

# CSE 564 - VISUALIZATION

## LAB 2

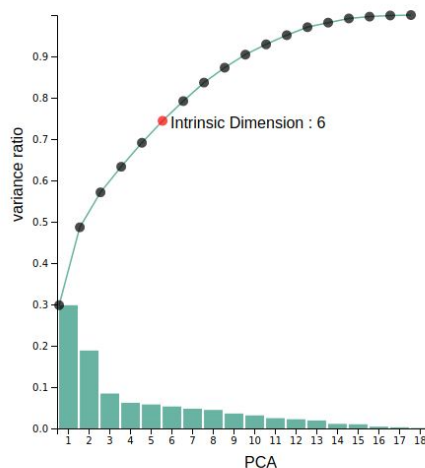
This is the report for lab 2 and contains a brief explanation of it.

**Dataset :** The dataset is taken from kaggle (FIFA 19 players), and has 18 features and 572 samples.

**Technologies used:** A client sever system is used, Having flask as a backend, all the data processing is done in python and then the result are sent to frontend Where they are visualized using D3.js.

**The Loading page :** I have taken a simple fixed side menu bar(only certain tasks were needed to be displayed), a bootstrap template.

LAB2
SCREE DATA
SCREE RANDOM
SCREE STRATIFIED
SCATTER DATA
SCATTER RANDOM
SCATTER STRATIFIED
MDSE DATA
MDSE RANDOM
MDSE STRATIFIED
MDSC DATA
MDSC RANDOM
MDSC STRATIFIED



Task `1:

1. *Implement random sampling and stratified sampling (remove 75% of data)*
2. *The latter includes the need for k-means clustering (optimize k using elbow)*

```

20 @app.route('/')
21 def hello_world():
22     clean_data()
23     return render_template("index.html")
24
25
26 def randomSample():
27     return data.sample(frac=0.25)
28
29 #Use of Kelbow visualizer
30 def elbowCheck():
31     mat = data.values
32     mat = mat.astype(float)
33     model = KMeans()
34     visualizer = KElbowVisualizer(model, k=(1,12))
35     visualizer.fit(mat)
36     visualizer.show(outputpath="static/images/kmeans.png")
37

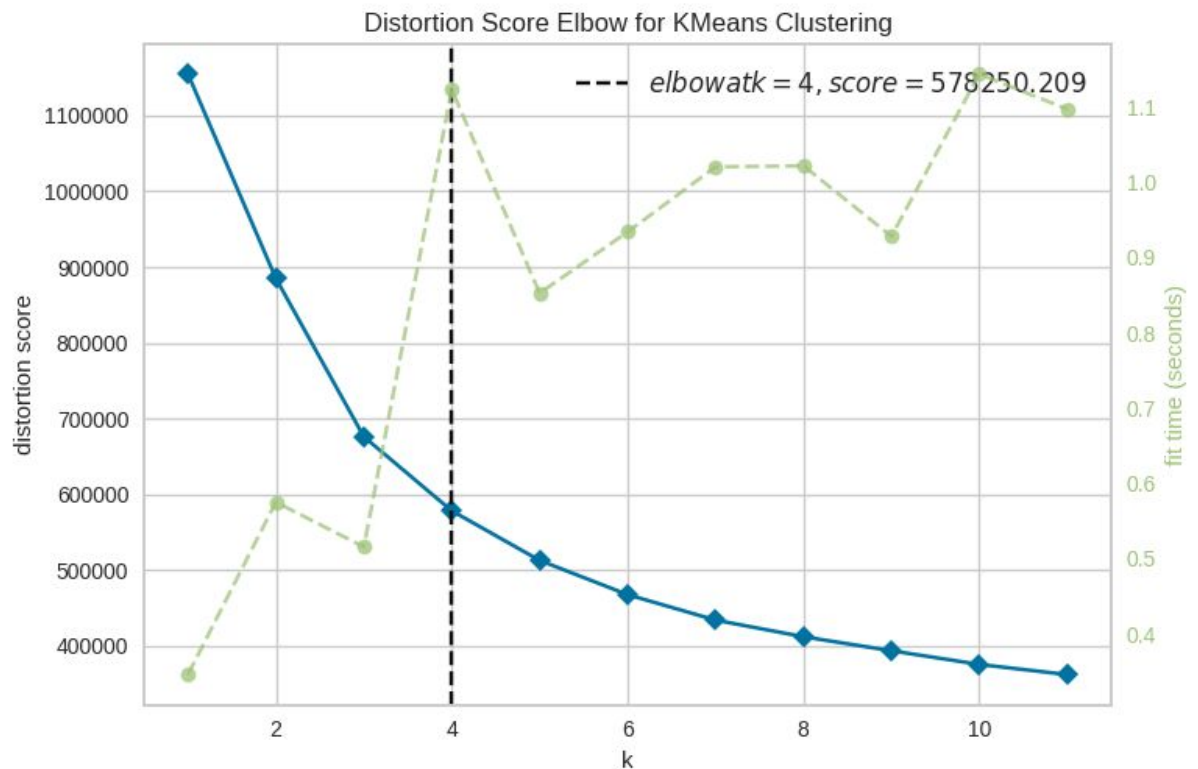
```

