'mh_discuss_Don't know', 'mh_resources_Don't know', 'mh_resources_No',
       'mh_resources_Yes', 'mh_serious_ph_Yes', 'mh_serious_ph_No',
       'mh_serious_ph_Don't know', 'country_United States',
       'country_United Kingdom', 'country_Canada', 'country_Netherlands',
       'country_Australia', 'country_France', 'country_Germany',
       'country_Ireland', 'country_India'],
      dtype='object')
```

```
1 df.drop(df.columns[[0]], axis=1, inplace=True)
2 df.head()
```

| | age | self_employed | family_history | mh_treatment | interfere | company_size | remote | tech_company | mh_negative_consequ |
|---|---|---|---|---|---|---|---|---|---|
| 0 | 3 | 1 | 1 | 0 | 3 | 1 | 1 | 1 | |
| 1 | 2 | 0 | 1 | 1 | 3 | 4 | 0 | 1 | |
| 2 | 2 | 1 | 0 | 0 | 1 | 1 | 1 | 1 | |
| 3 | 3 | 0 | 0 | 1 | 4 | 3 | 1 | 1 | |
| 4 | 3 | 0 | 0 | 1 | 1 | 6 | 0 | 0 | |

```
1 df.shape
```

```
(896, 50)
```

## Supervised ML

```
1 X=df.loc[:, df.columns != 'mh_treatment']
2 X.head()
```

| | age | self_employed | family_history | interfere | company_size | remote | tech_company | mh_negative_consequence_flag | ph_ |
|---|---|---|---|---|---|---|---|---|---|
| 0 | 3 | 1 | 1 | 3 | 1 | 1 | 1 | 0 | |
| 1 | 2 | 0 | 1 | 3 | 4 | 0 | 1 | 1 | |
| 2 | 2 | 1 | 0 | 1 | 1 | 1 | 1 | 0 | |
| 3 | 3 | 0 | 0 | 4 | 3 | 1 | 1 | 1 | |
| 4 | 3 | 0 | 0 | 1 | 6 | 0 | 0 | 1 | |

```
1 X.columns
```

```
Index(['age', 'self_employed', 'family_history', 'interfere', 'company_size',
       'remote', 'tech_company', 'mh_negative_consequence_flag',
       'ph_negative_consequence_flag', 'mh_disscuss_coworker',
       'mh_disscuss_supervisor', 'interview_mh_bringup',
       'interview_ph_bringup', 'witness_mh_nc', 'anonymity_protected_Yes',
       'anonymity_protected_No', 'anonymity_protected_Don't know',
       'awareness_mh_benefits_Not sure', 'awareness_mh_benefits_Yes',
       'awareness_mh_benefits_No', 'gender_M', 'gender_F', 'gender_T',
       'medical_leave_easy_Very easy', 'medical_leave_easy_Somewhat difficult',
       'medical_leave_easy_Don't know', 'medical_leave_easy_Very difficult',
       'medical_leave_easy_Somewhat easy', 'mh_benefits_Yes', 'mh_benefits_No',
       'mh_benefits_Don't know', 'mh_discuss_Yes', 'mh_discuss_No',
       'mh_discuss_Don't know', 'mh_resources_Don't know', 'mh_resources_No',
       'mh_resources_Yes', 'mh_serious_ph_Yes', 'mh_serious_ph_No',
       'mh_serious_ph_Don't know', 'country_United States',
       'country_United Kingdom', 'country_Canada', 'country_Netherlands',
       'country_Australia', 'country_France', 'country_Germany',
       'country_Ireland', 'country_India'],
      dtype='object')
```

```
1 y=df.mh_treatment
2 y.head()
```

```
0    0
1    1
2    0
3    1
4    1
Name: mh_treatment, dtype: int64
```

```
1 np.random.seed(888)
2 #Split X and y into train (70%) and test (30%) sets.
3 from sklearn.model_selection import train_test_split
4 # Let 20% of the data to be a test set
5 train_x, test_x, train_y, test_y = train_test_split(X, y, test_size=0.30, random_state=0)
```

### 1. Linear Regression

```
1 #fit the model
2 from sklearn.linear_model import LinearRegression
3 from sklearn.metrics import mean_squared_error
4 model1 = LinearRegression()
5 model1.fit(train_x,train_y)
```

```
LinearRegression(copy_X=True, fit_intercept=True, n_jobs=None, normalize=False)
```

```
1 # coefficeints of the trained model
2 print('\nCoefficient of model :', model1.coef_)
3
4 # intercept of the model
5 print('\nIntercept of model',model1.intercept_)
6
7 # predict the target on the test dataset
```

```
 8 predict_train1 = model1.predict(train_x)
 9
10 # Root Mean Squared Error on training dataset
11 rmse_train1 = mean_squared_error(train_y,predict_train1)**(0.5)
12 print('\nRMSE on train dataset : ', rmse_train1)
13
```

```
Coefficient of model : [ 7.94394178e-02 -3.13073938e-02  1.35247094e-01  2.00005298e-01
  2.16494895e-03  6.28433680e-03 -3.96517007e-03  2.29191247e-02
 -9.32070297e-03  6.02567060e-02 -8.94607436e-03  9.01932694e-03
  1.84655909e-02  6.08269687e-03  2.47073445e+11  2.47073445e+11
  2.47073445e+11 -6.17962248e+11 -6.17962248e+11 -6.17962248e+11
 -2.39548373e+11 -2.39548373e+11 -2.39548373e+11 -3.63133064e+11
 -3.63133064e+11 -3.63133064e+11 -3.63133064e+11 -3.63133064e+11
 -4.70124964e+11 -4.70124964e+11 -4.70124964e+11 -5.22730142e+11
 -5.22730142e+11 -5.22730142e+11 -7.48774855e+11 -7.48774855e+11
 -7.48774855e+11 -1.77157781e+11 -1.77157781e+11 -1.77157781e+11
 -1.33097841e+11 -1.33097841e+11 -1.33097841e+11 -1.33097841e+11
 -1.33097841e+11 -1.33097841e+11 -1.33097841e+11 -1.33097841e+11
 -1.33097841e+11]

Intercept of model 3025455823032.587

RMSE on train dataset :  0.38898088972415895
```

```
1 # predict the target on the testing dataset
2 predict_test1 = model1.predict(test_x)
3
4 # Root Mean Squared Error on testing dataset
5 rmse_test1 = mean_squared_error(test_y,predict_test1)**(0.5)
6 print('\nRMSE on test dataset : ', rmse_test1)
```

