

In [147]: ►

```
1 # Running the Libraries
2 from pdf2image import convert_from_path # for converting PDF to JPG
3 from matplotlib.image import imread
4 import numpy as np
5 import pandas as pd
6 import scipy as sc
7 import os
8 import matplotlib.pyplot as plt
9 from pylab import *
10 from skimage.transform import resize # for resizing the images
11 import statistics as stats
12
13 from tensorflow import keras
14 from sklearn.model_selection import train_test_split
15 from sklearn.metrics import classification_report, confusion_matrix
16 from sklearn.metrics import recall_score, f1_score, precision_score
17 from keras.models import load_model
18 from keras.layers.convolutional import Conv2D, MaxPooling2D
19 from keras.layers import Dense, Dropout, Flatten, Conv2D, MaxPool2D, BatchNormalization
20 from keras.models import Sequential
21 from pdf2image.exceptions import (PDFInfoNotInstalledError, PDFPageCountError)
22
23 # Packages for ANOVA Table
24 import statsmodels.api as sm
25 from statsmodels.formula.api import ols
26 import scikit_posthocs as sp
27 from scipy.stats import f_oneway
28 from scipy.stats import levene
29
30 #Chi-Square
31 from bioinfokit.analys import stat, get_data
32
33 # Fisher-Exact Test
34 import rpy2.robjects.numpy2ri
35 from rpy2.robjects.packages import importr
36 rpy2.robjects.numpy2ri.activate()
37
38 # KNN
39 from sklearn.neighbors import KNeighborsRegressor
40 from sklearn.metrics import mean_squared_error
41 from math import sqrt
42
43 # CNN
44 from sklearn.model_selection import train_test_split
45 from sklearn.metrics import confusion_matrix
46 from tensorflow import keras
47 from sklearn.model_selection import GridSearchCV
48 from scikeras.wrappers import KerasClassifier
49
50 # K-Means
51 from sklearn.cluster import KMeans
52
53 # K-Medoids
54 from sklearn_extra.cluster import KMedoids
```

We will import the dataset containing the FEV1, GOLD Stage, COPD Score, & Adjusted Clock Drawing Test Score for 1,716 unique Subject IDs.

```
In [16]: ► 1 Dataset_1 = pd.read_csv('C:/Users/12563/Documents/Res. Methods in Math &  
2 Dataset_1
```

Out[16]:

	sid	fev1pp_final	goldstage_final_V2	COPD_P3	Score
0	10015T	51.1	2	1	2
1	10017X	55.3	2	1	2
2	10031R	39.5	2	0	2
3	10032T	58.8	2	1	0
4	10049K	59.5	2	1	2
...
1712	23634Q	73.9	1	0	2
1713	23696M	72.4	2	1	2
1714	23785L	51.7	2	1	2
1715	23793K	108.6	0	0	2
1716	24637B	119.4	2	0	2

1717 rows × 5 columns

```
In [17]: ► 1 # Rename column in a dataframe  
2 Dataset_1.rename(columns = {'sid':'Clocks'}, inplace = True)
```

We will perform the K-Nearest Neighbors Algorithm to solve the classification problem of whether or not a particular person (Subject ID) has COPD or not.

```
In [4]: 1 df_knn = pd.DataFrame(Dataset_1, columns = [ 'fev1pp_final', 'goldstage_final_V2', 'COPD_P3' ]  
2 df_knn
```

Out[4]:

	fev1pp_final	goldstage_final_V2	COPD_P3	Score
0	51.1	2	1	2
1	55.3	2	1	2
2	39.5	2	0	2
3	58.8	2	1	0
4	59.5	2	1	2
...
1712	73.9	1	0	2
1713	72.4	2	1	2
1714	51.7	2	1	2
1715	108.6	0	0	2
1716	119.4	2	0	2

1717 rows × 4 columns

```
In [6]: 1 X_1 = df_knn.drop("COPD_P3", axis = 1)  
2 X_1 = X_1.values  
3 y_1 = df_knn['COPD_P3']  
4 y_1 = y_1.values
```

```
In [7]: 1 # Split the data into train and test sets  
2 X_train_1, X_test_1, y_train_1, y_test_1 = train_test_split(X_1, y_1, test_size=0.2, random_state=42)
```

```
In [8]: 1 # KNN is a supervised machine Learning method  
2 # Trying to predict whether or not someone will have COPD based  
3 # on their CDT Score (the predictor)  
4 knn_model_1 = KNeighborsRegressor(n_neighbors=2)  
5 knn_model_1.fit(X_train_1, y_train_1)
```

Out[8]:

```
▼ KNeighborsRegressor  
KNeighborsRegressor(n_neighbors=2)
```

```
In [9]: 1 # Calculate RMSE for Training Data  
2 train_preds_1 = knn_model_1.predict(X_train_1)  
3 mse_1 = mean_squared_error(y_train_1, train_preds_1)  
4 rmse_1 = sqrt(mse_1)  
5 rmse_1
```

Out[9]: 0.30741029992207347

```
In [10]: # Calculate RMSE for Testing Data
          2 test_preds_1_t = knn_model_1.predict(X_test_1)
          3 mse_1_t = mean_squared_error(y_test_1, test_preds_1_t)
          4 rmse_1_test = sqrt(mse_1_t)
          5 rmse_1_test
```

Out[10]: 0.4985443928356902

```
In [11]: # Calculate Accuracy
          2 y_pred_1 = knn_model_1.predict(X_test_1)
          3
          4 # Evaluate the training & testing Accuracy
          5 #
          6 print('Training accuracy score: %.3f' % knn_model_1.score(X_train_1, y_t
          7 print('Test accuracy score: %.3f' % knn_model_1.score(X_test_1, y_test_1)
```

Training accuracy score: 0.593

Test accuracy score: -0.132

Let us first import the Excel file containing 19,485 unique Subject IDs (SIDs) along with the associated Adjusted Clock Drawing Test Score (either 0 (High-Risk for Dementia) or 2 (not at High-Risk for Dementia)).

```
In [12]: # Reading clock scores from CSV score file
          2
          3 df = pd.read_csv('C:/Users/12563/Desktop/CDT/Adjudicated Mini-Cog Clock S
          4 scores1 = df.shape
          5 df = df[df['Adjudicated Clock Score'].notnull()]
          6 scores2 = df.shape
          7 df['Adjudicated Clock Score'] = df['Adjudicated Clock Score'].apply(lambda
          8 scores3 = df.shape
```

Now, we will create the Clock_List data frame that contains the SIDs of the Clock Drawing Test Images with the corresponding Adjusted Clock Score from the Excel file.

In [13]:

```

1 # Reading the image files from the source folder
2 path = r'C:/Users/12563/Documents/Test Clocks 2/'
3 files = os.listdir(path)
4
5
6 # The code below extracts the patient_id from the image names and match it
7 # add it to the file name (format : score_patientID)
8 scores = []
9 name_sel = []
10
11 nummatch = 0
12 numdiscard = 0
13
14 for i in files:
15
16     # Depending on the image file names different split patterns are needed
17     f_name_prim = os.path.basename(i).split('.')[0].split(' ')[0].split('_')[0]
18     #     f_name_prim = os.path.basename(i).split('.')[0].split(' ')[0].split('_')[0]
19     #     f_name_prim = os.path.basename(i).split(' ')[1].split('.')[0]
20     #     f_name_prim = os.path.basename(i).split(' ')[1].split('-')[1].split('.')[0]
21     #     f_name_prim = os.path.basename(i).split(' ')[1].split('#')[1].split('.')[0]
22     #     f_name_prim = os.path.basename(i).split(' ')[1].split('.')[0]
23     #     f_name_prim = os.path.basename(i).split(' ')[0].split('_')[1]
24     #     f_name_prim = os.path.basename(i).split(' ')[1].split('.')[0]
25     #     f_name_prim = os.path.basename(i).split('_')[1].split('_')[0]
26
27     for k in range(df.shape[0]):
28         if f_name_prim == df.iloc[k,0]:
29             nummatch += 1
30             scores.append(df.iloc[k,1])
31             name_sel.append(f_name_prim)
32
33     #         os.rename(os.path.join(path, i), os.path.join(path, ''.join(name_sel)))
34     else:
35         numdiscard += 1
36
37
38 Clock_list = pd.DataFrame([name_sel, scores], index = ['Clocks', 'Scores'])
39 Clock_list = Clock_list.T
40
41 # saving the results in the target folder
42 Clock_list.to_csv(path + 'Clock_list.csv')
43 Clock_list

```

Out[13]:

	Clocks	Scores
0	10071D	0
1	10098X	0
2	10111P	0
3	10286Y	0
4	10292T	0

	Clocks	Scores
...
538	30352X	2
539	30357H	2
540	30361Y	2
541	30363C	2
542	30368M	2

543 rows × 2 columns

Furthermore, we will create a new data frame that will merge the Clock_List Data Frame with the excel file that contains the associated FEV1 Score, GOLD Stage, & COPD Score for 370 of the the SIDs in which we have the Clock Drawing Test image & the associated score.

In [18]: ► 1 Dataset_2 = Clock_list.merge(Dataset_1,on = 'Clocks')
2 Dataset_2

Out[18]:

	Clocks	Scores	fev1pp_final	goldstage_final_V2	COPD_P3	Score
0	10071D	0	42.7	2	1	0
1	10098X	0	120.3	2	0	0
2	10111P	0	19.6	2	1	0
3	10286Y	0	77.7	2	0	0
4	10292T	0	35.6	2	1	0
...
365	23619U	2	83.6	2	0	2
366	23634Q	2	73.9	1	0	2
367	23696M	2	72.4	2	1	2
368	23785L	2	51.7	2	1	2
369	24637B	2	119.4	2	0	2

370 rows × 6 columns

We will now perform Statistical Analysis on Dataset 2, this Data Frame created by Merging Dataset 1 & Clock_List.