Random sampling and stratified sampling is done on the data and only 25% of the data is kept, for stratified sampling KMeans clustering is used and to get the best cluster size Kelbow visualizer is used.

```

38 def stratifiedSample():
39     #From elbow check the cluster size is best when K=4
40     global frame
41     global sample_done
42     global colors
43
44     if not sample_done:
45         nCluster = 4
46         mat = data.values
47         mat = mat.astype(float)
48         kmeans = KMeans(n_clusters=nCluster)
49         kmeans.fit(mat)
50         cInfo = kmeans.labels_
51         #save clusters and recreate dataframe
52         cluster = [[],[],[],[]]
53         for i in range(0,len(data.index)):
54             cluster[cInfo[i]].append(data.loc[i]);
55
56         temp = []
57         colors = []
58
59         for i in range(0,nCluster):
60             tempFrame = pd.DataFrame(cluster[i]).sample(frac=0.25)
61             dfSize = tempFrame.shape[0]
62             for j in range(0,dfSize):
63                 colors.append(i)
64                 temp.append(tempFrame)
65         frame = pd.concat(temp)
66         sample_done = True
67
68     return frame
69

```



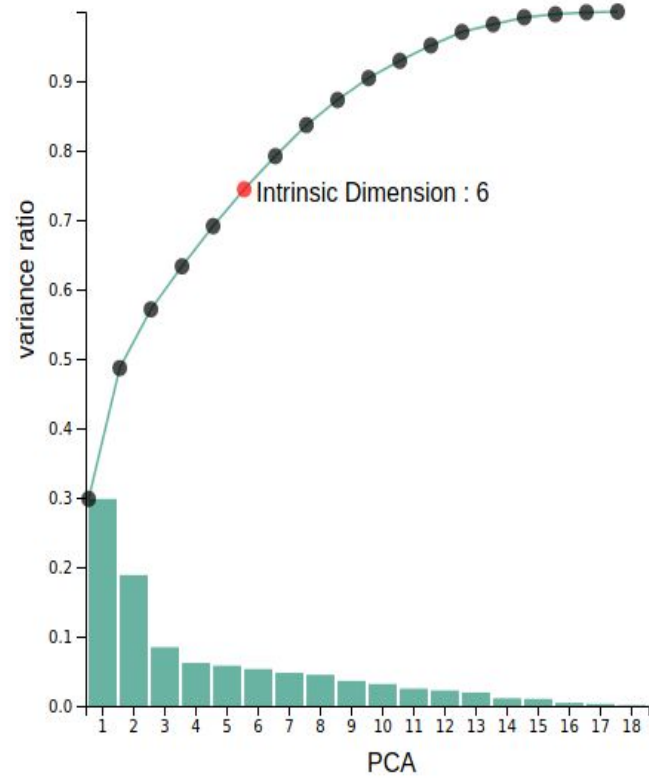
Task 2:

