# CSE 564 - Visualization

## Mini Project 2 Report

This report contains the brief implementation details of all the tasks of Mini Project 2 of CSE 564.

**Practice the three basic tasks of visual data analytics**
*(a) use data from mini project #1 (or other), begin with |N|≥500, |D|≥10)*
*(b) client-server system: python for processing (server), D3 for VIS (client)*

**Dataset**
I have used *FlightData2008.csv* for this lab. This data has 30 dimensions and has 10000 rows. Half of the columns have non numerical or redundant data. So I have deleted those columns in the start so that the input file is has 15 dimensions. Also the data is standardized to scale the units.

**Client-Server System**
I have used Python 3.6 for server side program and employed Python Flask Framework. For the client side visualization, I have used D3 library and JavaScript. AJAX calls are made to fetch data from the server.

**Task 1 - Data clustering and decimation**

*(a) Implement random sampling and stratified sampling.*
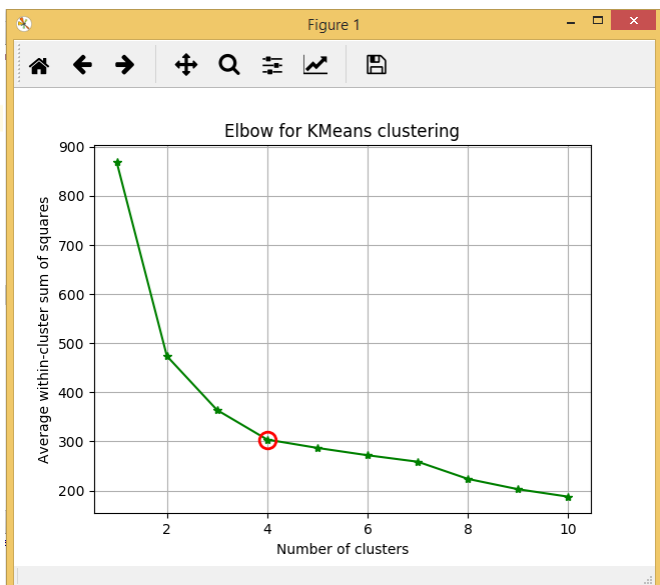*(b) The latter includes the need for k-means clustering (optimize k using elbow)*