```
RMSE on test dataset :  0.3846918650527029
```

## 2. Tree

```
1 from sklearn.tree import DecisionTreeClassifier
2 from sklearn.metrics import accuracy_score
3 model2 = DecisionTreeClassifier()
4
5 # fit the model with the training data
6 model2.fit(train_x,train_y)
7
8 # depth of the decision tree
9 print('Depth of the Decision Tree :', model2.get_depth())
10
```
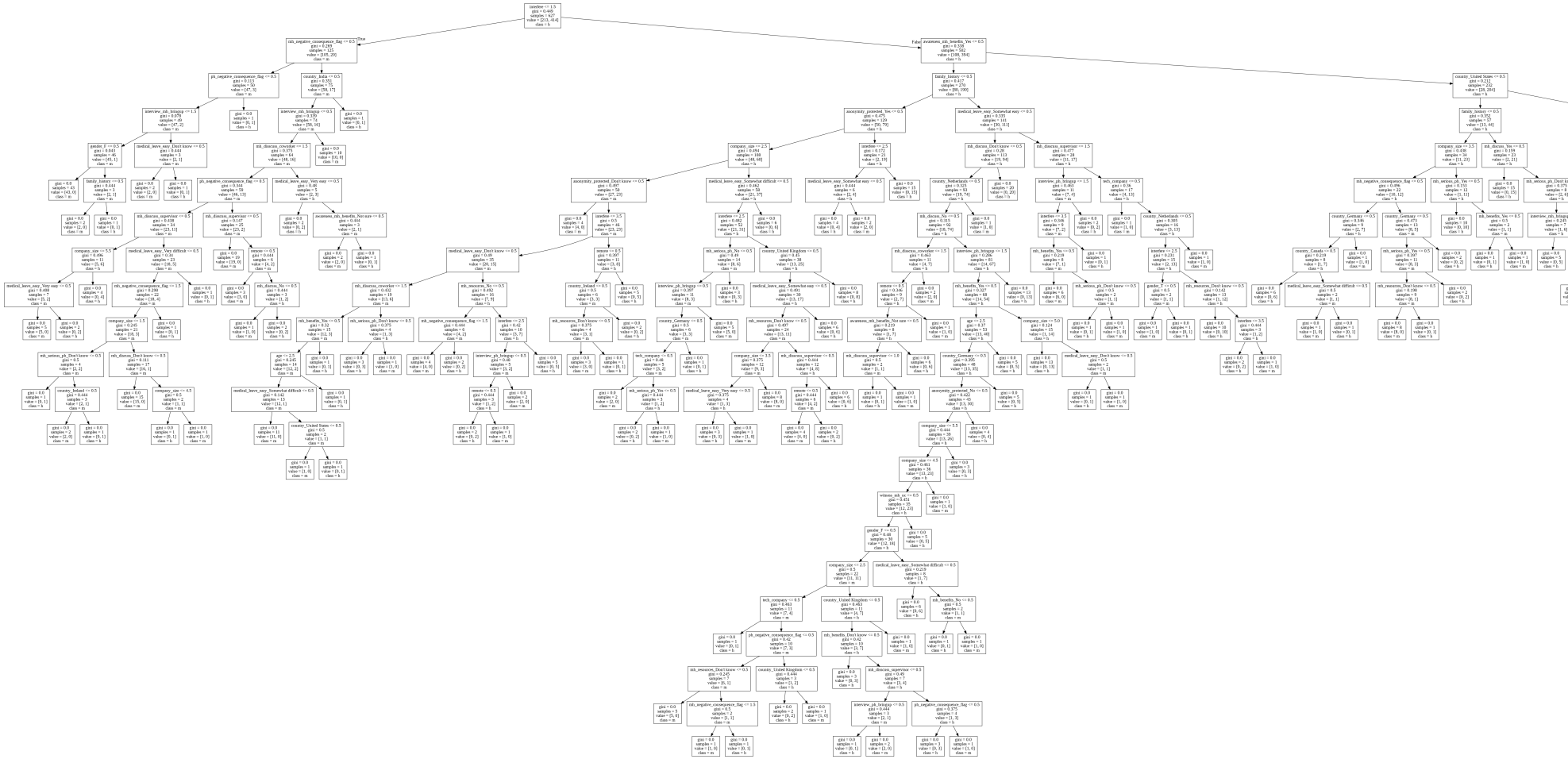
```
Depth of the Decision Tree : 21
```

```
1 from sklearn import datasets
2 from IPython.display import Image
3 from sklearn import tree
4 import pydotplus
5 # Create DOT data
6 dot_data = tree.export_graphviz(model2.fit(train_x,train_y), out_file=None,
7                                 feature_names=train_x.columns,
8                                 class_names=train_y.name)
9
10 # Draw graph
11 graph = pydotplus.graph_from_dot_data(dot_data)
12
13 # Show graph
14 Image(graph.create_png())
```

```
1 # predict the target on the train dataset
2 predict_train2 = model2.predict(train_x)
3
4 # Accuray Score on train dataset
5 accuracy_train2 = accuracy_score(train_y,predict_train2)
6 print('accuracy_score on train dataset : ', accuracy_train2)
7
```

    accuracy_score on train dataset :  1.0

```
1 # predict the target on the test dataset
2 predict_test2 = model2.predict(test_x)
3
4 # Accuracy Score on test dataset
5 accuracy_test2 = accuracy_score(test_y,predict_test2)
6 print('accuracy_score on test dataset : ', accuracy_test2)
```

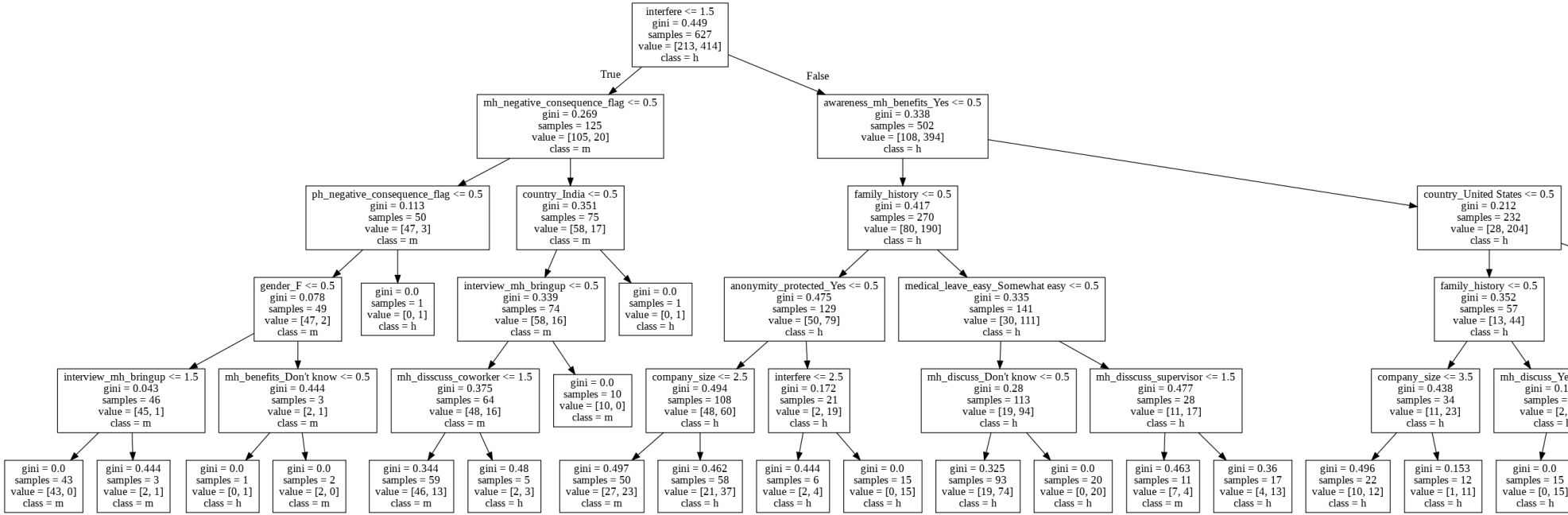    accuracy_score on test dataset :  0.7026022304832714

### 3. Decision Tree Calssifier