- 1 We begin the statistical Analysis by performing the Chi-Squared Test for Independence on the Adjusted Clock Drawing Test Score by the Self-Reported COPD Score

```
In [20]: ► 1 # Chi-Square Test  
2  
3 df_use_2_copd = pd.DataFrame(Dataset_2, columns = ['Score', 'COPD_P3'])  
4 df_use_2_copd
```

Out[20]:

	Score	COPD_P3
0	0	1
1	0	0
2	0	1
3	0	0
4	0	1
...
365	2	0
366	2	0
367	2	1
368	2	1
369	2	0

370 rows × 2 columns

```
In [22]: ► 1 # create contingency table  
2 data_crosstab_2_copd = pd.crosstab(df_use_2_copd['Score'],  
3                                     df_use_2_copd['COPD_P3'],  
4                                     )  
5  
6 data_crosstab_2_copd
```

Out[22]:

COPD_P3	0	1
Score		
0	31	26
2	218	95

```
In [23]: ► 1 # Chi-Squared Test
2
3 res_2 = stat()
4 res_2.chisq(df=data_crosstab_2_copd)
5 # output
6 print(res_2.summary)
```

Chi-squared test for independence

Test	Df	Chi-square	P-value
Pearson	1	4.43385	0.0352329
Log-likelihood	1	4.25976	0.0390255

We will now observe the Expected Frequency Counts for the Chi-Squared Test for Independence executed on the Adjusted Clock Drawing Test Score by the Self-Reported COPD Score variable.

```
In [24]: ► 1 print(res_2.expected_df)
2
3 # Need at least 5, preferably 15
4 # Assumptions met.
```

Expected frequency counts

	0	1
0	38.3595	18.6405
1	210.641	102.359

Chi-Squared Test for Independence on the Adjusted Clock Drawing Test Score by the Final GOLD Baseline Variable

In [28]: ►

```
1 # Chi-Square Test  
2  
3 df_use_2_gold = pd.DataFrame(Dataset_2, columns = ['Score', 'goldstage_final_V2'])  
4 df_use_2_gold
```

Out[28]:

	Score	goldstage_final_V2
0	0	2
1	0	2
2	0	2
3	0	2
4	0	2
...
365	2	2
366	2	1
367	2	2
368	2	2
369	2	2

370 rows × 2 columns

In [30]: ►

```
1 # create contingency table  
2 data_crosstab_2_gold = pd.crosstab(df_use_2_gold['Score'],  
3                                     df_use_2_gold['goldstage_final_V2'],  
4                                     )  
5  
6 data_crosstab_2_gold
```

Out[30]:

	goldstage_final_V2		
Score	0	1	2
0	20	6	31
2	141	37	135

```
In [31]: ► 1 # Chi-Squared Test
2
3 res_2_gold = stat()
4 res_2_gold.chisq(df=data_crosstab_2_gold)
5 # output
6 print(res_2_gold.summary)
```

Chi-squared test for independence

Test	Df	Chi-square	P-value
Pearson	2	2.53034	0.282192
Log-likelihood	2	2.52584	0.282827

- 1 Finally, we will perform the K-Nearest Neighbors Algorithm to solve the classification problem of a patients Clock Drawing Test Score.

```
In [38]: ► 1 df_knn_2 = pd.DataFrame(Dataset_2, columns = ['fev1pp_final', 'goldstage']
2 df_knn_2
```

Out[38]:

	fev1pp_final	goldstage_final_V2	COPD_P3	Score
0	42.7	2	1	0
1	120.3	2	0	0
2	19.6	2	1	0
3	77.7	2	0	0
4	35.6	2	1	0
...
365	83.6	2	0	2
366	73.9	1	0	2
367	72.4	2	1	2
368	51.7	2	1	2
369	119.4	2	0	2

370 rows × 4 columns

```
In [45]: ► 1 X_2 = df_knn_2.drop("Score", axis = 1)
2 X_2 = X_2.values
3 y_2 = df_knn_2['Score']
4 y_2 = y_2.values
```

```
In [46]: ► train and test sets
1 X_2_train, X_2_test, y_2_train, y_2_test = train_test_split(X_2, y_2, test_size=0.2, random_state=42)
```

```
In [47]: ► 1 # KNN is a supervised machine Learning method  
2 # Trying to predict weather or not someone will have COPD based  
3 # on their CDT Score (the predictor)  
4 knn_model_2_score = KNeighborsRegressor(n_neighbors=2)  
5 knn_model_2_score.fit(X_train_2_score, y_train_2_score)
```

Out[47]:

KNeighborsRegressor
KNeighborsRegressor(n_neighbors=2)

```
In [48]: ► 1 # Calculate RMSE for Training Data  
2 train_preds_2_score = knn_model_2_score.predict(X_train_2_score)  
3 mse_2_score = mean_squared_error(y_train_2_score, train_preds_2_score)  
4 rmse_2_score = sqrt(mse_2_score)  
5 rmse_2_score
```

Out[48]: 0.48629876257963733

```
In [49]: ► 1 # Calculate RMSE for Testing Data  
2 test_preds_2_test_s = knn_model_2_score.predict(X_test_2_score)  
3 mse_2_test_s = mean_squared_error(y_test_2_score, test_preds_2_test_s)  
4 rmse_2_test_s = sqrt(mse_2_test_s)  
5 rmse_2_test_s
```

Out[49]: 0.86991767240168

```
In [50]: ► 1 # Calculate Accuracy  
2 y_pred_2_s = knn_model_2_score.predict(X_test_2_score)  
3  
4 # Evaluate the training & testing Accuracy  
5 #  
6 print('Training accuracy score: %.3f' % knn_model_2_score.score(X_train_2_score))  
7 print('Test accuracy score: %.3f' % knn_model_2_score.score(X_test_2_score))
```

Training accuracy score: 0.541
Test accuracy score: -0.392

1 Analysis of the Factor Variables on the Adjusted Clock Drawing Test Score

Now, we will pixelate the 565 Clock Drawing Test images

In [57]: ►

```
1 # pixelating PDF clock images and converting them to data matrix
2
3 # Reading in Images with an Adjusted Clock Drawing Test Score
4 # of 0
5
6 path = r'C:/Users/12563/Documents/Test Clocks 2/'
7 #path = r'C:/Users/12563/Documents/Clocks_0/'
8 files = os.listdir(path)
9 dmat1 = np.zeros((len(files), 288, 288))
10
11 for index, file in enumerate(files):
12     f_name_prim = os.path.basename(file)
13     print(f_name_prim)
14     print(path)
15     print(path+f_name_prim)
16
17 _, ext = os.path.splitext(file)
18 if ext.lower() == '.pdf':
19     # pdf processing
20     #print(file)
21     images = convert_from_path(path + f_name_prim, poppler_path = r"C:/Program Files/PDF2Image/Poppler/bin")
22     image = np.array(images[0])
23     y,x,_ = image.shape
24
25     # cropping the clocks from the images
26     #print(image.shape)
27     # saving the clocks in the target folder
28     #fig, ax = plt.subplots()
29     #plt.title('File name: {}, Index: {}'.format(file, index))
30     #ax.imshow(image, cmap='gray')
31     #Cropping
32     # I decided on the cropping bounds by looking at printed images
33     image = image[350:1250,350:1350]
34     #print(image.shape)
35     # saving the clocks in the target folder
36     fig, ax = plt.subplots()
37     plt.title('File name: {}, Index: {}'.format(file, index))
38     ax.imshow(image, cmap='gray')
39     print()
40
41     # Grayscaleing the images
42     image = np.mean(image, axis = 2)
43
44     # resizing the cropped clocks
45     image_resized = resize(image, (288, 288))
46
47     dmat1[index] = image_resized
48
49     # saving the clocks into the target folder
50     fig, ax = plt.subplots()
51     plt.title('File name: {}, Index: {}'.format(file, index))
52     ax.imshow(image_resized, cmap='gray')
53     #fig.savefig(path + f_name_prim + '.jpeg')
54
55     # saving the datamatrix in the target folder
56     #np.save(path + 'dmat2', dmat2)
```

```

57
58     elif ext.lower() == '.jpeg':
59         # jpeg processing
60
61         image = imread(path + f_name_prim)
62
63         # cropping the clocks from the images
64         image = image[550:2400, 550:2400]
65
66         # Grayscaleing the images
67         image = np.mean(image, axis = 2)
68
69         # resizing the cropped clocks
70         image_resized = resize(image, (288, 288))
71
72         dmat1[index] = image_resized
73
74
75         # saving the clocks in the target folder
76         fig, ax = plt.subplots()
77         plt.title('File name: {}, Index: {}'.format(file, index))
78         ax.imshow(image, cmap='gray')
79
80         #fig.savefig(path + f_name_prim)
81     else:
82         #print(file, ext)
83
84         if file == "JPEG":
85             #continue
86             #print(file)
87             #images = convert_from_path(path + f_name_prim, poppler_path = r'')
88             #image = np.array(images[0])
89             #y,x,_ = image.shape

```

10071D.PDF

C:/Users/12563/Documents/Test Clocks 2/
C:/Users/12563/Documents/Test Clocks 2/10071D.PDF

10098X.PDF

C:/Users/12563/Documents/Test Clocks 2/
C:/Users/12563/Documents/Test Clocks 2/10098X.PDF

10111P.PDF

C:/Users/12563/Documents/Test Clocks 2/
C:/Users/12563/Documents/Test Clocks 2/10111P.PDF

10286Y.jpeg

C:/Users/12563/Documents/Test Clocks 2/
C:/Users/12563/Documents/Test Clocks 2/10286Y.jpeg

10292T.PDF

C:/Users/12563/Documents/Test Clocks 2/
C:/Users/12563/Documents/Test Clocks 2/10292T.PDF