Random sampling and Stratified sampling is implemented with a sample size of 1000.
To implement stratified sampling, we need to cluster the data in k groups. The optimum value of k is calculated using the elbow method in k clustering which comes as 4. The plot and code snippet is shown below -

```python
def random_sampling():
    # Randomized sampling
    print("Getting random samples");
    global data
    global input_file
    global random_samples
    global sample_size
    features = input_file[columns]
    data = np.array(features)
    rnd = random.sample(range(len(input_file)), sample_size)
    for j in rnd:
        random_samples.append(data[j])

def clustering():
    # Clustering the data
    print("Clustering data with K = 4");
    global input_file
    features = input_file[columns]
    kmeans = KMeans(n_clusters=4)
    kmeans.fit(features)
    labels = kmeans.labels_
    input_file['kcluster'] = pandas.Series(labels)

def adaptive_sampling():
    # Adaptive sampling
    print("Getting adaptive samples");
    global input_file
    global adaptive_samples
    global sample_size
```
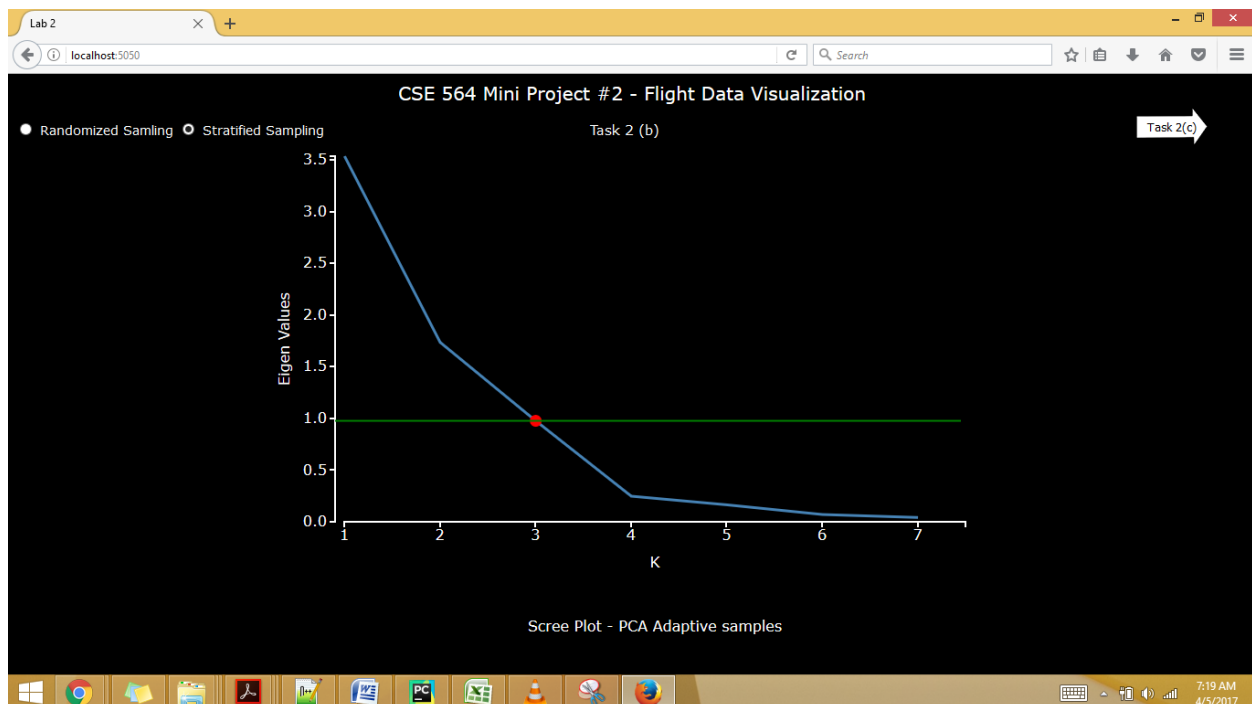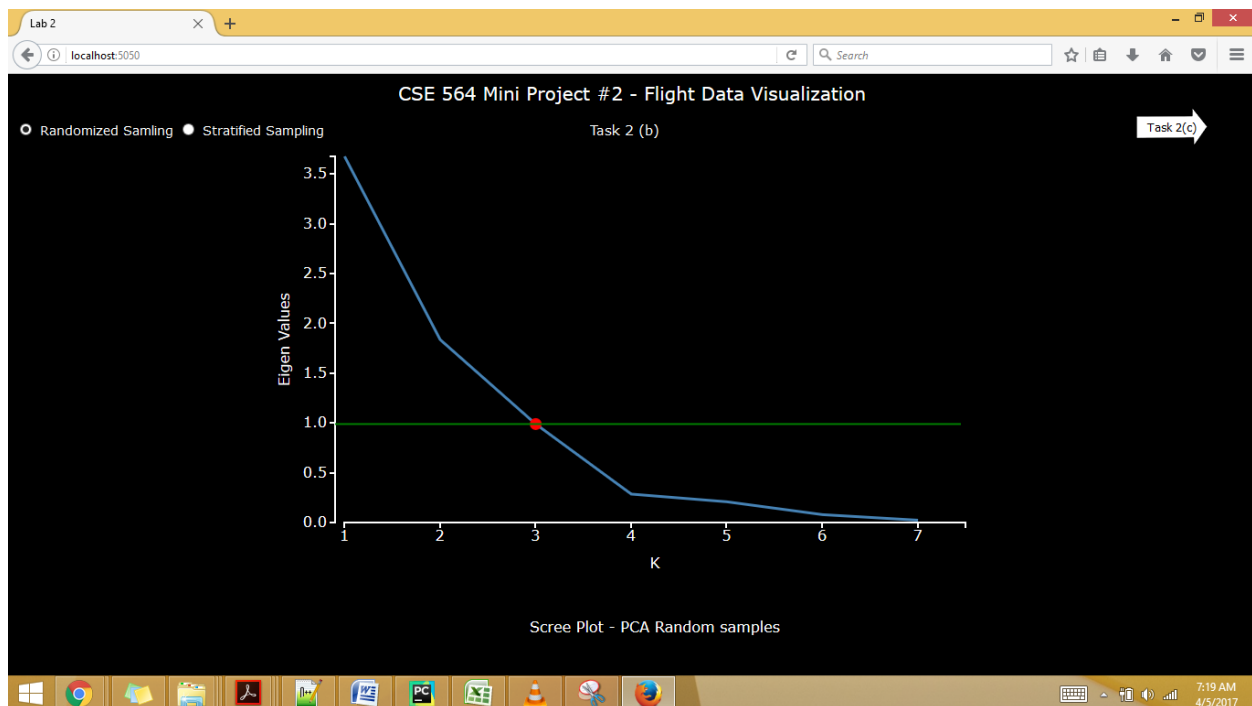