```
1 model3 = DecisionTreeClassifier(max_depth=5)
2
3 # fit the model with the training data
4 model3.fit(train_x,train_y)
5
6 # depth of the decision tree
7 print('Depth of the Decision Tree :', model3.get_depth())
```

    Depth of the Decision Tree : 5

```
1 # Create DOT data
2 dot_data3 = tree.export_graphviz(model3.fit(train_x,train_y), out_file=None,
3                                  feature_names=train_x.columns,
4                                  class_names=train_y.name)
5
6 # Draw graph
7 graph3 = pydotplus.graph_from_dot_data(dot_data3)
8
9 # Show graph
10 Image(graph3.create_png())
```

interfere <= 1.5
gini = 0.449
samples = 627
value = [213, 414]
class = h

True — False

mh_negative_consequence_flag <= 0.5
gini = 0.269
samples = 125
value = [105, 20]
class = m

awareness_mh_benefits_Yes <= 0.5
gini = 0.338
samples = 502
value = [108, 394]
class = h

ph_negative_consequence_flag <= 0.5
gini = 0.113
samples = 50
value = [47, 3]
class = m

country_India <= 0.5
gini = 0.351
samples = 75
value = [58, 17]
class = m

family_history <= 0.5
gini = 0.417
samples = 270
value = [80, 190]
class = h

country_United States <= 0.5
gini = 0.212
samples = 232
value = [28, 204]
class = h

gender_F <= 0.5
gini = 0.078
samples = 49
value = [47, 2]
class = m

gini = 0.0
samples = 1
value = [0, 1]
class = h

interview_mh_bringup <= 0.5
gini = 0.339
samples = 74
value = [58, 16]
class = m

gini = 0.0
samples = 1
value = [0, 1]
class = h

anonymity_protected_Yes <= 0.5
gini = 0.475
samples = 129
value = [50, 79]
class = h

medical_leave_easy_Somewhat easy <= 0.5
gini = 0.335
samples = 141
value = [30, 111]
class = h

family_history <= 0.5
gini = 0.352
samples = 57
value = [13, 44]
class = h

interview_mh_bringup <= 1.5
gini = 0.043
samples = 46
value = [45, 1]
class = m

mh_benefits_Don't know <= 0.5
gini = 0.444
samples = 3
value = [2, 1]
class = m

mh_disscuss_coworker <= 1.5
gini = 0.375
samples = 64
value = [48, 16]
class = m

gini = 0.0
samples = 10
value = [10, 0]
class = m

company_size <= 2.5
gini = 0.494
samples = 108
value = [48, 60]
class = h

interfere <= 2.5
gini = 0.172
samples = 21
value = [2, 19]
class = h

mh_discuss_Don't know <= 0.5
gini = 0.28
samples = 113
value = [19, 94]
class = h

mh_disscuss_supervisor <= 1.5
gini = 0.477
samples = 28
value = [11, 17]
class = h

company_size <= 3.5
gini = 0.438
samples = 34
value = [11, 23]
class = h

mh_discuss_Ye...
gini = 0.1...
samples = ...
value = [2, ...]
class = h

gini = 0.0
samples = 43
value = [43, 0]
class = m

gini = 0.444
samples = 3
value = [2, 1]
class = m

gini = 0.0
samples = 1
value = [0, 1]
class = h

gini = 0.0
samples = 2
value = [2, 0]
class = m

gini = 0.344
samples = 59
value = [46, 13]
class = m

gini = 0.48
samples = 5
value = [2, 3]
class = h

gini = 0.497
samples = 50
value = [27, 23]
class = m

gini = 0.462
samples = 58
value = [21, 37]
class = h

gini = 0.444
samples = 6
value = [2, 4]
class = h

gini = 0.0
samples = 15
value = [0, 15]
class = h

gini = 0.325
samples = 93
value = [19, 74]
class = h

gini = 0.0
samples = 20
value = [0, 20]
class = h

gini = 0.463
samples = 11
value = [7, 4]
class = m

gini = 0.36
samples = 17
value = [4, 13]
class = h

gini = 0.496
samples = 22
value = [10, 12]
class = h

gini = 0.153
samples = 12
value = [1, 11]
class = h

gini = 0.0
samples = 15
value = [0, 15]
class = h

```
1 # predict the target on the train dataset
2 predict_train3 = model3.predict(train_x)
3
4 # Accuray Score on train dataset
5 accuracy_train3 = accuracy_score(train_y,predict_train3)
6 print('accuracy_score on train dataset : ', accuracy_train3)
```

accuracy_score on train dataset :  0.8133971291866029

```
1 # predict the target on the test dataset
2 predict_test3 = model3.predict(test_x)
3
4 # Accuracy Score on test dataset
5 accuracy_test3 = accuracy_score(test_y,predict_test3)
6 print('accuracy_score on test dataset : ', accuracy_test3)
```

accuracy_score on test dataset :  0.758364312267658

## 4. SVM

```
1 from sklearn.svm import SVC
2 model4 = SVC()
3
4 # fit the model with the training data
5 model4.fit(train_x,train_y)
```

SVC(C=1.0, break_ties=False, cache_size=200, class_weight=None, coef0=0.0,
    decision_function_shape='ovr', degree=3, gamma='scale', kernel='rbf',
    max_iter=-1, probability=False, random_state=None, shrinking=True,
    tol=0.001, verbose=False)

```
1 # predict the target on the train dataset
2 predict_train4 = model4.predict(train_x)
3
4 # Accuray Score on train dataset
5 accuracy_train4 = accuracy_score(train_y,predict_train4)
6 print('accuracy_score on train dataset : ', accuracy_train4)
```

accuracy_score on train dataset :  0.8181818181818182

```
1 # predict the target on the test dataset
2 predict_test4 = model4.predict(test_x)
3
4 # Accuracy Score on test dataset
5 accuracy_test4 = accuracy_score(test_y,predict_test4)
6 print('accuracy_score on test dataset : ', accuracy_test4)
```

accuracy_score on test dataset :  0.8066914498141264

## 5. KNN

```
1 from sklearn.neighbors import KNeighborsClassifier
2 model5 = KNeighborsClassifier()
3
4 # fit the model with the training data
5 model5.fit(train_x,train_y)
```

```
5 model5.fit(train_x,train_y)
6
7 # Number of Neighbors used to predict the target
8 print('\nThe number of neighbors used to predict the target : ',model5.n_neighbors)
```

> The number of neighbors used to predict the target :  5

```
1 # predict the target on the train dataset
2 predict_train5 = model5.predict(train_x)
3
4 # Accuray Score on train dataset
5 accuracy_train5 = accuracy_score(train_y,predict_train5)
6 print('accuracy_score on train dataset : ', accuracy_train5)
```

> accuracy_score on train dataset :  0.8197767145135566