1. Find the intrinsic dimensionality of the data using PCA
2. Produce scree plot visualization and mark the intrinsic dimensionality
3. Show the scree plots before/after sampling to assess the bias introduced
4. Obtain the three attributes with highest PCA loadings

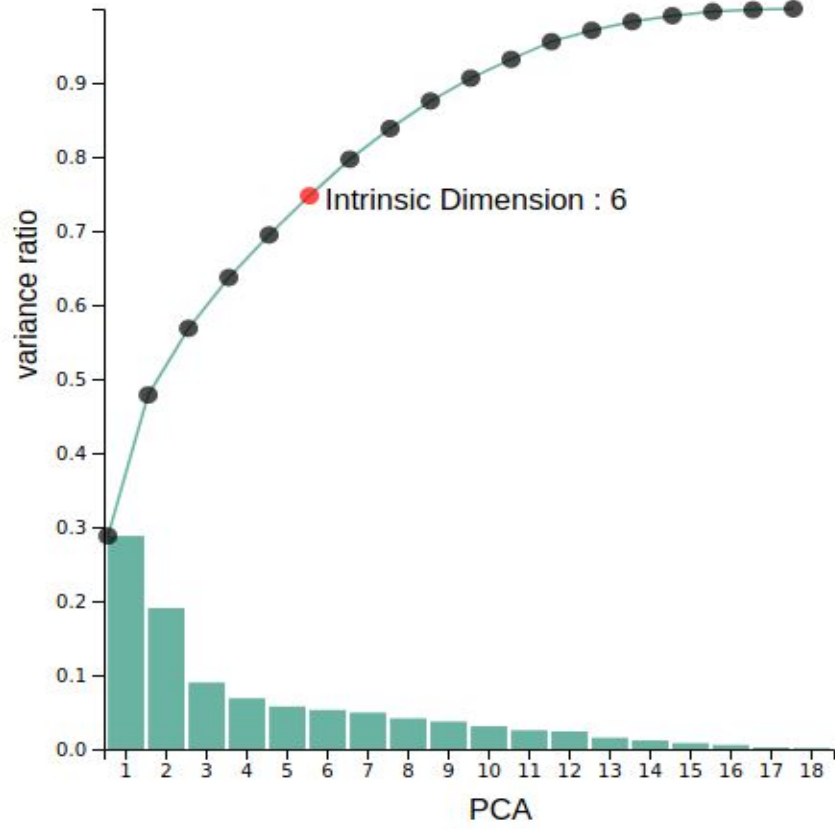
By using the PCA technique Intrinsic dimension of the data is found, when 75% of variance is accumulated we stop and consider it to be the intrinsic dimension.