Elbow for KMeans clustering

## Task 2 - Dimension reduction (use decimated data)

*(a) find the intrinsic dimensionality of the data using PCA*
*(b) produce scree plot visualization and mark the intrinsic dimensionality*
*(c) obtain the three attributes with highest PCA loadings*

Using the Principle component analysis technique, eigen values for each feature is calculated and it is plotted on a scree plot. The components with eigen values more than 1 are considered as the intrinsic dimensions. After drawing the scree plot for both random samples and adaptive samples, I came to conclusion that *there are 3 intrinsic dimensions* in this data. This can be seen in the scree plot as marked with the red dot in figures below -
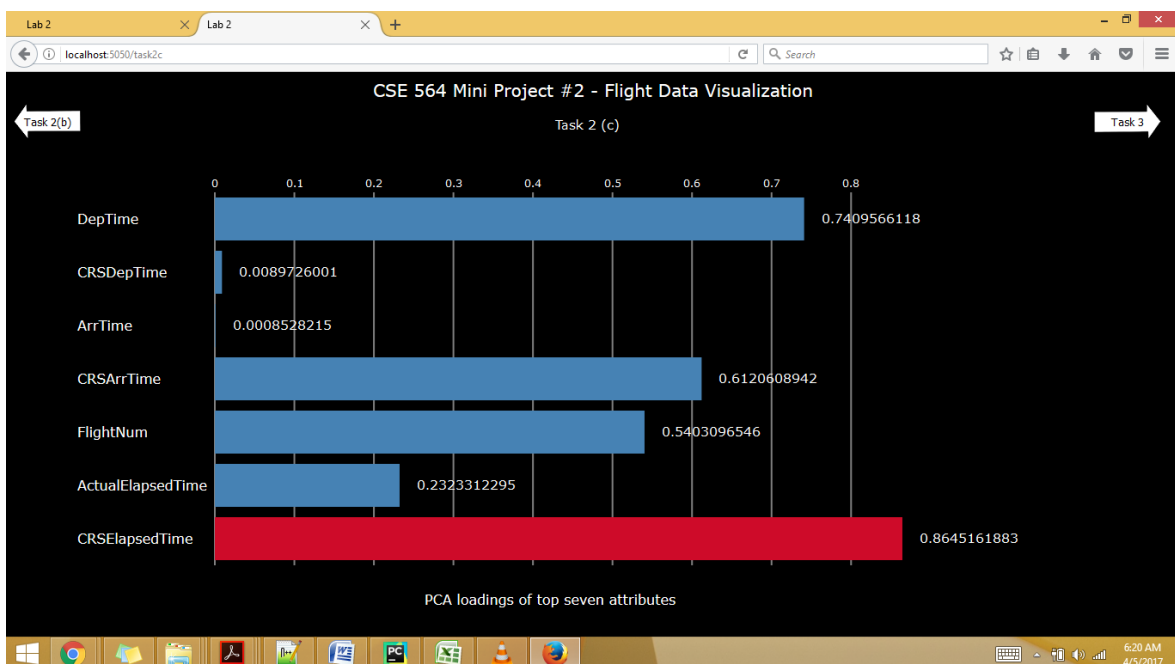
Once I get the intrinsic dimensionality of the data as 3, I calculated the eigen vector for all the features and their squared loadings. The code snippet to calculate the square loadings is shown below -

```
133    def generate_eig_values(data):
134        centered_matrix = data - np.mean(data, axis=0)
135        cov = np.dot(centered_matrix.T, centered_matrix)
136        eig_values, eig_vectors = np.linalg.eig(cov)
137
138        idx = eig_values.argsort()[::-1]
139        eig_values = eig_values[idx]
140        eig_vectors = eig_vectors[:, idx]
141        return eig_values, eig_vectors
142
143    def plot_intrinsic_dimensionality_pca(data, k):
144        # print("Inside plot_intrinsic_dimensionality_pca")
145        global loadingVector
146        [eigenValues, eigenVectors] = generate_eig_values(data)
147        idx = eigenValues.argsort()[::-1]
148        eigenValues = eigenValues[idx]
149        eigenVectors = eigenVectors[:, idx]
150        squaredLoadings = []
151        ftrCount = len(eigenVectors)
152        for ftrId in range(0,ftrCount):
153            loadings = 0
154            temp = []
155            for compId in range(0, k):
156                loadings = loadings + eigenVectors[compId][ftrId] * eigenVectors[compId][ftrId]
157            loadingVector[columns[ftrId]] = loadings
158            squaredLoadings.append(loadings)
159
```