```
1 # predict the target on the test dataset
2 predict_test5 = model5.predict(test_x)
3
4 # Accuracy Score on test dataset
5 accuracy_test5 = accuracy_score(test_y,predict_test5)
6 print('accuracy_score on test dataset : ', accuracy_test5)
```

> accuracy_score on test dataset :  0.7211895910780669

### 6. GBM

```
1 from sklearn.ensemble import GradientBoostingClassifier
2 model6 = GradientBoostingClassifier(n_estimators=100,max_depth=5)
3
4 # fit the model with the training data
5 model6.fit(train_x,train_y)
```

> GradientBoostingClassifier(ccp_alpha=0.0, criterion='friedman_mse', init=None,
>                            learning_rate=0.1, loss='deviance', max_depth=5,
>                            max_features=None, max_leaf_nodes=None,
>                            min_impurity_decrease=0.0, min_impurity_split=None,
>                            min_samples_leaf=1, min_samples_split=2,
>                            min_weight_fraction_leaf=0.0, n_estimators=100,
>                            n_iter_no_change=None, presort='deprecated',
>                            random_state=None, subsample=1.0, tol=0.0001,
>                            validation_fraction=0.1, verbose=0,
>                            warm_start=False)

```
1 # predict the target on the train dataset
2 predict_train6 = model6.predict(train_x)
3
4 # Accuray Score on train dataset
5 accuracy_train6 = accuracy_score(train_y,predict_train6)
6 print('\naccuracy_score on train dataset : ', accuracy_train6)
```

> accuracy_score on train dataset :  0.9904306220095693

```
1 # predict the target on the test dataset
2 predict_test6 = model6.predict(test_x)
3
4 # Accuracy Score on test dataset
5 accuracy_test6 = accuracy_score(test_y,predict_test6)
6 print('\naccuracy_score on test dataset : ', accuracy_test6)
```
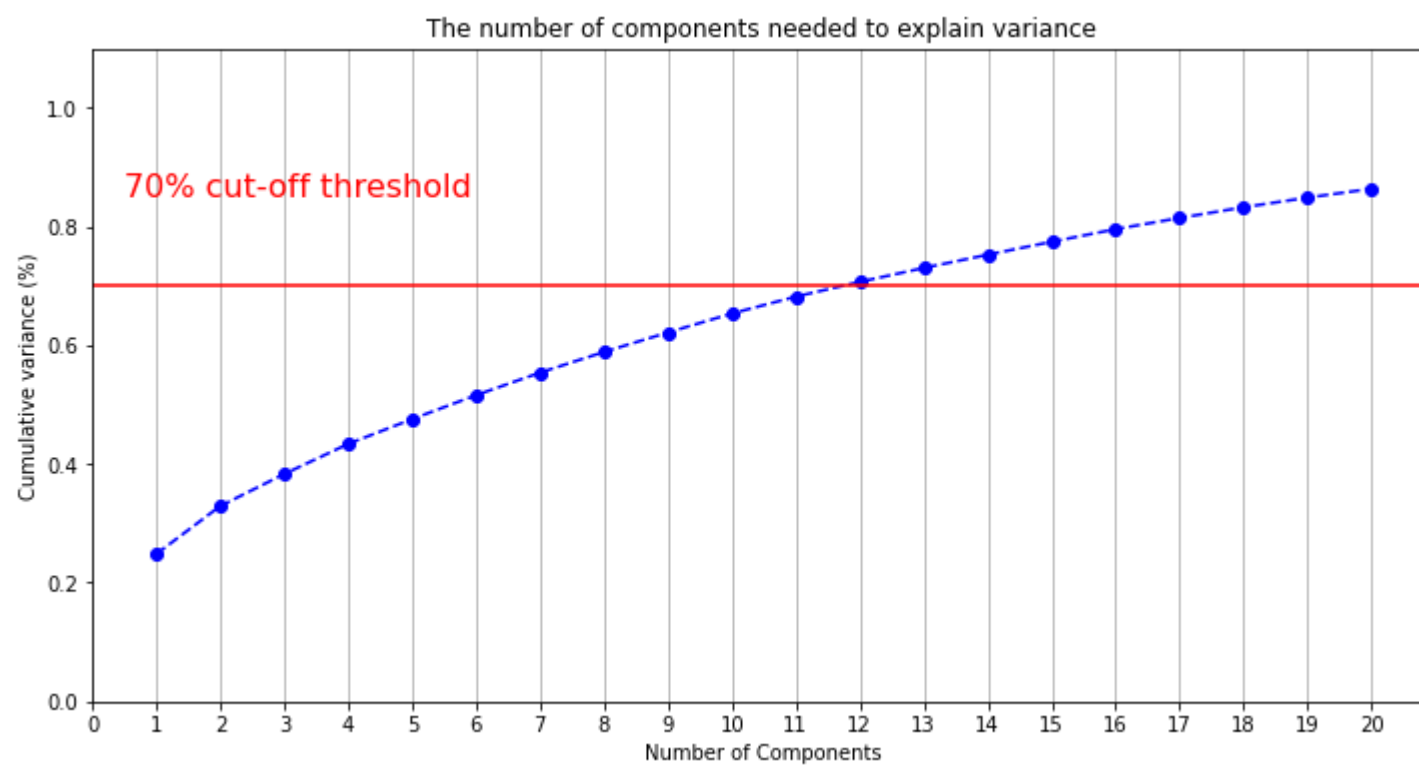
> accuracy_score on test dataset :  0.724907063197026

### 7. XGBoost

```
1 from xgboost import XGBClassifier
2 model7 = XGBClassifier()
3
4 # fit the model with the training data
5 model7.fit(train_x,train_y)
```

>

```
XGBClassifier(base_score=0.5, booster='gbtree', colsample_bylevel=1,
              colsample_bynode=1, colsample_bytree=1, gamma=0,
              learning_rate=0.1, max_delta_step=0, max_depth=3,
              min_child_weight=1, missing=None, n_estimators=100, n_jobs=1,
              nthread=None, objective='binary:logistic', random_state=0,
              reg_alpha=0, reg_lambda=1, scale_pos_weight=1, seed=None,
              silent=None, subsample=1, verbosity=1)
```

```
1 # predict the target on the train dataset
2 predict_train7 = model7.predict(train_x)
3
4 # Accuray Score on train dataset
5 accuracy_train7 = accuracy_score(train_y,predict_train7)
6 print('\naccuracy_score on train dataset : ', accuracy_train7)
```

accuracy_score on train dataset :  0.8389154704944178

```
1 # predict the target on the test dataset
2 predict_test7 = model7.predict(test_x)
3
4 # Accuracy Score on test dataset
5 accuracy_test7 = accuracy_score(test_y,predict_test7)
6 print('\naccuracy_score on test dataset : ', accuracy_test7)
```

accuracy_score on test dataset :  0.7695167286245354

Unsupervised ML

PCA + kmeans

```
1 #scaler data
2 from sklearn.preprocessing import MinMaxScaler
3 scaler = MinMaxScaler()
4 data_rescaled = scaler.fit_transform(df)
```

```
1 # find the suitable number of components
2 from sklearn.decomposition import PCA
3 pca = PCA().fit(data_rescaled)
4
5 % matplotlib inline
6 import matplotlib.pyplot as plt
7 plt.rcParams["figure.figsize"] = (12,6)
8
9 fig, ax = plt.subplots()
10 xi = np.arange(1, 21, step=1)
11 y = np.cumsum(pca.explained_variance_ratio_)[1:21]
12
13 plt.ylim(0.0,1.1)
14 plt.plot(xi, y, marker='o', linestyle='--', color='b')
15
16 plt.xlabel('Number of Components')
17 plt.xticks(np.arange(0, 21, step=1))
18 plt.ylabel('Cumulative variance (%)')
19 plt.title('The number of components needed to explain variance')
20
21 plt.axhline(y=0.70, color='r', linestyle='-')
22 plt.text(0.5, 0.85, '70% cut-off threshold', color = 'red', fontsize=16)
23
24 ax.grid(axis='x')
25 plt.show()
```

The number of components needed to explain variance



In order to get aroud 70%, we decided to choose the first 12 components.