Now, we will rotate the images in need of rotation as well as delete images that need to be deleted.

```
In [62]: ► 1 # This code removes the noise and rotates the images
2
3 # Images that need rotation (Indexes via printing in the cell above)
4 S_rot = [31, 216, 445, 554]
5
6 # Noise images that will be removed, Clock not shown on images
7 S_noise = [4, 12, 14, 15, 27, 33, 35, 41, 44, 47, 51, 52, 93, 103, 125, 1
8 251, 259, 270, 271, 274, 278, 281, 285, 287, 288, 293, 296, 3
9 449, 450, 451, 467, 468, 526, 559, 561]
10
11 #S_rot = []
12 #S_noise= [4, 12, 13, 34, 46, 65]
13 # function for rotating the JPG images
14
15 def img_rotate (df, rotate) :
16     for i in rotate:
17         df[i] = np.rot90(df[i],3)
18     return(df)
19
20 dmat1 = img_rotate (dmat1, S_rot)
21
22
23 # removing noise images
24 dmat1_nonoise = np.delete(dmat1, S_noise , axis = 0)
```

```
In [63]: ► 1 # Loading X(clock matrix) and Y (labels) data (already suffeled and noise
2 X = np.load('/Users/12563/Desktop/Batch 1 (2)/Xall_nonois_shfl_use_0.npy'
3 y = np.load('/Users/12563/Desktop/Batch 1 (2)/Yall_nonois_shfl_use_0.npy')
```

```
In [64]: ► 1 # Creating array of the files to create
2 files = np.array(files)
3
4 files_nonoise = np.delete(files, S_noise, axis = 0)
```

Now, we are going to create a data frame that consists of the Subject ID's in which we have successfully pixeled their Clock Drawing Test Image.

In [65]: ►

```
1 # Generating the List of scores associated with the files_nonoise
2 score_list = np.zeros(files_nonoise.shape[0])
3
4 blank = []
5
6 # for Loop to match score of df_join to Pixelated clock images
7 for indices in range(files_nonoise.shape[0]):
8     ext, _ = os.path.splitext(files_nonoise[indices])
9     #print(ext)
10    if ext in Clock_list['Clocks'].tolist():
11        score_update = Clock_list[Clock_list['Clocks'] == ext]['Scores']
12        print(score_update)
13        score_list[indices] = score_update
14    else:
15        blank.append(indices)
16
17 # Cleaned Image Data
18 dmat1_delete = np.delete(dmat1_nonoise, blank, axis = 0)
19
20 # Cleaned Score Data
21 Scores_delete = np.delete(score_list, blank, axis = 0)
```

```
0
0
0
0
2
2
2
0
2
2
2
0
0
2
0
2
2
2
2
2
2
~
```

In [87]: ►

```
1 # Deleting the SIDs of the files that were deleted
2 files_delete = np.delete(files_nonoise, blank, axis = 0)
3
4 # Convert the NumPy Array to Pandas data Frame
5 files_delete = pd.DataFrame(files_delete)
6
7 # Rename column in a dataframe
8 files_delete.columns = ['Clocks']
```

```
In [88]: ► 1 # We need to get rid of the .PDF, .pdf, & .jpeg
  2 # Strip .pdf
  3 files_delete['Clocks'] = files_delete['Clocks'].str.rstrip('.pdf')
  4 # Strip .PDF
  5 files_delete['Clocks'] = files_delete['Clocks'].str.rstrip('.PDF')
  6 # Strip .jpeg
  7 files_delete['Clocks'] = files_delete['Clocks'].str.rstrip('.jpeg')
  8
  9 files_delete
```

Out[88]:

Clocks	
0	10071
1	10098X
2	10111
3	10286Y
4	10294X
...	...
480	30352X
481	30357H
482	30361Y
483	30363C
484	30368M

485 rows × 1 columns

```
In [86]: ► 1 files_delete
```

```
In [92]: ► 1 # Rename column in a dataframe
  2 df_factors.rename(columns = {'sid':'Clocks' }, inplace = True)
  3
```

```
In [93]: ┆ 1 df = files_delete.merge(df_factors, on = 'Clocks')  
2 df
```

Out[93]:

	Clocks	F1compP2	F2compP2	F3compP2	F4compP2
0	10098X	0.396002	0.868078	0.451179	-0.729966
1	10286Y	0.792556	-0.630638	0.904748	-0.785195
2	10294X	1.988744	0.705892	1.891007	0.835416
3	10348U	-0.819512	-0.414752	-0.379585	-2.070531
4	10528W	1.043932	-0.255540	0.788534	0.838686
...
325	24160C	-0.985130	-0.566078	-1.393219	-0.892061
326	24196X	-1.222794	-1.192649	-0.664929	-0.932192
327	24205Y	0.067990	-2.257025	-0.196353	-1.515259
328	24555Z	0.356886	-0.348367	0.187666	0.204181
329	25340J	-0.577903	-1.390396	-0.710814	0.085333

330 rows × 5 columns

```
In [94]: ┆ 1 dataset_factors_use = df.merge(Clock_list, on = 'Clocks')  
2 dataset_factors_use
```

Out[94]:

	Clocks	F1compP2	F2compP2	F3compP2	F4compP2	Scores
0	10098X	0.396002	0.868078	0.451179	-0.729966	0
1	10286Y	0.792556	-0.630638	0.904748	-0.785195	0
2	10294X	1.988744	0.705892	1.891007	0.835416	2
3	10348U	-0.819512	-0.414752	-0.379585	-2.070531	2
4	10528W	1.043932	-0.255540	0.788534	0.838686	2
...
328	24160C	-0.985130	-0.566078	-1.393219	-0.892061	2
329	24196X	-1.222794	-1.192649	-0.664929	-0.932192	2
330	24205Y	0.067990	-2.257025	-0.196353	-1.515259	2
331	24555Z	0.356886	-0.348367	0.187666	0.204181	2
332	25340J	-0.577903	-1.390396	-0.710814	0.085333	2

333 rows × 6 columns

```
In [97]: 1 df_box = df.drop('Clocks', axis = 1)
2 df_box
```

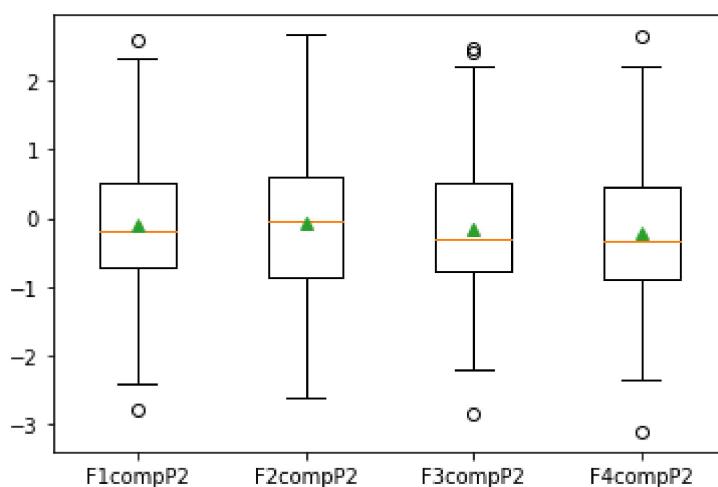
Out[97]:

	F1compP2	F2compP2	F3compP2	F4compP2
0	0.396002	0.868078	0.451179	-0.729966
1	0.792556	-0.630638	0.904748	-0.785195
2	1.988744	0.705892	1.891007	0.835416
3	-0.819512	-0.414752	-0.379585	-2.070531
4	1.043932	-0.255540	0.788534	0.838686
...
325	-0.985130	-0.566078	-1.393219	-0.892061
326	-1.222794	-1.192649	-0.664929	-0.932192
327	0.067990	-2.257025	-0.196353	-1.515259
328	0.356886	-0.348367	0.187666	0.204181
329	-0.577903	-1.390396	-0.710814	0.085333

330 rows × 4 columns

```
In [98]: 1 # Creating plot
2 plt.boxplot(df_box, showmeans = True)
3 plt.xticks([1, 2, 3, 4], ['F1compP2', 'F2compP2', 'F3compP2', 'F4compP2'])
```

Out[98]: ([<matplotlib.axis.XTick at 0x234daf6a370>,
<matplotlib.axis.XTick at 0x234daf6a400>,
<matplotlib.axis.XTick at 0x234cd9cef0>,
<matplotlib.axis.XTick at 0x234b8acec70>],
[Text(1, 0, 'F1compP2'),
Text(2, 0, 'F2compP2'),
Text(3, 0, 'F3compP2'),
Text(4, 0, 'F4compP2')])



```
In [99]: 1 # T-test Factor 1  
2 sc.stats.ttest_ind(dataset_factors_use['Scores'], dataset_factors_use['F
```

```
Out[99]: Ttest_indResult(statistic=17.992744029470185, pvalue=2.9525947380683684e-59)
```

```
In [100]: 1 # T-test Factor 2  
2 sc.stats.ttest_ind(dataset_factors_use['Scores'], dataset_factors_use['F
```

```
Out[100]: Ttest_indResult(statistic=16.599459858946588, pvalue=5.094635135350925e-52)
```

```
In [101]: 1 # T-test Factor 3  
2 sc.stats.ttest_ind(dataset_factors_use['Scores'], dataset_factors_use['F
```

```
Out[101]: Ttest_indResult(statistic=18.668544585336004, pvalue=7.720186926657116e-63)
```

```
In [102]: 1 # T-test Factor 4  
2 sc.stats.ttest_ind(dataset_factors_use['Scores'], dataset_factors_use['F
```

```
Out[102]: Ttest_indResult(statistic=18.664412707937306, pvalue=8.121916090616809e-63)
```

```
In [103]: 1 df_fac_adj = pd.read_csv('C:/Users/12563/Documents/Res. Methods in Math &  
2 df_fac_adj  
3
```

```
Out[103]:
```

	sid	age_final	gender_final	race_final	education_finalV1	income_final	kidney_dise
0	10015T	69.9		1	1	3	5
1	10017X	77.6		1	1	3	4
2	10031R	71.7		1	1	3	2
3	10032T	71.1		2	1	3	2
4	10049K	60.0		1	2	1	1
...
1712	23634Q	66.1		2	1	1	3
1713	23696M	68.7		2	1	3	2
1714	23785L	66.9		1	1	3	1
1715	23793K	55.8		2	1	3	5
1716	24637B	68.1		1	1	3	2