Once the square loading are calculated, we can determine the top three attributes with highest loadings. These attributes are -

1. CRSElapsedTime
2. DepTime
3. CRSArrTime

This is visualized with the help of a horizontal bar chart as shown in the figure below. The value with highest squared loading is highlighted -

## Task 3 - Visualization (use dimension reduced data)

*(a) visualize data projected into the top two PCA vectors via 2D scatter plot*
*(b) visualize data via MDS (Euclidean & correlation distance) in 2D scatter plots*
*(c) visualize scatter plot matrix of the three highest PCA loaded attributes*
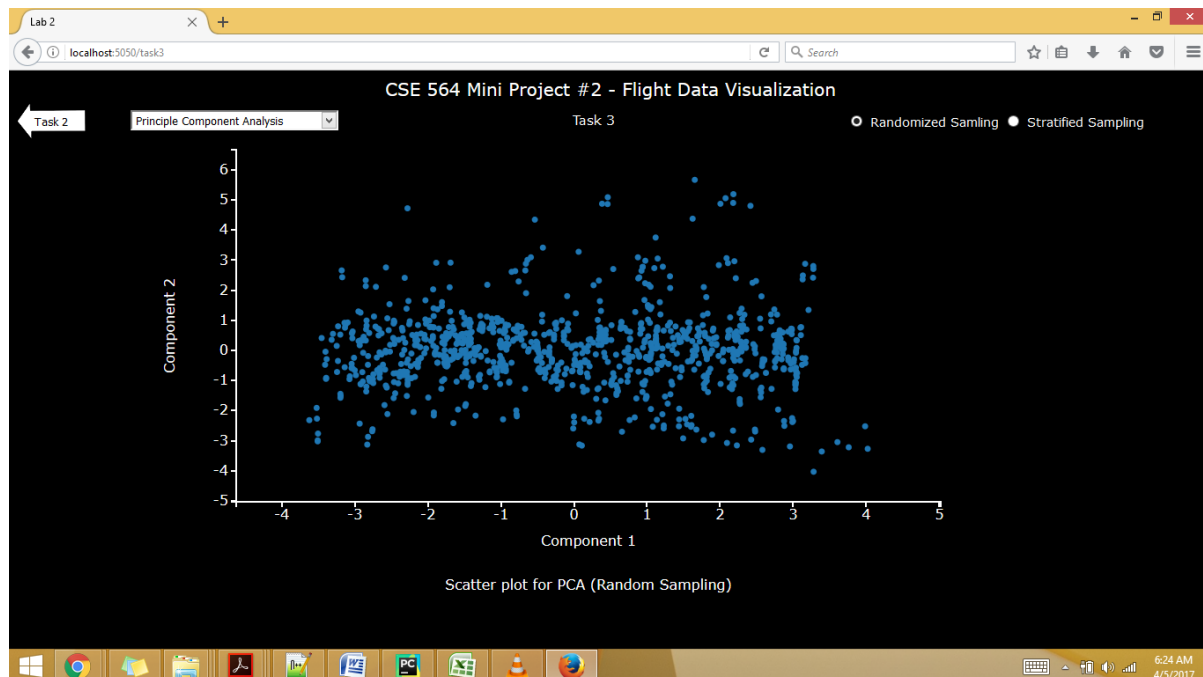
Once I know the squared loadings of all the features, I reduced the dimension of the data taking the top seven attributes which contribute almost 98% in the principal components. These components are -
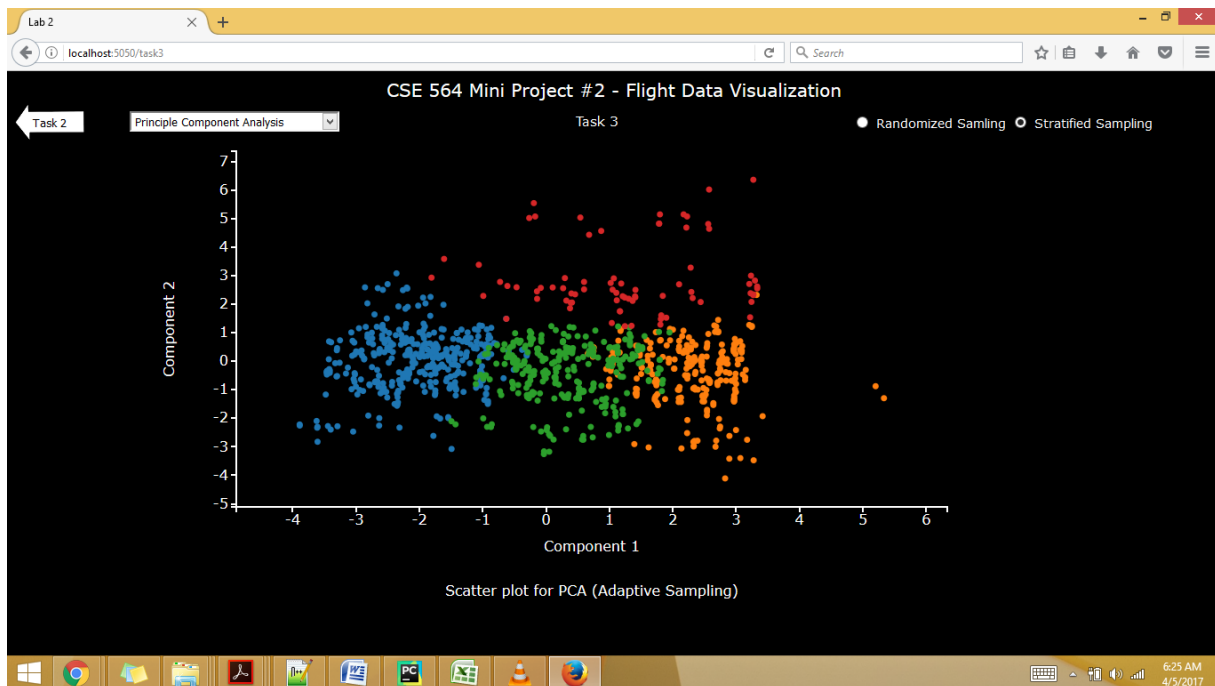
'DepTime', 'CRSDepTime', 'ArrTime', 'CRSArrTime', 'FlightNum', 'ActualElapsedTime', 'CRSElapsedTime'