```
1 #fit the model
2 pca_final = PCA(n_components = 12)
3 pca_final.fit(data_rescaled)
4 data_pca = pca_final.fit_transform(data_rescaled)
5 data_pca.shape
```
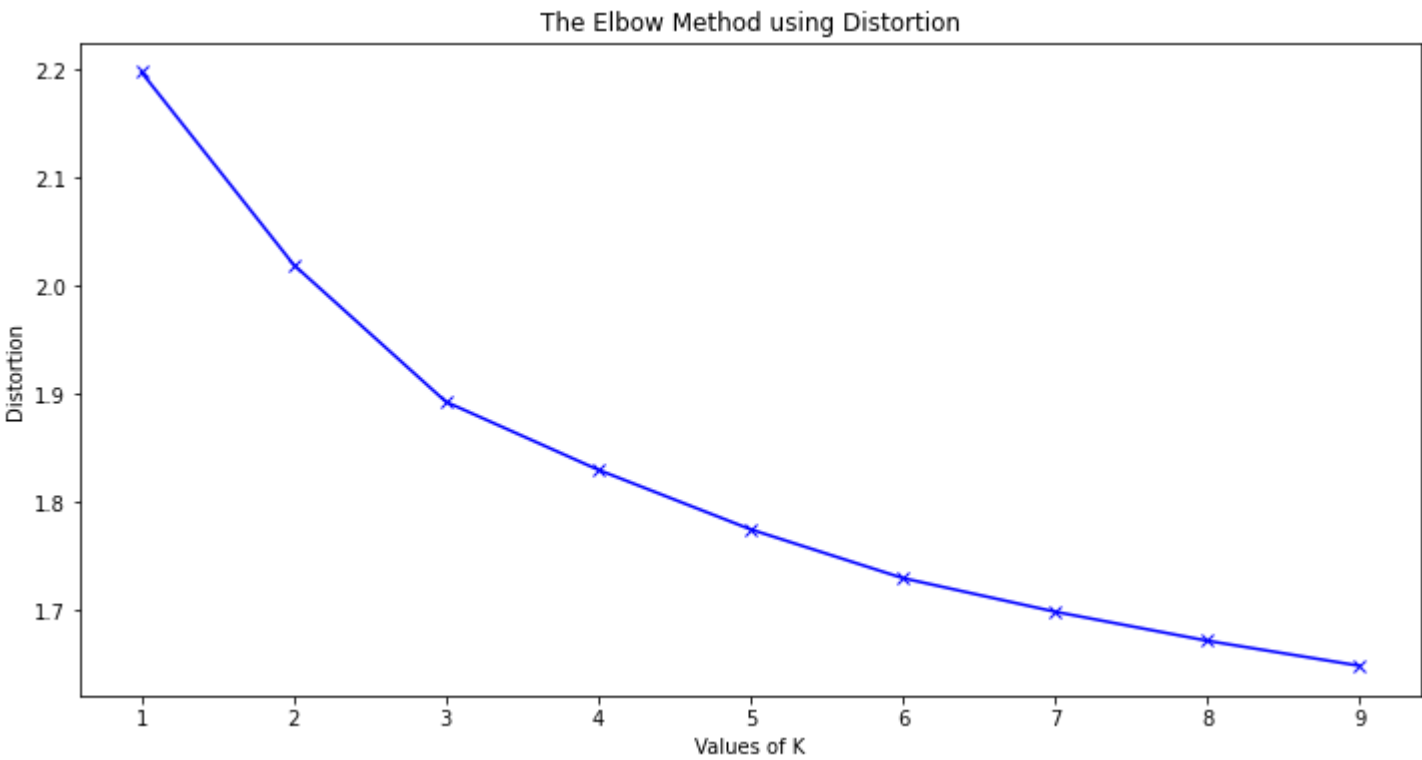
```
(896, 12)
```