1717 rows × 8 columns

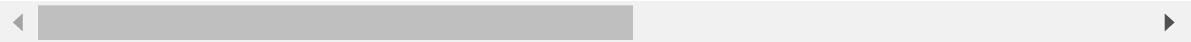
```
In [106]: 1 df_fac_adj = df_fac_adj.rename(columns = {'sid' : 'Clocks'})
```

```
In [109]: ► 1 df_fac_adj_ana = dataset_factors_use.merge(df_fac_adj, on = 'Clocks')
2 df_fac_adj_ana
```

Out[109]:

	Clocks	F1compP2	F2compP2	F3compP2	F4compP2	Scores	age_final	gender_final
0	10098X	0.396002	0.868078	0.451179	-0.729966	0	88.3	2
1	10286Y	0.792556	-0.630638	0.904748	-0.785195	0	59.6	2
2	10348U	-0.819512	-0.414752	-0.379585	-2.070531	2	84.1	1
3	10528W	1.043932	-0.255540	0.788534	0.838686	2	85.7	1
4	10647E	1.361600	0.373160	1.250930	0.534924	2	77.8	2
...
223	23467X	-0.456152	-1.794621	-0.599533	-0.607022	2	62.9	2
224	23478C	-0.905084	1.815225	-0.993261	-0.841087	2	59.8	2
225	23619U	0.723385	0.281395	0.804737	-0.071033	2	62.5	2
226	23634Q	-0.533470	-0.131728	-0.553965	0.532302	2	66.1	2
227	23785L	0.771995	1.026923	1.211564	1.812055	2	66.9	1

228 rows × 13 columns



```
In [112]: ► 1 B1 = pd.read_csv('C:/Users/12563/Documents/Res. Methods in Math & Stats/I
2 B1 = B1.rename(columns = {'sid' : 'Clocks'})
```

```
In [113]: ► 1 df_b = B1.merge(df_fac_adj_ana, on = 'Clocks')
2 df_b
```

Out[113]:

compP2	F3compP2	F4compP2	Scores	age_final	gender_final	race_final	education_finalV1	ir
0.868078	0.451179	-0.729966	0	88.3	2	1		3
-0.630638	0.904748	-0.785195	0	59.6	2	2		3
-0.414752	-0.379585	-2.070531	2	84.1	1	1		3
0.255540	0.788534	0.838686	2	85.7	1	1		3
0.373160	1.250930	0.534924	2	77.8	2	1		3
...
-0.794621	-0.599533	-0.607022	2	62.9	2	1		3
0.815225	-0.993261	-0.841087	2	59.8	2	1		3
0.281395	0.804737	-0.071033	2	62.5	2	1		3
-0.131728	-0.553965	0.532302	2	66.1	2	1		1
0.026923	1.211564	1.812055	2	66.9	1	1		3



In [115]:

```

1 # F1 ANOVA
2 # w/o GOLD
3 F1_anova_model = ols('F1compP2 ~ C(Scores) + age_final + C(gender_final)
4 anova_table_F1 = sm.stats.anova_lm(F1_anova_model, typ=1)
5 anova_table_F1

```

Out[115]:

	df	sum_sq	mean_sq	F	PR(>F)
C(Scores)	1.0	7.081691	7.081691	10.137368	0.001668
C(gender_final)	1.0	2.616250	2.616250	3.745135	0.054273
C(race_final)	1.0	1.154399	1.154399	1.652510	0.200001
C(education_finalV1)	2.0	3.000253	1.500126	2.147416	0.119284
C(income_final)	5.0	3.659889	0.731978	1.047819	0.390533
C(kidney_disease_final)	1.0	3.185581	3.185581	4.560126	0.033856
age_final	1.0	12.580973	12.580973	18.009533	0.000033
Residual	215.0	150.193184	0.698573	NaN	NaN

In [118]:

```

1 # F2 ANOVA
2 F2_anova_model = ols('F2compP2 ~ C(Scores) + age_final + C(gender_final)
3 anova_table_F2 = sm.stats.anova_lm(F2_anova_model, typ=1)
4 anova_table_F2

```

Out[118]:

	df	sum_sq	mean_sq	F	PR(>F)
C(Scores)	1.0	1.254987	1.254987	1.488475	0.223790
C(gender_final)	1.0	1.714421	1.714421	2.033386	0.155328
C(race_final)	1.0	0.337690	0.337690	0.400517	0.527495
C(education_finalV1)	2.0	2.649832	1.324916	1.571414	0.210128
C(income_final)	5.0	8.449333	1.689867	2.004263	0.079221
C(kidney_disease_final)	1.0	0.012088	0.012088	0.014338	0.904801
age_final	1.0	1.270923	1.270923	1.507375	0.220883
Residual	215.0	181.274321	0.843136	NaN	NaN

In [120]: ►

```

1 # F3 ANOVA
2 # Gas-Trapping
3 F3_anova_model = ols('F3compP2 ~ C(Scores) + age_final + C(gender_final)
4 anova_table_F3 = sm.stats.anova_lm(F3_anova_model, typ=1)
5 anova_table_F3

```

Out[120]:

	df	sum_sq	mean_sq	F	PR(>F)
C(Scores)	1.0	7.432433	7.432433	10.702248	0.001246
C(gender_final)	1.0	0.700826	0.700826	1.009146	0.316237
C(race_final)	1.0	0.000024	0.000024	0.000034	0.995331
C(education_finalV1)	2.0	6.052799	3.026399	4.357830	0.013956
C(income_final)	5.0	4.369009	0.873802	1.258221	0.283208
C(kidney_disease_final)	1.0	3.222713	3.222713	4.640509	0.032338
age_final	1.0	16.417315	16.417315	23.639926	0.000002
Residual	215.0	149.311918	0.694474	NaN	NaN

In [122]: ►

```

1 # F4 ANOVA
2 # TLC FRC
3 F4_anova_model = ols('F4compP2 ~ C(Scores) + age_final + C(gender_final)
4 anova_table_F4 = sm.stats.anova_lm(F4_anova_model, typ=1)
5 anova_table_F4

```

Out[122]:

	df	sum_sq	mean_sq	F	PR(>F)
C(Scores)	1.0	3.506831	3.506831	4.191373	0.041846
C(gender_final)	1.0	0.002536	0.002536	0.003031	0.956145
C(race_final)	1.0	0.264772	0.264772	0.316456	0.574331
C(education_finalV1)	2.0	14.658253	7.329127	8.759791	0.000220
C(income_final)	5.0	6.207951	1.241590	1.483952	0.196206
C(kidney_disease_final)	1.0	4.268622	4.268622	5.101867	0.024903
age_final	1.0	6.710542	6.710542	8.020457	0.005065
Residual	215.0	179.885834	0.836678	NaN	NaN

In [123]: ►

```
1 # pixelating PDF clock images and converting them to data matrix
2
3 # Reading in Images with an Adjusted Clock Drawing Test Score
4 # of 0
5
6 #path = r'C:/Users/12563/Documents/Test_Clocks_2/'
7 path = r'C:/Users/12563/Documents/Clocks_0/'
8 files = os.listdir(path)
9 dmat1 = np.zeros((len(files), 288, 288))
10
11 for index, file in enumerate(files):
12     f_name_prim = os.path.basename(file)
13     print(f_name_prim)
14     print(path)
15     print(path+f_name_prim)
16
17 _, ext = os.path.splitext(file)
18 if ext.lower() == '.pdf':
19     # pdf processing
20     #print(file)
21     images = convert_from_path(path + f_name_prim, poppler_path = r'C:/Program Files/Poppler/bin')
22     image = np.array(images[0])
23     y,x,_ = image.shape
24
25     # cropping the clocks from the images
26     #print(image.shape)
27     # saving the clocks in the target folder
28     #fig, ax = plt.subplots()
29     #plt.title('File name: {}, Index: {}'.format(file, index))
30     #ax.imshow(image, cmap='gray')
31     #Cropping
32     # I decided on the cropping bounds by looking at printed images
33     image = image[350:1250,350:1350]
34     #print(image.shape)
35     # saving the clocks in the target folder
36     fig, ax = plt.subplots()
37     plt.title('File name: {}, Index: {}'.format(file, index))
38     ax.imshow(image, cmap='gray')
39     print()
39
40     # Grayscaleing the images
41     image = np.mean(image, axis = 2)
42
43     # resizing the cropped clocks
44     image_resized = resize(image, (288, 288))
45
46     dmat1[index] = image_resized
47
48     # saving the clocks into the target folder
49     fig, ax = plt.subplots()
50     plt.title('File name: {}, Index: {}'.format(file, index))
51     ax.imshow(image_resized, cmap='gray')
52     #fig.savefig(path + f_name_prim + '.jpeg')
53
54     # saving the datamatrix in the target folder
55     #np.save(path + 'dmat2', dmat2)
```

```
57
58 elif ext.lower() == '.jpeg':
59     # jpeg processing
60
61     image = imread(path + f_name_prim)
62
63     # cropping the clocks from the images
64     image = image[550:2400, 550:2400]
65
66     # Grayscaleing the images
67     image = np.mean(image, axis = 2)
68
69     # resizing the cropped clocks
70     image_resized = resize(image, (288, 288))
71
72     dmat1[index] = image_resized
73
74     # saving the clocks in the target folder
75     fig, ax = plt.subplots()
76     plt.title('File name: {}, Index: {}'.format(file, index))
77     ax.imshow(image, cmap='gray')
78
79     #fig.savefig(path + f_name_prim)
80 else:
81     #print(file, ext)
82
83     #if file == "JPEG":
84         #continue
85     #print(file)
86     #images = convert_from_path(path + f_name_prim, poppler_path = r
87     #image = np.array(images[0])
88     #y,x,_ = image.shape
```