Now I pass these columns of the random and adaptive samples from task 1 to the function written to calculate the principle components. PCA library of sklearn.decomposition in Python is used to calculate the output for 2 PCA. The code snippet is shown below -

```
220    @app.route("/pca_random")
221    def pca_random():
222        print("Inside PCA Random");
223        # PCA reduction with random sampling
224        data_col = []
225        try:
226            global random_samples
227            global imp_fetures
228            pca_data = PCA(n_components=2)
229            X = random_samples
230            pca_data.fit(X)
231            X = pca_data.transform(X)
232            data_col = pandas.DataFrame(X)
233
234            for i in range(0, 2):
235                data_col[columns[imp_fetures[i]]] = orig_file[columns[imp_fetures[i]]][:sample_size]
236
237            data_col['clusterid'] = input_file['kcluster'][:sample_size]
238
239        except:
240            e = sys.exc_info()[0]
241            print(e)
242        return pandas.json.dumps(data_col)
```

The scatter plots of PCA for random and adaptive samples are shown below -



Scatter plot for PCA (Random Sampling)

These plots show that the PCA technique can be used to reduce the dimensions of the data without losing much information. The data is transformed such that component with the maximum variation is plotted on the X-axis and component with the second best variation is plotted on the Y-axis. The clusters in the scatter plot with adaptive sampling clearly shows that the data with similar features are clustered together.