```
 1 #find reasonable number of k with elbow method
 2 from sklearn.cluster import KMeans
 3 from sklearn import metrics
 4 from scipy.spatial.distance import cdist
 5 import numpy as np
 6 import matplotlib.pyplot as plt
 7
 8 distortions = []
 9 inertias = []
10 mapping1 = {}
11 mapping2 = {}
12 K = range(1,10)
13
14 for k in K:
15     #Building and fitting the model
16     kmeanModel = KMeans(n_clusters=k).fit(data_pca)
17     kmeanModel.fit(data_pca)
18
19     distortions.append(sum(np.min(cdist(data_pca, kmeanModel.cluster_centers_,
20                     'euclidean'),axis=1)) / data_pca.shape[0])
21     inertias.append(kmeanModel.inertia_)
22
23     mapping1[k] = sum(np.min(cdist(data_pca, kmeanModel.cluster_centers_,
24                 'euclidean'),axis=1)) / data_pca.shape[0]
25     mapping2[k] = kmeanModel.inertia_
26
```
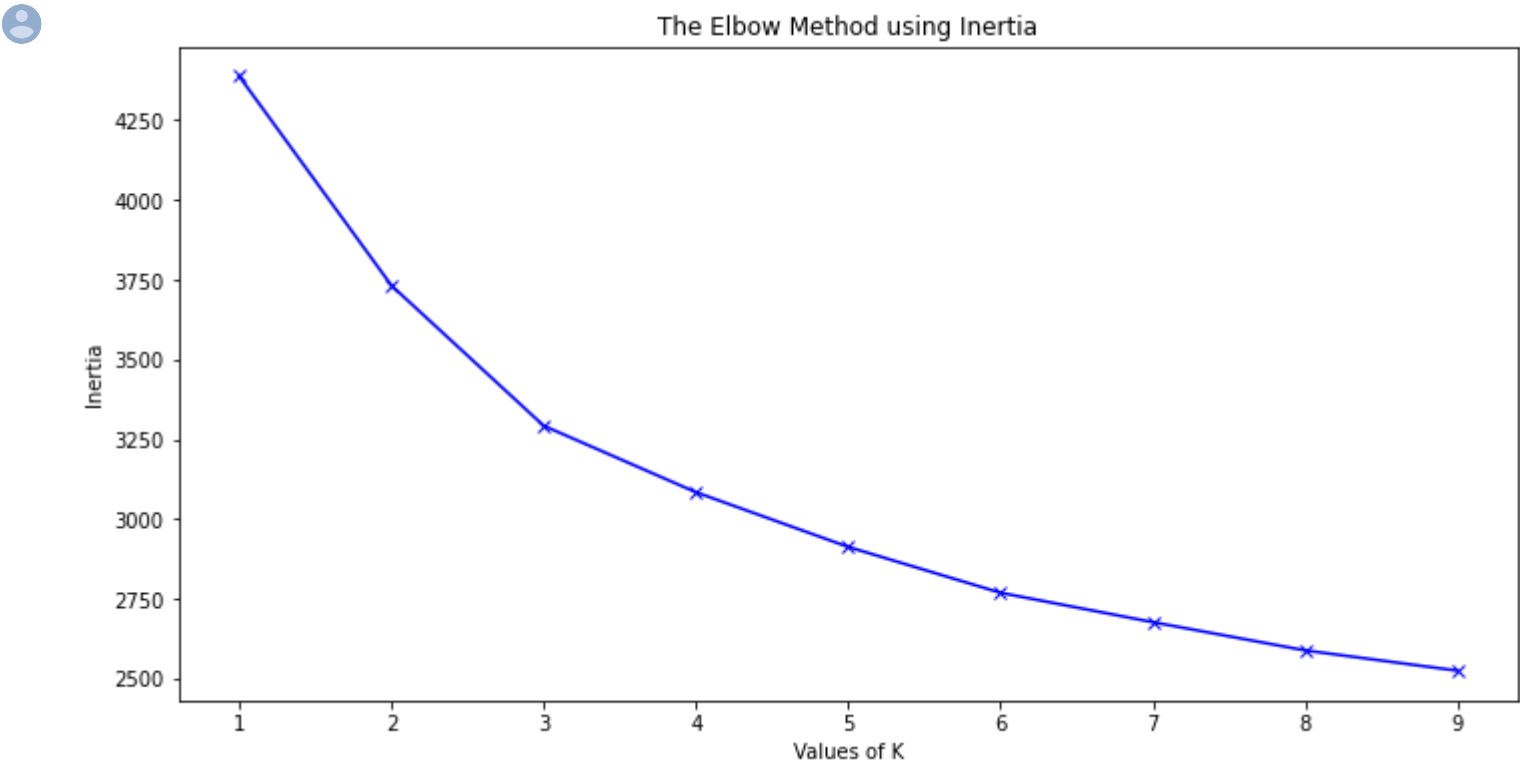
```
1 plt.plot(K, distortions, 'bx-')
2 plt.xlabel('Values of K')
3 plt.ylabel('Distortion')
4 plt.title('The Elbow Method using Distortion')
5 plt.show()
```

The Elbow Method using Distortion

```
1 plt.plot(K, inertias, 'bx-')
2 plt.xlabel('Values of K')
3 plt.ylabel('Inertia')
4 plt.title('The Elbow Method using Inertia')
5 plt.show()
```



The Elbow Method using Inertia
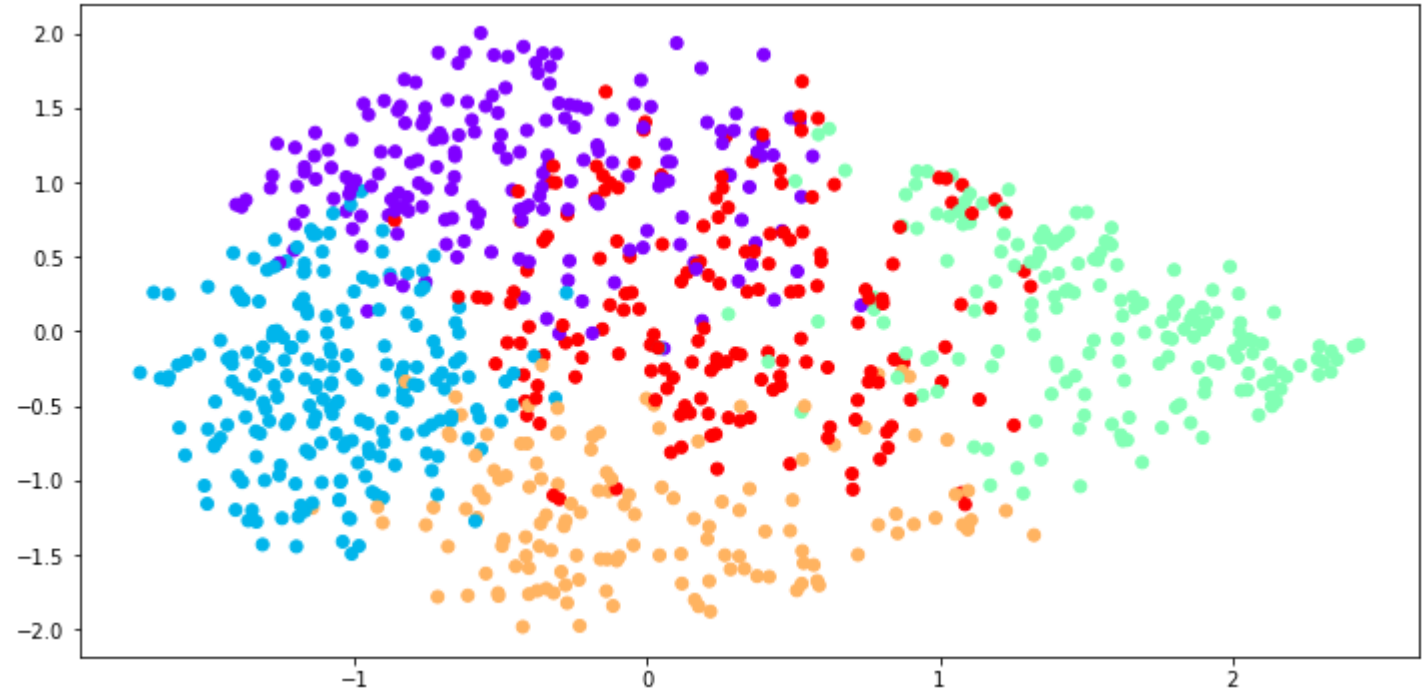
Based on the results, we chose to use k=5.