### Scree Plots with Intrinsic Dimension.

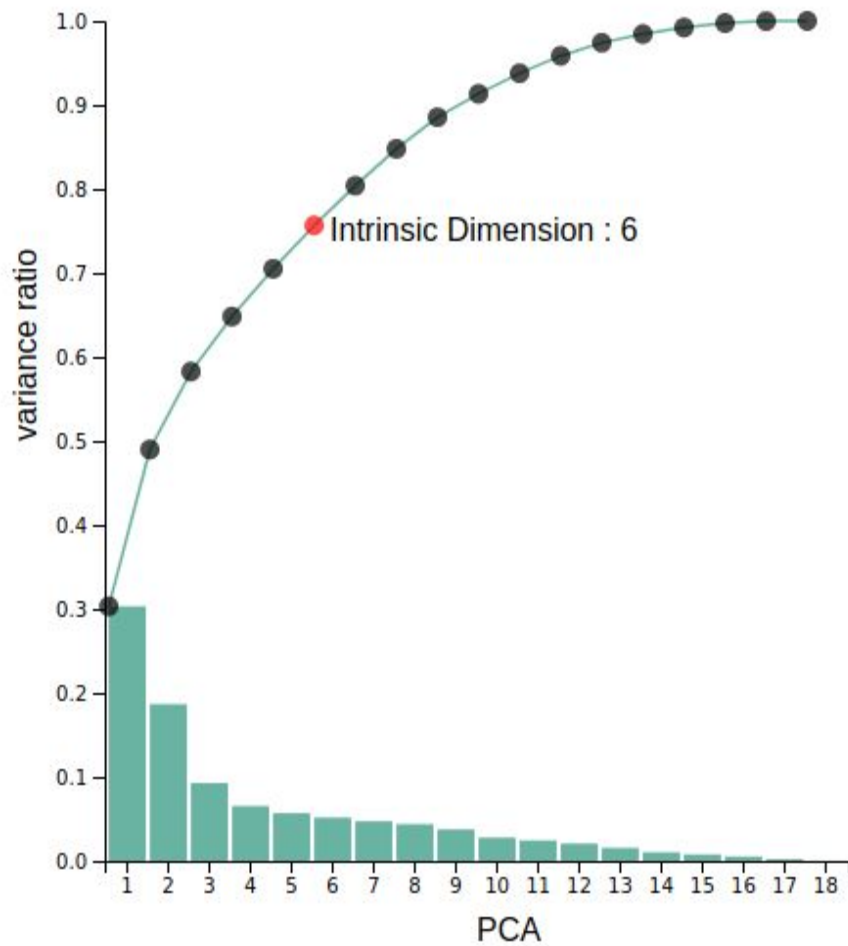
1. Scree Plot original data



## 2. Scree Plot Random sampling



### 3. Scree Plot stratified sampling



**Observation :** We can see that the Intrinsic dimension for the three kind of sampling is Same, we can for now think that the samples are in a good way representing the data.

**Code snippets:**

```

70
71 def pcaAnalysis(sample):
72
73     global sq_load
74     global intD
75
76     mat = sample.values
77     mat = mat.astype(float)
78     mat = StandardScaler().fit_transform(mat)
79     nComponents = 18
80     pca = PCA(n_components=nComponents)
81     pca.fit_transform(mat)
82     count = 0
83     cumsum = 0
84     for eigV in pca.explained_variance_ratio_:
85         cumsum += eigV
86         if cumsum > 0.78:
87             break
88
89         count += 1
90     intD = count
91     #get the loading matrix
92     loadings = pca.components_.T * np.sqrt(pca.explained_variance_)
93     sq_loadings = np.square(loadings)
94     sq_load = np.sum(sq_loadings[:,0:3],axis=1)
95     topPcaLoad()
96     return pca.explained_variance_ratio_
97

```

Common code for doing pca, various types of samples are passed to this function.

Loading Matrix is calculated using the formula.

**loadings = pca.components\_.T \* np.sqrt(pca.explained\_variance\_)**

The top three attributes with highest pca loadings come out to be

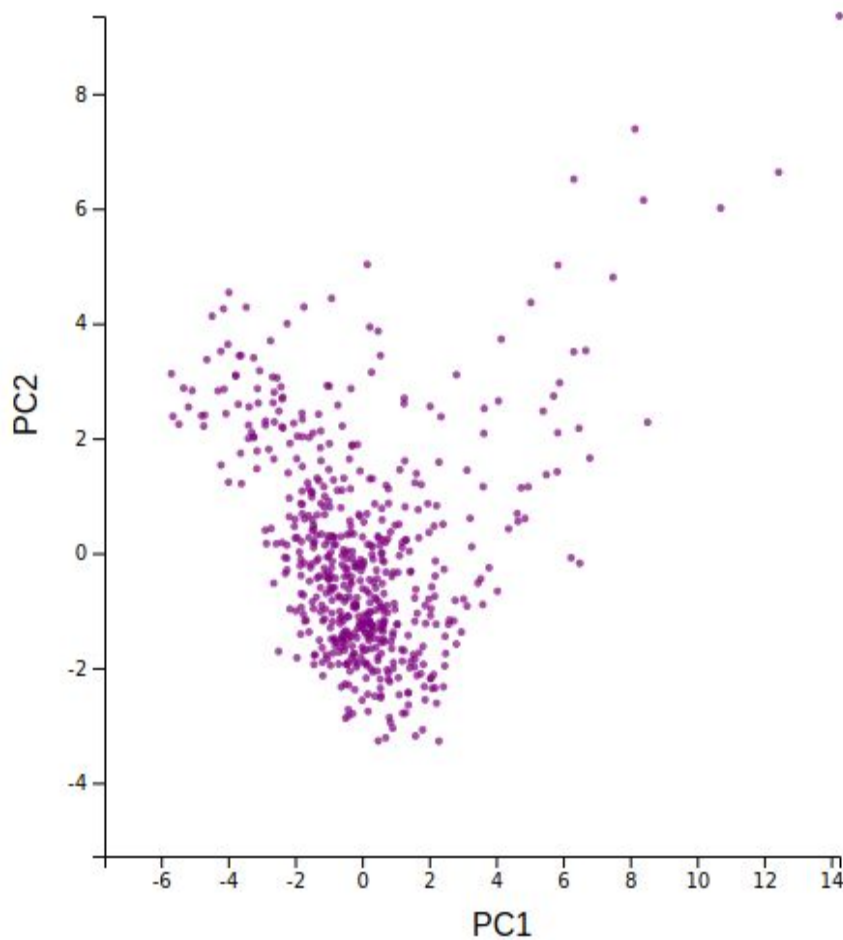
1. 'Value',
2. 'Acceleration',
3. 'Release Clause'