In [124]: ►

```

1 # This code removes the noise and rotates the images
2
3 # Images that need rotation (Indexes via printing in the cell above)
4 #S_rot = [31, 216, 445, 554]
5
6 # Noise images that will be removed, Clock not shown on images
7 #S_noise = [4, 12, 14, 15, 27, 33, 35, 41, 44, 47, 51, 52, 93, 103, 125, :
8 #      251, 259, 270, 271, 274, 278, 281, 285, 287, 288, 293, 296, :
9 #      449, 450, 451, 467, 468, 526, 559, 561]
10
11 S_rot = []
12 S_noise= [4, 12, 13, 34, 46, 65]
13 # function for rotating the JPG images
14
15 def img_rotate (df, rotate) :
16     for i in rotate:
17         df[i] = np.rot90(df[i],3)
18     return(df)
19
20 dmat1 = img_rotate (dmat1, S_rot)
21
22
23 # removing noise images
24 dmat1_nonoise = np.delete(dmat1, S_noise , axis = 0)

```

In [126]: ►

```

1 # Loading X(clock matrix) and Y (labels) data (already suffeled and noise
2 X = np.load('/Users/12563/Desktop/Batch 1 (2)/Xall_nonois_hfl_use_0.npy'
3 y = np.load('/Users/12563/Desktop/Batch 1 (2)/Yall_nonois_hfl_use_0.npy'

```

In [140]: ►

```

1 # Add an extra Dimension for Channel for CNN
2 X = X.reshape([len(X), 288, 288, 1])
3
4 # Split the data into train and test sets
5 X_train, X_test, y_train, y_test = train_test_split(X, y, test_size= 0.2)
6
7 print (len(y_train[y_train == 0]) , len(y_train[y_train == 1]))
8 print (len(y_test[y_test == 0]) , len(y_test[y_test == 1]))
9
10 # Scaling train and test datasets
11 X_train = X_train/255
12 X_test = X_test/255

```

52 0
13 0

In [141]:

```
1 # Function for individual CNN models (2 structures: Conv_Pool and Conv_Pool_
2 def model(activation = 'relu'):
3     model = Sequential()
4     # Input Image with dimensions (288x288x1) an input image
5     # 288 pixels wide by 288 pixels in height with 1 channel
6     model.add(Conv2D(32, kernel_size=(3,3), activation= activation , input_shape=(288,288,1)))
7     model.add(MaxPooling2D(pool_size=(2,2),name="maxpool_1"))
8     #     model.add(Dropout(0.25))
9     #     model.add(Conv2D(64, kernel_size=(3, 3), activation= activation , name="conv_2"))
10    #     model.add(MaxPooling2D(pool_size=(2,2),name="maxpool_2"))
11    #     model.add(Dropout(0.25))
12    model.add(Flatten())
13    model.add(Dense(32,activation="relu", name="dense_layer_1"))
14    model.add(Dense(2,activation="sigmoid"))
15    model.compile(loss="sparse_categorical_crossentropy", metrics=['accuracy'])
16    return(model)
17
18 conv_model = model(activation = 'relu')
19 conv_hist = conv_model.fit(X_train,y_train ,epochs= 20,validation_data=(X_val,y_val))
20 conv_hist_dict = conv_hist.history
```

```
Epoch 1/20
2/2 [=====] - 2s 591ms/step - loss: 0.4427 - accuracy: 0.6346 - val_loss: 0.0030 - val_accuracy: 1.0000
Epoch 2/20
2/2 [=====] - 1s 519ms/step - loss: 7.7510e-04 - accuracy: 1.0000 - val_loss: 4.4656e-06 - val_accuracy: 1.0000
Epoch 3/20
2/2 [=====] - 1s 449ms/step - loss: 1.1508e-06 - accuracy: 1.0000 - val_loss: 0.0000e+00 - val_accuracy: 1.0000
Epoch 4/20
2/2 [=====] - 1s 446ms/step - loss: 0.0000e+00 - accuracy: 1.0000 - val_loss: 0.0000e+00 - val_accuracy: 1.0000
Epoch 5/20
2/2 [=====] - 1s 511ms/step - loss: 0.0000e+00 - accuracy: 1.0000 - val_loss: 0.0000e+00 - val_accuracy: 1.0000
Epoch 6/20
2/2 [=====] - 1s 526ms/step - loss: 0.0000e+00 - accuracy: 1.0000 - val_loss: 0.0000e+00 - val_accuracy: 1.0000
Epoch 7/20
2/2 [=====] - 1s 576ms/step - loss: 0.0000e+00 - accuracy: 1.0000 - val_loss: 0.0000e+00 - val_accuracy: 1.0000
Epoch 8/20
2/2 [=====] - 1s 548ms/step - loss: 0.0000e+00 - accuracy: 1.0000 - val_loss: 0.0000e+00 - val_accuracy: 1.0000
Epoch 9/20
2/2 [=====] - 1s 560ms/step - loss: 0.0000e+00 - accuracy: 1.0000 - val_loss: 0.0000e+00 - val_accuracy: 1.0000
Epoch 10/20
2/2 [=====] - 1s 513ms/step - loss: 0.0000e+00 - accuracy: 1.0000 - val_loss: 0.0000e+00 - val_accuracy: 1.0000
Epoch 11/20
2/2 [=====] - 1s 483ms/step - loss: 0.0000e+00 - accuracy: 1.0000 - val_loss: 0.0000e+00 - val_accuracy: 1.0000
Epoch 12/20
2/2 [=====] - 1s 482ms/step - loss: 0.0000e+00 - accuracy: 1.0000 - val_loss: 0.0000e+00 - val_accuracy: 1.0000
```

```
accuracy: 1.0000 - val_loss: 0.0000e+00 - val_accuracy: 1.0000
Epoch 13/20
2/2 [=====] - 1s 505ms/step - loss: 0.0000e+00 - accuracy: 1.0000 - val_loss: 0.0000e+00 - val_accuracy: 1.0000
Epoch 14/20
2/2 [=====] - 1s 537ms/step - loss: 0.0000e+00 - accuracy: 1.0000 - val_loss: 0.0000e+00 - val_accuracy: 1.0000
Epoch 15/20
2/2 [=====] - 1s 479ms/step - loss: 0.0000e+00 - accuracy: 1.0000 - val_loss: 0.0000e+00 - val_accuracy: 1.0000
Epoch 16/20
2/2 [=====] - 1s 470ms/step - loss: 0.0000e+00 - accuracy: 1.0000 - val_loss: 0.0000e+00 - val_accuracy: 1.0000
Epoch 17/20
2/2 [=====] - 1s 550ms/step - loss: 0.0000e+00 - accuracy: 1.0000 - val_loss: 0.0000e+00 - val_accuracy: 1.0000
Epoch 18/20
2/2 [=====] - 1s 480ms/step - loss: 0.0000e+00 - accuracy: 1.0000 - val_loss: 0.0000e+00 - val_accuracy: 1.0000
Epoch 19/20
2/2 [=====] - 1s 464ms/step - loss: 0.0000e+00 - accuracy: 1.0000 - val_loss: 0.0000e+00 - val_accuracy: 1.0000
Epoch 20/20
2/2 [=====] - 1s 478ms/step - loss: 0.0000e+00 - accuracy: 1.0000 - val_loss: 0.0000e+00 - val_accuracy: 1.0000
```

In [144]: ►

```
1 import tensorflow as tf
2 feature_extractor = tf.keras.Model(
3     inputs=conv_model.inputs,
4     outputs=[
5         conv_model.output, # < Last Layer output
6         conv_model.layers[2].output # < your convolution Layer output
7     ]
8 )
9
10 model_output, conv_layer_output = feature_extractor(X)
11 conv_layer_output
```

Out[144]: <tf.Tensor: shape=(65, 654368), dtype=float32, numpy=
array([[0. , 113.47755 , 40.749527, ... , 0. , 186.66112 ,
 0.],
 [0. , 34.348255, 18.957294, ... , 0. , 158.27054 ,
 0.],
 [0. , 113.02469 , 40.945564, ... , 0. , 185.92625 ,
 0.],
 ... ,
 [0. , 113.92434 , 40.90993 , ... , 0. , 187.39598 ,
 0.],
 [0. , 113.92434 , 40.90993 , ... , 0. , 187.39598 ,
 0.],
 [0. , 110.74818 , 39.708916, ... , 0. , 183.49847 ,
 0.]], dtype=float32)>

In [145]: ►

```
1 # Perform K-Means on Conv-Layer Output
2 model = KMeans(n_clusters=2, random_state=42)
3 kmeans = model.fit(conv_layer_output)
```

C:\Users\12563\OneDrive\New folder\lib\site-packages\sklearn\cluster_kmeans.py:1334: UserWarning: KMeans is known to have a memory leak on Windows with MKL, when there are less chunks than available threads. You can avoid it by setting the environment variable OMP_NUM_THREADS=1.