```
1 #fit the kmeans model
2 kmeans= KMeans(n_clusters=5)
3 kmeans5=kmeans.fit_predict(data_pca)
```

```
1 kmeans.cluster_centers_
```

```
array([[-0.4429739 ,  1.05666511,  0.43913159,  0.28931182,  0.11870278,
        -0.09887435, -0.00617988, -0.07980882, -0.07975342, -0.05690273,
         0.12475788,  0.15104256],
       [-1.07536002, -0.34658112, -0.08408743, -0.11913515,  0.10029392,
         0.18634079, -0.02428671,  0.00897236,  0.06673901,  0.05804446,
        -0.03544   , -0.17652532],
       [ 1.53302396,  0.07519759,  0.0984913 , -0.14740826,  0.2511941 ,
         0.16342133, -0.03881158, -0.01913487, -0.05577262,  0.09761498,
         0.04129115, -0.1241347 ],
       [-0.00363126, -1.22165991,  0.41078203,  0.09182208, -0.23131532,
        -0.18612103,  0.15502751, -0.00565208, -0.06049708, -0.10816998,
         0.06534425,  0.0312912 ],
       [ 0.26556962,  0.15832695, -0.79598708, -0.08762594, -0.33144407,
        -0.14719943, -0.04429081,  0.0998988 ,  0.10781135, -0.02589955,
        -0.18514208,  0.15876326]])
```

```
1 #visualizing clustering
2 plt.scatter(data_pca[:,0],data_pca[:,1],c=kmeans5,cmap="rainbow")
```

```
<matplotlib.collections.PathCollection at 0x7f457c7b1198>
```



```
1 # Get cluster assignment labels
2 labels = kmeans.labels_
3 df['Group']=labels
4 df.head()
```

| | age | self_employed | family_history | mh_treatment | interfere | company_size | remote | tech_company | mh_negative_consequ |
|---|---|---|---|---|---|---|---|---|---|
| **0** | 3 | 1 | 1 | 0 | 3 | 1 | 1 | 1 | |
| **1** | 2 | 0 | 1 | 1 | 3 | 4 | 0 | 1 | |
| **2** | 2 | 1 | 0 | 0 | 1 | 1 | 1 | 1 | |
| **3** | 3 | 0 | 0 | 1 | 4 | 3 | 1 | 1 | |
| **4** | 3 | 0 | 0 | 1 | 1 | 6 | 0 | 0 | |

```
1 df.groupby('Group').describe()
```

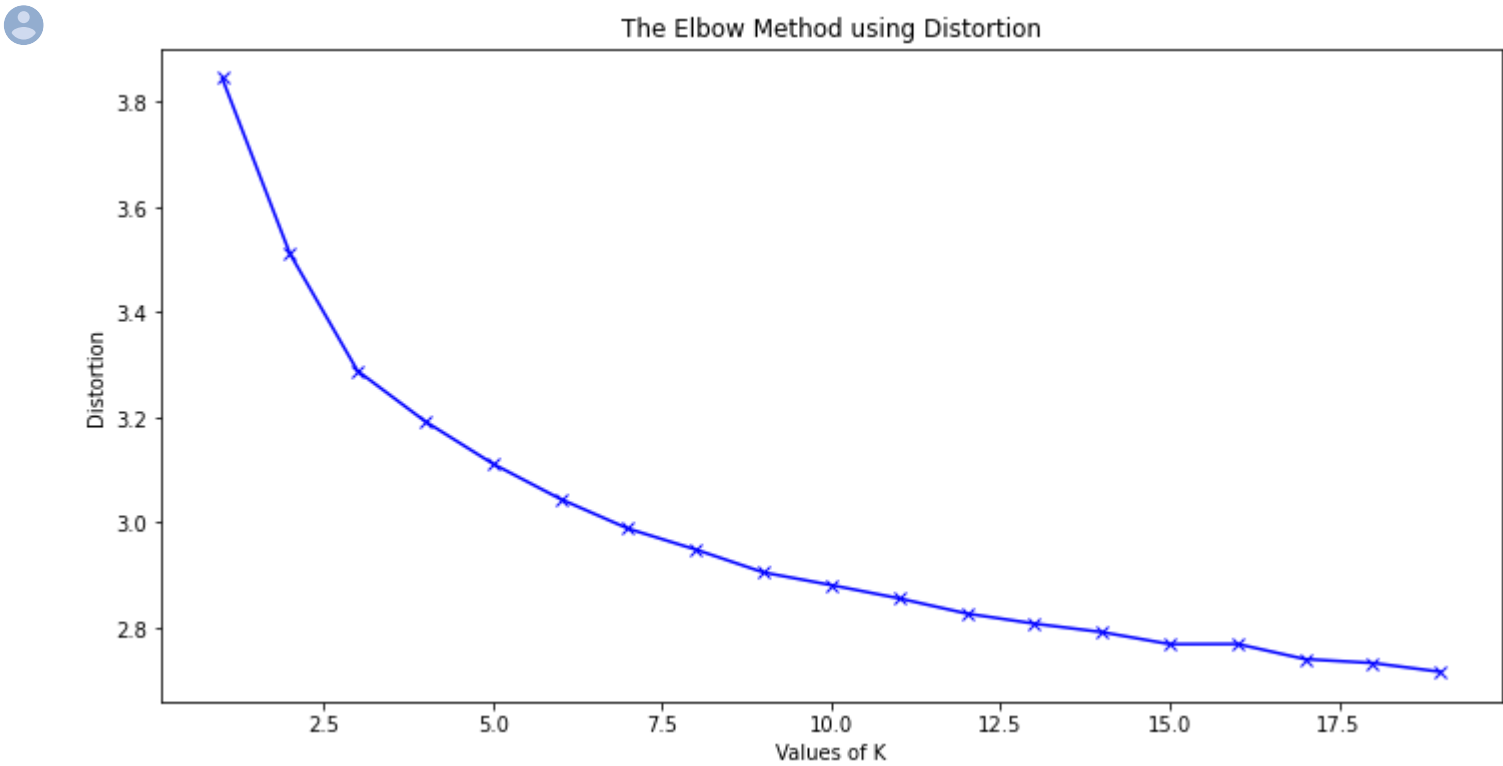| | age | | | | | | | self_employed | | | | | | | family_history | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| | count | mean | std | min | 25% | 50% | 75% | max | count | mean | std | min | 25% | 50% | 75% | max | count | mean |
| **Group** | | | | | | | | | | | | | | | | |
| **0** | 190.0 | 2.100000 | 0.378664 | 1.0 | 2.0 | 2.0 | 2.0 | 3.0 | 190.0 | 0.100000 | 0.300793 | 0.0 | 0.0 | 0.0 | 0.0 | 1.0 | 190.0 | 0.384211 |
| **1** | 218.0 | 2.059633 | 0.347637 | 1.0 | 2.0 | 2.0 | 2.0 | 3.0 | 218.0 | 0.114679 | 0.319367 | 0.0 | 0.0 | 0.0 | 0.0 | 1.0 | 218.0 | 0.371560 |
| **2** | 178.0 | 2.207865 | 0.433805 | 2.0 | 2.0 | 2.0 | 2.0 | 4.0 | 178.0 | 0.061798 | 0.241467 | 0.0 | 0.0 | 0.0 | 0.0 | 1.0 | 178.0 | 0.500000 |
| **3** | 136.0 | 2.080882 | 0.323297 | 1.0 | 2.0 | 2.0 | 2.0 | 3.0 | 136.0 | 0.367647 | 0.483947 | 0.0 | 0.0 | 0.0 | 1.0 | 1.0 | 136.0 | 0.433824 |
| **4** | 174.0 | 2.132184 | 0.339668 | 2.0 | 2.0 | 2.0 | 2.0 | 3.0 | 174.0 | 0.011494 | 0.106901 | 0.0 | 0.0 | 0.0 | 0.0 | 1.0 | 174.0 | 0.637931 |

5 rows × 400 columns