Similarly, the MDS technique is employed to show its dimension reduction capabilities. Both Euclidean distance and correlation distance methods are used on both random and adaptive samples from task 1. MDS and pairwise_distance form sklearn.metrices are used to calculate the pair-wise distance and the MDS. The code snippets are shown below -

```python
278    @app.route("/mds_euclidean_random")
279    def mds_euclidean_random():
280        print("Inside MDS Random using Euclidean Distance")
281        # MSD reduction with random sampling and using euclidean distance
282        data_col = []
283        try:
284            global random_samples
285            global imp_fetures
286            mds_data = manifold.MDS(n_components=2, dissimilarity='precomputed')
287            similarity = pairwise_distances(random_samples, metric='euclidean')
288            X = mds_data.fit_transform(similarity)
289            data_col = pandas.DataFrame(X)
290
291            for i in range(0, 2):
292                data_col[columns[imp_fetures[i]]] = orig_file[columns[imp_fetures[i]]][:sample_size]
293
294            data_col['clusterid'] = input_file['kcluster'][:sample_size]
295        except:
296            e = sys.exc_info()[0]
297            print(e)
298        return pandas.json.dumps(data_col)
```
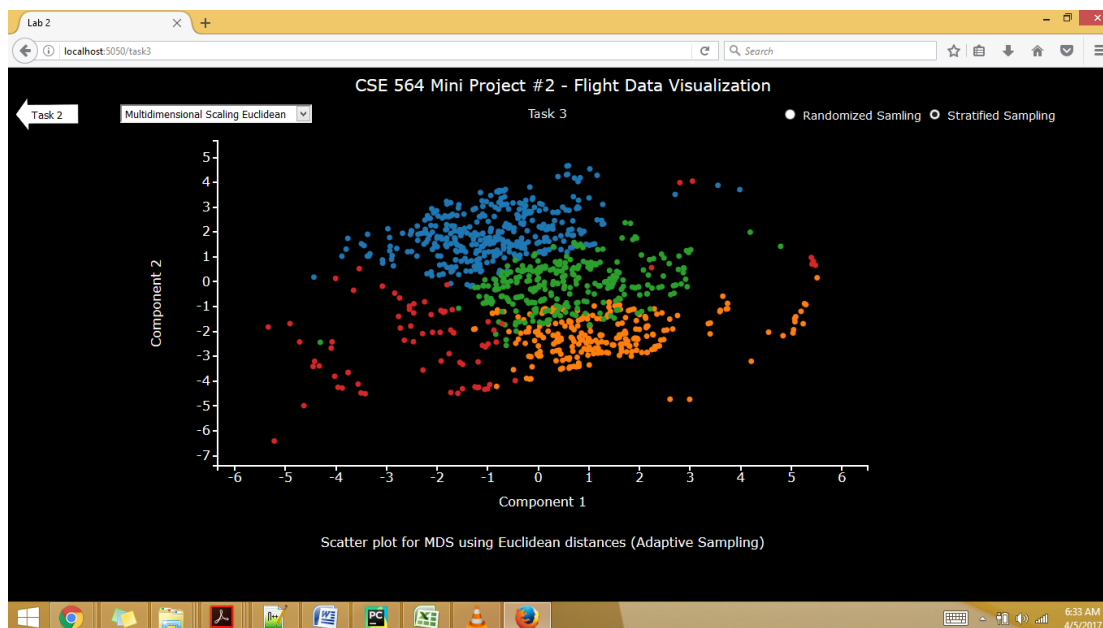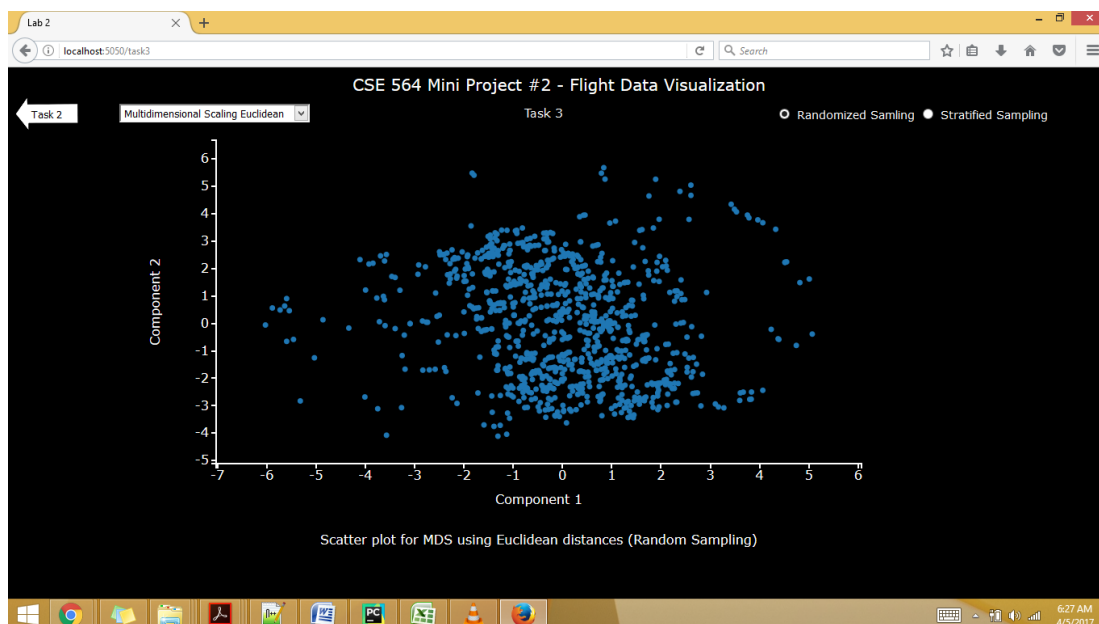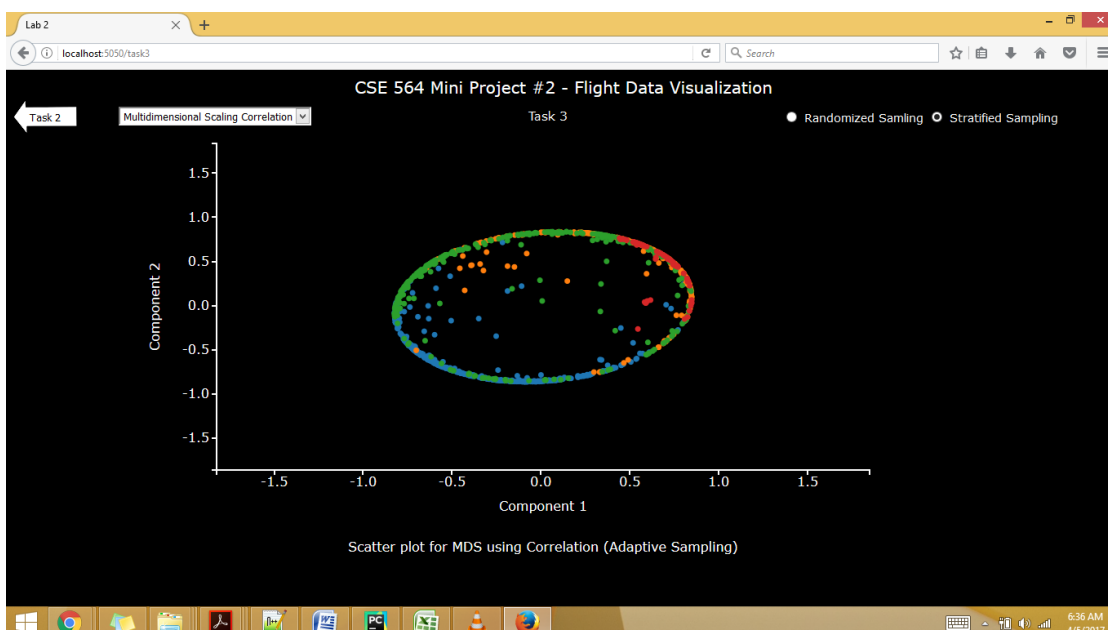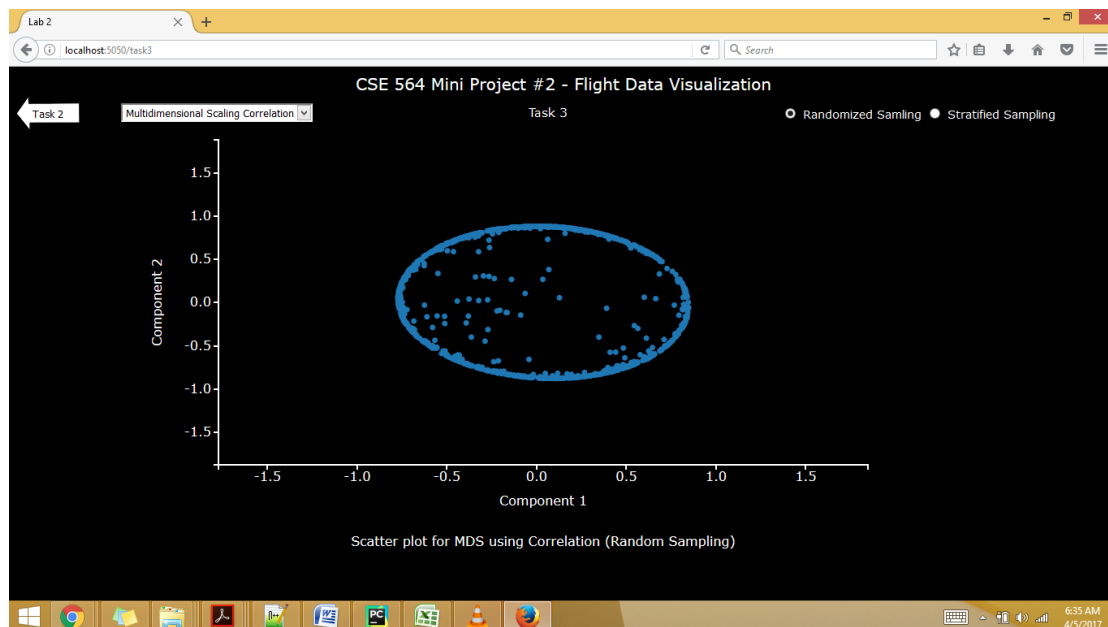
```
328    @app.route("/mds_correlation_random")
329    def mds_correlation_random():
330        print("Inside MDS Random using Correlation")
331        # MSD reduction with random sampling and using Correlation
332        data_col = []
333        try:
334            global random_samples
335            mds_data = manifold.MDS(n_components=2, dissimilarity='precomputed')
336            similarity = pairwise_distances(random_samples, metric='correlation')
337            X = mds_data.fit_transform(similarity)
338            data_col = pandas.DataFrame(X)
339
340            for i in range(0, 2):
341                data_col[columns[imp_fetures[i]]] = orig_file[columns[imp_fetures[i]]][:sample_size]
342
343            data_col['clusterid'] = input_file['kcluster'][:sample_size]
344
345        except:
346            e = sys.exc_info()[0]
347            print(e)
348        return pandas.json.dumps(data_col)
```