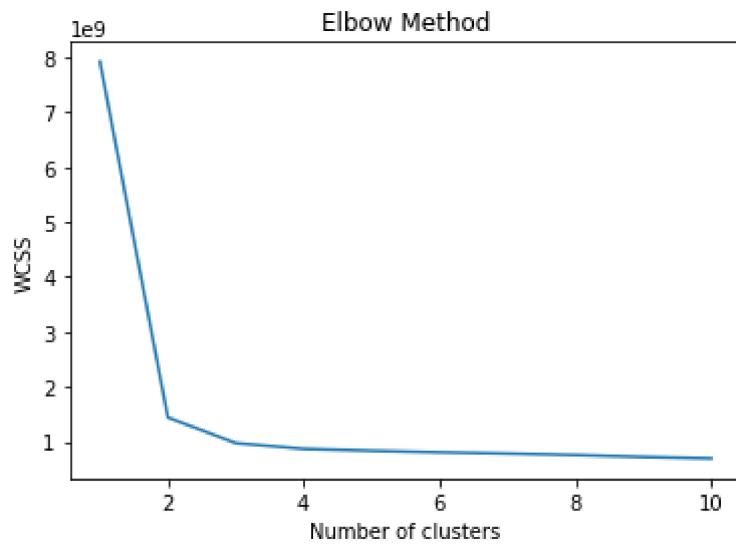
```
warnings.warn(
```

In [134]:

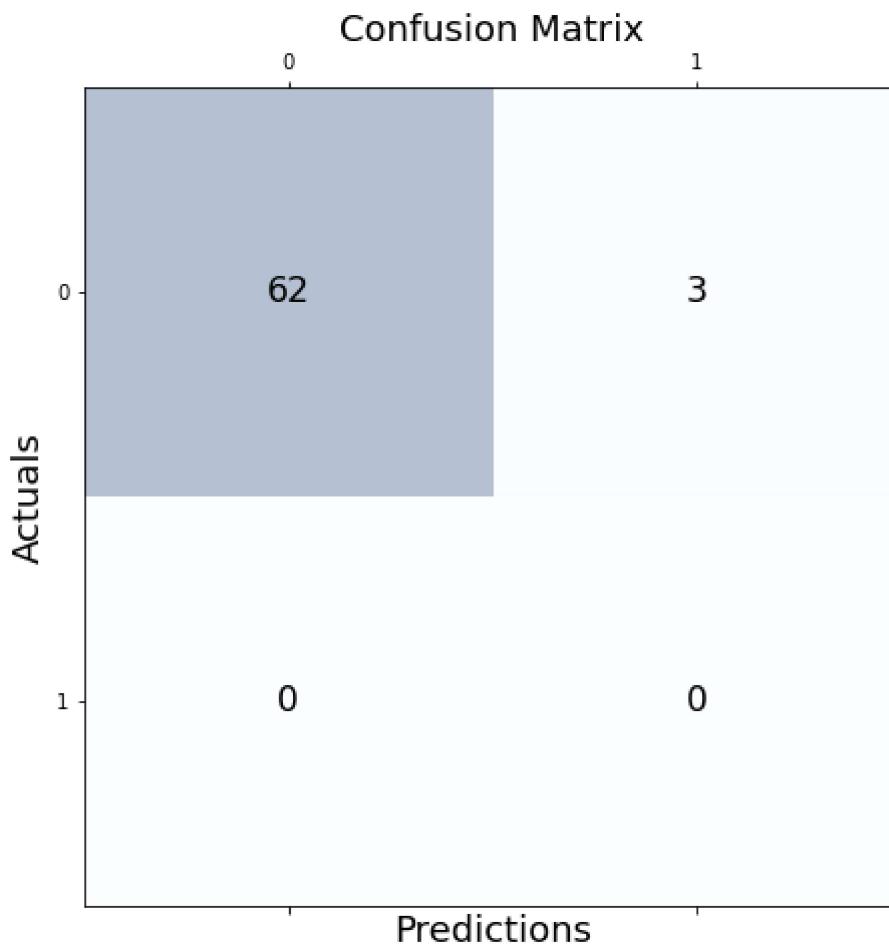
```
1 wcss = []
2 for i in range(1, 11):
3     kmeans = KMeans(n_clusters=i, init='k-means++', max_iter=300, n_init=1)
4     kmeans.fit(conv_layer_output)
5     wcss.append(kmeans.inertia_)
6 plt.plot(range(1, 11), wcss)
7 plt.title('Elbow Method')
8 plt.xlabel('Number of clusters')
9 plt.ylabel('WCSS')
10 plt.show()
```

```
C:\Users\12563\OneDrive\New folder\lib\site-packages\sklearn\cluster\_kmeans.py:1334: UserWarning: KMeans is known to have a memory leak on Windows with MKL, when there are less chunks than available threads. You can avoid it by setting the environment variable OMP_NUM_THREADS=1.
    warnings.warn(
C:\Users\12563\OneDrive\New folder\lib\site-packages\sklearn\cluster\_kmeans.py:1334: UserWarning: KMeans is known to have a memory leak on Windows with MKL, when there are less chunks than available threads. You can avoid it by setting the environment variable OMP_NUM_THREADS=1.
    warnings.warn(
C:\Users\12563\OneDrive\New folder\lib\site-packages\sklearn\cluster\_kmeans.py:1334: UserWarning: KMeans is known to have a memory leak on Windows with MKL, when there are less chunks than available threads. You can avoid it by setting the environment variable OMP_NUM_THREADS=1.
    warnings.warn(
C:\Users\12563\OneDrive\New folder\lib\site-packages\sklearn\cluster\_kmeans.py:1334: UserWarning: KMeans is known to have a memory leak on Windows with MKL, when there are less chunks than available threads. You can avoid it by setting the environment variable OMP_NUM_THREADS=1.
    warnings.warn(
C:\Users\12563\OneDrive\New folder\lib\site-packages\sklearn\cluster\_kmeans.py:1334: UserWarning: KMeans is known to have a memory leak on Windows with MKL, when there are less chunks than available threads. You can avoid it by setting the environment variable OMP_NUM_THREADS=1.
    warnings.warn(
C:\Users\12563\OneDrive\New folder\lib\site-packages\sklearn\cluster\_kmeans.py:1334: UserWarning: KMeans is known to have a memory leak on Windows with MKL, when there are less chunks than available threads. You can avoid it by setting the environment variable OMP_NUM_THREADS=1.
    warnings.warn(
C:\Users\12563\OneDrive\New folder\lib\site-packages\sklearn\cluster\_kmeans.py:1334: UserWarning: KMeans is known to have a memory leak on Windows with MKL, when there are less chunks than available threads. You can avoid it by setting the environment variable OMP_NUM_THREADS=1.
    warnings.warn(
C:\Users\12563\OneDrive\New folder\lib\site-packages\sklearn\cluster\_kmeans.py:1334: UserWarning: KMeans is known to have a memory leak on Windows with MKL, when there are less chunks than available threads. You can avoid it by setting the environment variable OMP_NUM_THREADS=1.
    warnings.warn(
C:\Users\12563\OneDrive\New folder\lib\site-packages\sklearn\cluster\_kmeans.py:1334: UserWarning: KMeans is known to have a memory leak on Windows with MKL, when there are less chunks than available threads. You can avoid it by setting the environment variable OMP_NUM_THREADS=1.
```

```
C:\Users\12563\OneDrive\New folder\lib\site-packages\sklearn\cluster\_kmeans.py:1334: UserWarning: KMeans is known to have a memory leak on Windows with MKL, when there are less chunks than available threads. You can avoid it by setting the environment variable OMP_NUM_THREADS=1.  
warnings.warn(
```



```
In [146]: ┏━
1 from sklearn.metrics import confusion_matrix
2 conf_matrix=confusion_matrix(y, kmeans.labels_)
3
4 fig, ax = plt.subplots(figsize=(7.5, 7.5))
5 ax.matshow(conf_matrix, cmap=plt.cm.Blues, alpha=0.3)
6 for i in range(conf_matrix.shape[0]):
7     for j in range(conf_matrix.shape[1]):
8         ax.text(x=j, y=i,s=conf_matrix[i, j], va='center',
9                  ha='center', size='xx-large')
10
11 plt.xlabel('Predictions', fontsize=18)
12 plt.ylabel('Actuals', fontsize=18)
13 plt.title('Confusion Matrix', fontsize=18)
14 plt.show()
```



```
In [135]: ┏━
1 from more_itertools import locate
2 list_to_check = kmeans.labels_
3 def find_indices(list_to_check, item_to_find):
4     indices = locate(list_to_check,lambda x: x == item_to_find)
5     return list(indices)
6
7 indices_of_intrest = find_indices(list_to_check, 1)
8 indices_of_intrest
```

Out[135]: [12, 23, 36, 50, 57]

```
In [136]: ┏━ 1 files_nonoise  
         2 files_nonoise[indices_of_intrest]
```

```
Out[136]: array(['10973R.PDF', '11715A.PDF', '12719N.PDF', '13529N.PDF',  
                 '13840N.pdf'], dtype='<U46')
```

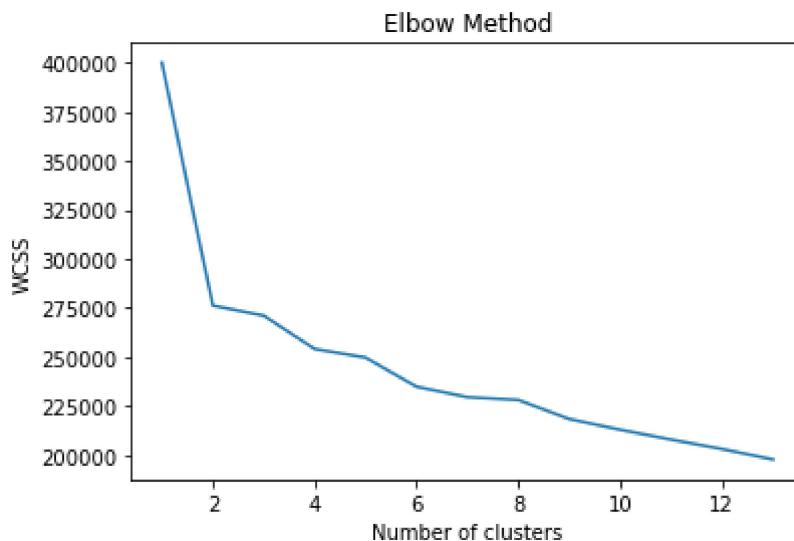
```
In [137]: ┏━ 1 # The Lowest SSE Value  
         2 kmeans.inertia_
```

```
Out[137]: 694187557.1665838
```