```
1 df.groupby('Group').mean()
```

| | age | self_employed | family_history | mh_treatment | interfere | company_size | remote | tech_company | mh_negativ |
|---|---|---|---|---|---|---|---|---|---|
| **Group** | | | | | | | | | |
| **0** | 2.100000 | 0.100000 | 0.384211 | 0.505263 | 2.257895 | 3.573684 | 0.315789 | 0.821053 | |
| **1** | 2.059633 | 0.114679 | 0.371560 | 0.481651 | 2.481651 | 2.834862 | 0.288991 | 0.834862 | |
| **2** | 2.207865 | 0.061798 | 0.500000 | 0.735955 | 2.477528 | 4.516854 | 0.224719 | 0.707865 | |
| **3** | 2.080882 | 0.367647 | 0.433824 | 0.727941 | 2.801471 | 2.066176 | 0.514706 | 0.897059 | |
| **4** | 2.132184 | 0.011494 | 0.637931 | 0.890805 | 2.816092 | 4.097701 | 0.229885 | 0.816092 | |

```
1 df.groupby('Group').mean().rank()
```

| | age | self_employed | family_history | mh_treatment | interfere | company_size | remote | tech_company | mh_negative_con |
|---|---|---|---|---|---|---|---|---|---|
| **Group** | | | | | | | | | |
| **0** | 3.0 | 3.0 | 2.0 | 2.0 | 1.0 | 3.0 | 4.0 | 3.0 | |
| **1** | 1.0 | 4.0 | 1.0 | 1.0 | 3.0 | 2.0 | 3.0 | 4.0 | |
| **2** | 5.0 | 2.0 | 4.0 | 4.0 | 2.0 | 5.0 | 1.0 | 1.0 | |
| **3** | 2.0 | 5.0 | 3.0 | 3.0 | 4.0 | 1.0 | 5.0 | 5.0 | |
| **4** | 4.0 | 1.0 | 5.0 | 5.0 | 5.0 | 4.0 | 2.0 | 2.0 | |

## PCA

```
1 distortions = []
2 inertias = []
3 mapping3 = {}
4 mapping4 = {}
5 K2 = range(1,20)
6
7 for k in K2:
8     #Building and fitting the model
9     kmeanModel2 = KMeans(n_clusters=k).fit(df)
10    kmeanModel2.fit(df)
11
12    distortions.append(sum(np.min(cdist(df, kmeanModel2.cluster_centers_,
13                      'euclidean'),axis=1)) / df.shape[0])
14    inertias.append(kmeanModel2.inertia_)
15
16    mapping3[k] = sum(np.min(cdist(df, kmeanModel2.cluster_centers_,
17                  'euclidean'),axis=1)) / df.shape[0]
18    mapping4[k] = kmeanModel2.inertia_
```

```
1 plt.plot(K2, distortions, 'bx-')
2 plt.xlabel('Values of K')
3 plt.ylabel('Distortion')
4 plt.title('The Elbow Method using Distortion')
5 plt.show()
```



```
1 plt.plot(K2, inertias, 'bx-')
2 plt.xlabel('Values of K')
3 plt.ylabel('Inertia')
4 plt.title('The Elbow Method using Inertia')
5 plt.show()
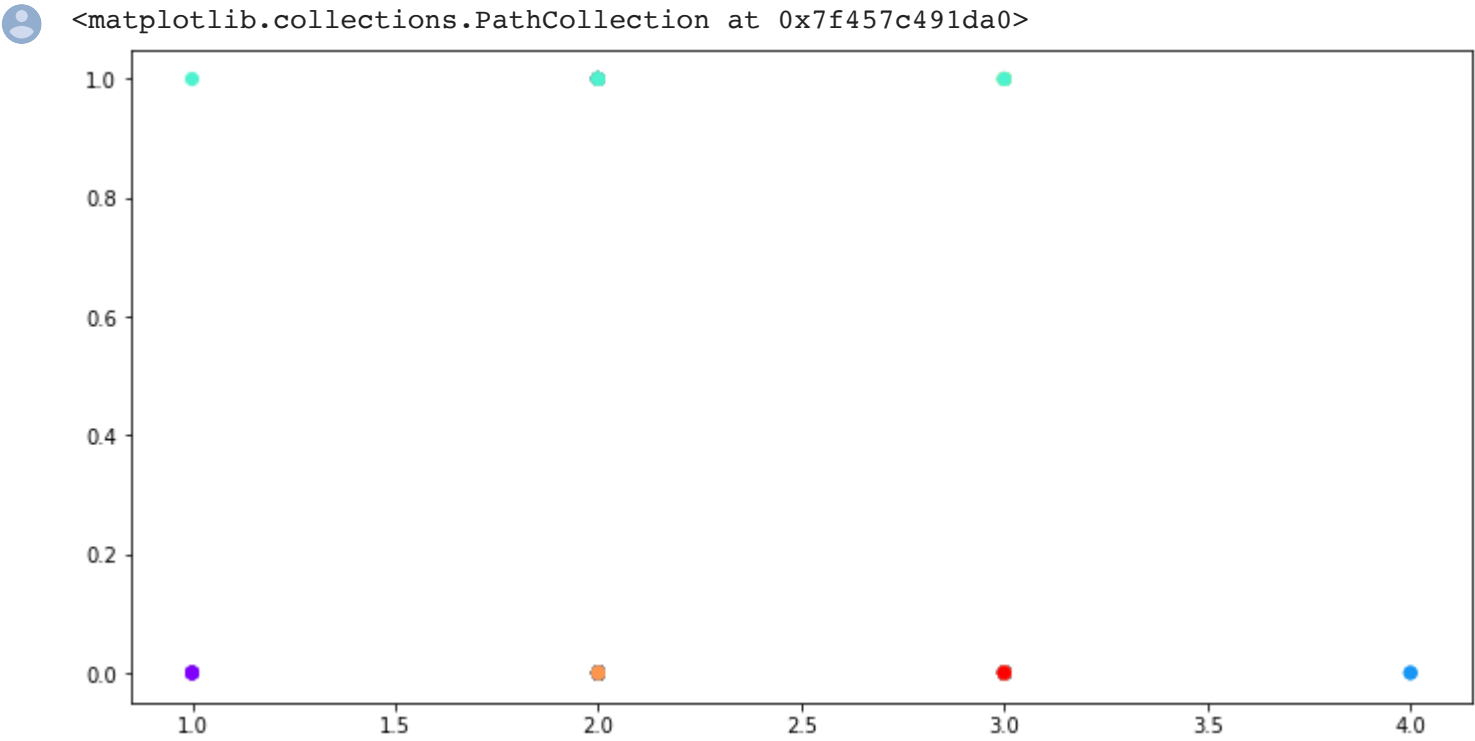```

The Elbow Method using Inertia



```
1 We decided to use 6 clusters.
```

```
1 #fit the kmeans model
2 kmeans2= KMeans(n_clusters=6)
3 kmeans6=kmeans2.fit_predict(df)
```

```
1 df=np.array(df)
```

```
1 #visualizing clustering
2 plt.scatter(df[:,0],df[:,1], c=kmeans6, cmap='rainbow')
```

<matplotlib.collections.PathCollection at 0x7f457c491da0>



The reason why we got imperfect plot is that there are two many dimensions here.