We get the following four scatter plots for the MDS techniques -



Scatter plot for MDS using Euclidean distances (Random Sampling)



Scatter plot for MDS using Euclidean distances (Adaptive Sampling)

Scatter plot for MDS using Correlation (Random Sampling)



Scatter plot for MDS using Correlation (Adaptive Sampling)

MDS finds a set of vectors in p-dimensional space such that the matrix of Euclidean distances among them corresponds as closely as possible to some function of the input matrix according to a criterion function called stress. The plot depicts the same behavior as all the points are closely packed together. The Euclidean distance groups the points according to the Euclidean distances while the correlation method tries to put variables with high positive correlations near each other, and variables with strong negative correlations far apart. The plots with the adaptive samples depict that the points in the same clusters are placed closed to each other as they will have very less Euclidean distance (in the Euclidean plots) and high Correlation (in the correlation plot).

Finally, I have visualized the scatterplot matrix of the three highest loaded attributes calculated in the task 2(c). The scatterplot shows the pair-wise relation with one another. There are 9 sub plots showing the variance and covariance between the parameters placed on x-axis and y-axis. The adaptive plot again shows that points in the same clusters are closely placed as they are similar.

Scatterplot Matrix for top 3 PCA loaded attributes (Random Sampling)



Scatterplot Matrix for top 3 PCA loaded attributes (Adaptive Sampling)