K-Medoids

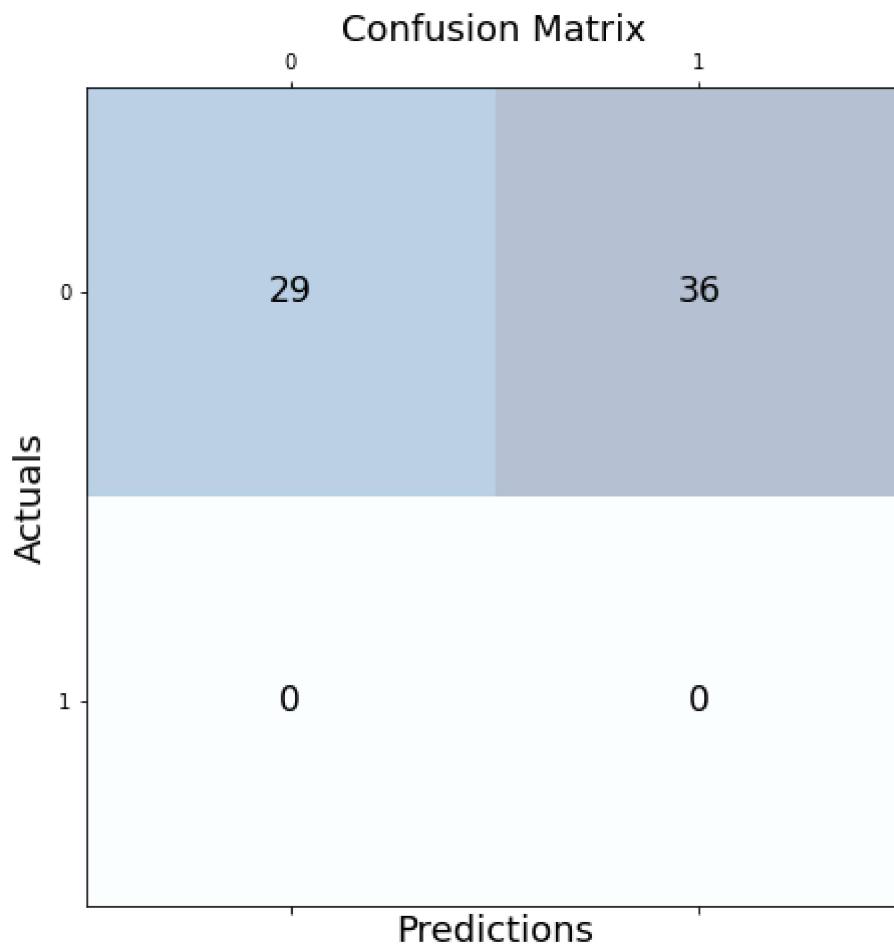
```
In [148]: ┏━ 1 # 2-Clusters  
         2  
         3 kmedoids = KMedoids(n_clusters=2, random_state=0).fit(conv_layer_output)
```

```
In [149]: ┏━ 1 wcss = []  
         2 for i in range(1, 14):  
         3     kmed = KMedoids(n_clusters=i, init='k-medoids++', max_iter=300, random_state=0)  
         4     kmed.fit(conv_layer_output)  
         5     wcss.append(kmed.inertia_)  
         6 plt.plot(range(1, 14), wcss)  
         7 plt.title('Elbow Method')  
         8 plt.xlabel('Number of clusters')  
         9 plt.ylabel('WCSS')  
        10 plt.show()
```



In [150]:

```
1 from sklearn.metrics import confusion_matrix
2 conf_matrix=confusion_matrix(y, kmedoids.labels_)
3
4 fig, ax = plt.subplots(figsize=(7.5, 7.5))
5 ax.matshow(conf_matrix, cmap=plt.cm.Blues, alpha=0.3)
6 for i in range(conf_matrix.shape[0]):
7     for j in range(conf_matrix.shape[1]):
8         ax.text(x=j, y=i,s=conf_matrix[i, j], va='center',
9                  ha='center', size='xx-large')
10
11 plt.xlabel('Predictions', fontsize=18)
12 plt.ylabel('Actuals', fontsize=18)
13 plt.title('Confusion Matrix', fontsize=18)
14 plt.show()
```



```
In [151]: ► 1 # The Lowest SSE Value  
2 # A good model is one with low inertia AND a small number  
3 # of clusters  
4  
5 kmmedoids.inertia_
```

Out[151]: 393572.34

```
In [158]: ► 1 df_ = pd.read_csv("C:/Users/12563/Documents/Res. Methods in Math & Stats/  
2 df_
```

Out[158]:

	sid	COPD_P3
0	10015T	1
1	10017X	1
2	10031R	0
3	10032T	1
4	10049K	1
...
1712	23634Q	0
1713	23696M	1
1714	23785L	1
1715	23793K	0
1716	24637B	0

1717 rows × 2 columns

In [159]: ►

```
1 from more_itertools import locate
2 list_to_check = kmmedoids.labels_
3 def find_indices(list_to_check, item_to_find):
4     indices = locate(list_to_check, lambda x: x == item_to_find)
5     return list(indices)
6
7 indices_of_intrest = find_indices(list_to_check, 0)
8 indices_of_intrest
```

Out[159]: [0,

3,

7,

11,

13,

16,

17,

18,

21,

22,

25,

27,

28,

29,

31,

33,

36,

38,

44,

45]

```
In [163]: ┏━
1 df_p = pd.DataFrame(files_nonoise[indices_of_intrest], columns = ['sid'])
2
3 # We need to get rid of the .PDF, .pdf, & .jpeg
4 # Strip .pdf
5 # Remove Leading and Treailing spaces
6 df_p['sid'] = df_p['sid'].str.strip('.pdf')
7 # Remove Leading and Treailing spaces
8 df_p['sid'] = df_p['sid'].str.strip('.PDF')
9 # Remove Leading and Treailing spaces
10 df_p['sid'] = df_p['sid'].str.strip('.jpeg')
11 df_p
12 #email_
```

Out[163]:

	sid
0	10071
1	10286Y
2	10482Y
3	10796V
4	11001N
5	11239U
6	11240
7	11246R
8	11292Y
9	11308N
10	11724B
11	11750C
12	11944
13	12258
14	12330L
15	12426Y
16	12719N
17	12857Z
18	13215Q
19	13334Y
20	13484R
21	13529N
22	13610W
23	13666X
24	13765Z
25	13848D
26	13937C

	sid
27	13954C
28	14006M

```
In [165]: ► 1 # Obtaining the Clock Drawing images obtaining a label of 1
  2 indices_of_intrest_1 = find_indices(list_to_check, 1)
  3 indices_of_intrest_1
```

```
Out[165]: [1,
 2,
 4,
 5,
 6,
 8,
 9,
 10,
 12,
 14,
 15,
 19,
 20,
 23,
 24,
 26,
 30,
 32,
 34,
 35,
 37,
 39,
 40,
 41,
 42,
 43,
 45,
 47,
 48,
 51,
 53,
 56,
 57,
 59,
 62,
 63]
```

```
In [166]: ┏━ df_p_1 = pd.DataFrame(files_nonoise[indices_of_intrest_1], columns = ['sid'])
  2
  3 # We need to get rid of the .PDF, .pdf, & .jpeg
  4 # Strip .pdf
  5 # Remove Leading and Treailing spaces
  6 df_p_1['sid'] = df_p_1['sid'].str.strip('.pdf')
  7 # Remove Leading and Treailing spaces
  8 df_p_1['sid'] = df_p_1['sid'].str.strip('.PDF')
  9 # Remove Leading and Treailing spaces
10 df_p_1['sid'] = df_p_1['sid'].str.strip('.jpeg')
11 df_p_1
12 #email_
```

Out[166]:

	sid
0	10098X
1	10111
2	10294X
3	10348U
4	10473X
5	10528W
6	10630N
7	10647E
8	10973R
9	11161J
10	11185X
11	11251K
12	11252M
13	11715A
14	11716C
15	11734E
16	12281Y
17	12343U
18	12581K
19	12635H
20	12721A
21	12862S
22	12885E
23	13083Z
24	13105J
25	13118S

	sid
26	13328
27	13378S
28	13384N
29	13550E
30	13629R
31	13824
32	13840N
33	13853W
34	13997U
35	14005K

In [168]: ► 1 fac = pd.read_csv("C:/Users/12563/Documents/Res. Methods in Math & Stats,

In [169]: ► 1 # Obtaining Average Value of Factor Scores in Cluster A
2 df_fac = fac.merge(df_p, on = 'sid')
3 df_fac
4
5 df_fac.mean()

C:\Users\12563\AppData\Local\Temp\ipykernel_19464/2776651171.py:5: FutureWarning: Dropping of nuisance columns in DataFrame reductions (with 'numeric_only=None') is deprecated; in a future version this will raise TypeError.
Select only valid columns before calling the reduction.
df_fac.mean()

Out[169]: F1compP2 0.234492
F2compP2 0.096026
F3compP2 0.331235
F4compP2 0.365538
dtype: float64

In [171]: ► 1 # Obtaining Average Value of Factor Scores in Cluster B
2 df_fac_1 = fac.merge(df_p_1, on = 'sid')
3 df_fac_1
4
5 df_fac_1.mean()

C:\Users\12563\AppData\Local\Temp\ipykernel_19464/4284279653.py:5: FutureWarning: Dropping of nuisance columns in DataFrame reductions (with 'numeric_only=None') is deprecated; in a future version this will raise TypeError.
Select only valid columns before calling the reduction.
df_fac_1.mean()

Out[171]: F1compP2 -0.210619
F2compP2 0.417673
F3compP2 -0.250752
F4compP2 -0.190061
dtype: float64

```
In [172]: ► 1 gold = pd.read_csv('C:/Users/12563/Documents/Res. Methods in Math & Stat:
```

```
In [175]: ► 1 # Obtaining Percentage of people with FINAL Golde Baseline 1-4 Cluster A
  2 df_gold = gold.merge(df_p, on = 'sid')
  3 df_gold
  4
  5 (np.sum(df_gold['goldstage_final_V2'] != 0)/29) *100
```

Out[175]: 41.37931034482759

```
In [181]: ► 1 # Obtaining Percentage of people with Self-Reported COPD Cluster A
  2 df_copd = gold.merge(df_p, on = 'sid')
  3 df_copd
  4
  5 (np.sum(df_copd['COPD_P3'] == 1)/29) *100
```

Out[181]: 24.137931034482758

```
In [174]: ► 1 # Obtaining Percentage of people with FINAL Golde Baseline 1-4 Cluster B
  2 df_gold_1 = gold.merge(df_p_1, on = 'sid')
  3 df_gold_1
  4
  5 (np.sum(df_gold_1['goldstage_final_V2'] != 0)/36) *100
```

Out[174]: 36.11111111111111

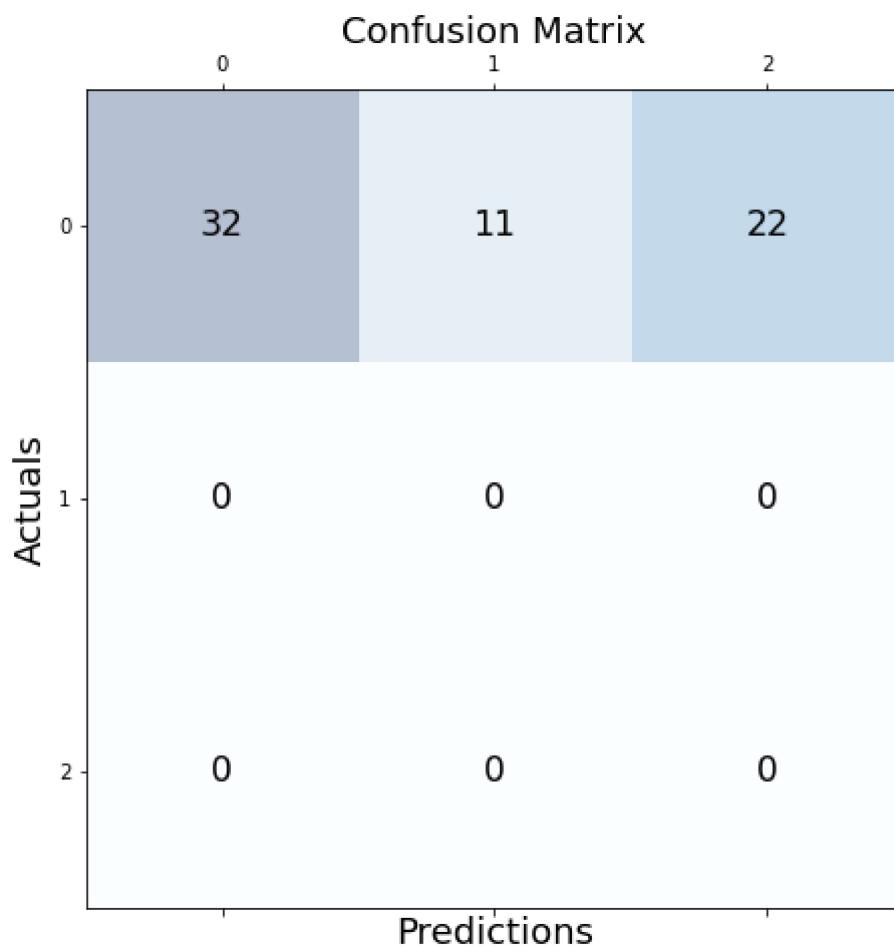
```
In [180]: ► 1 # Obtaining Percentage of people with Self-Reported COPD Cluster A
  2 df_copd_1 = gold.merge(df_p_1, on = 'sid')
  3 df_copd_1
  4
  5 (np.sum(df_copd_1['COPD_P3'] == 1)/29) *100
```

Out[180]: 10.344827586206897

```
In [152]: ► 1 # 3 Clusters
  2 kmedoids_3 = KMedoids(n_clusters=3, random_state=0).fit(conv_layer_outpu
```

In [153]:

```
1 from sklearn.metrics import confusion_matrix
2 conf_matrix=confusion_matrix(y, kmedoids_3.labels_)
3
4 fig, ax = plt.subplots(figsize=(7.5, 7.5))
5 ax.matshow(conf_matrix, cmap=plt.cm.Blues, alpha=0.3)
6 for i in range(conf_matrix.shape[0]):
7     for j in range(conf_matrix.shape[1]):
8         ax.text(x=j, y=i,s=conf_matrix[i, j], va='center',
9                  ha='center', size='xx-large')
10
11 plt.xlabel('Predictions', fontsize=18)
12 plt.ylabel('Actuals', fontsize=18)
13 plt.title('Confusion Matrix', fontsize=18)
14 plt.show()
```

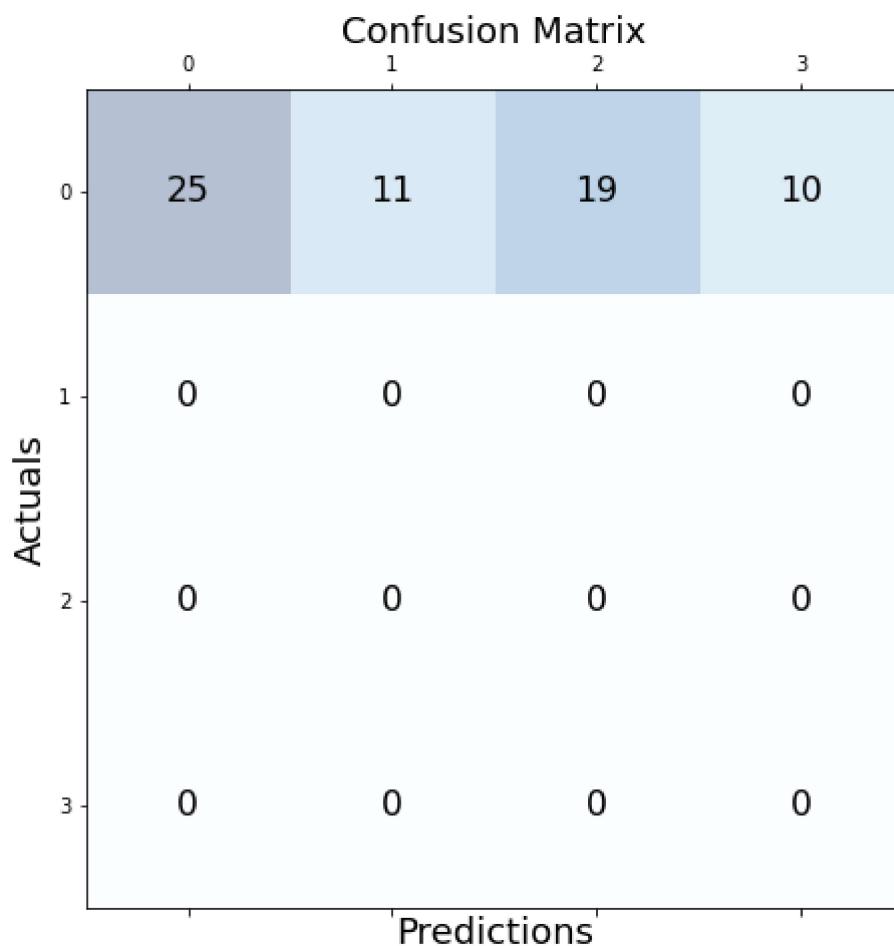


```
In [154]: ► 1 # The Lowest SSE Value  
2 # A good model is one with low inertia AND a small number  
3 # of clusters  
4  
5 kmmedoids_3.inertia_
```

Out[154]: 387474.84

```
In [155]: ► 1 # 4 Clusters  
2  
3 kmmedoids_4 = KMedoids(n_clusters=4, random_state=0).fit(conv_layer_output)
```

```
In [156]: ► 1 from sklearn.metrics import confusion_matrix
  2 conf_matrix=confusion_matrix(y, kmmedoids_4.labels_)
  3
  4 fig, ax = plt.subplots(figsize=(7.5, 7.5))
  5 ax.matshow(conf_matrix, cmap=plt.cm.Blues, alpha=0.3)
  6 for i in range(conf_matrix.shape[0]):
  7     for j in range(conf_matrix.shape[1]):
  8         ax.text(x=j, y=i,s=conf_matrix[i, j], va='center',
  9                 ha='center', size='xx-large')
 10
 11 plt.xlabel('Predictions', fontsize=18)
 12 plt.ylabel('Actuals', fontsize=18)
 13 plt.title('Confusion Matrix', fontsize=18)
 14 plt.show()
```



```
In [157]: ► 1 # The Lowest SSE Value
  2 # A good model is one with low inertia AND a small number
  3 # of clusters
  4
  5 kmmedoids_4.inertia_
```

Out[157]: 383652.62

```
In [ ]: ► 1
```