Task 3:

1. Visualize the data projected into the top two PCA vectors via 2D scatterplot.
2. Visualize the data via MDS (Euclidian & correlation distance) in 2D scatterplots.
3. Visualize the scatterplot matrix of the three highest PCA loaded attributes.

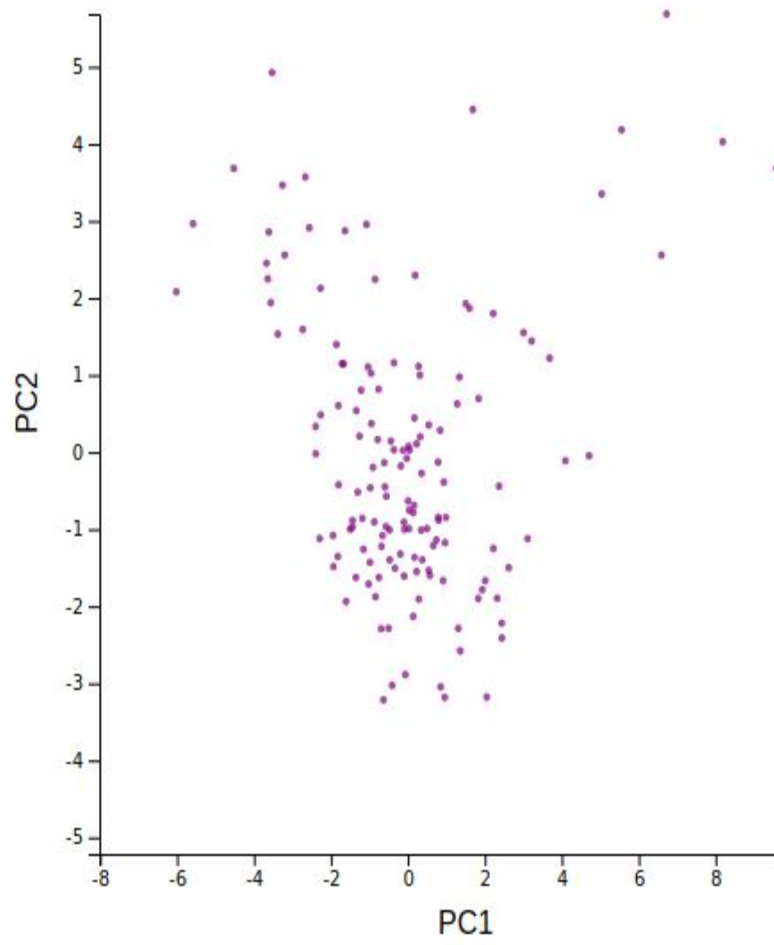
Now using PCA and components to be 2 , can produce a scatter plot.

1. Scatter plot original data.

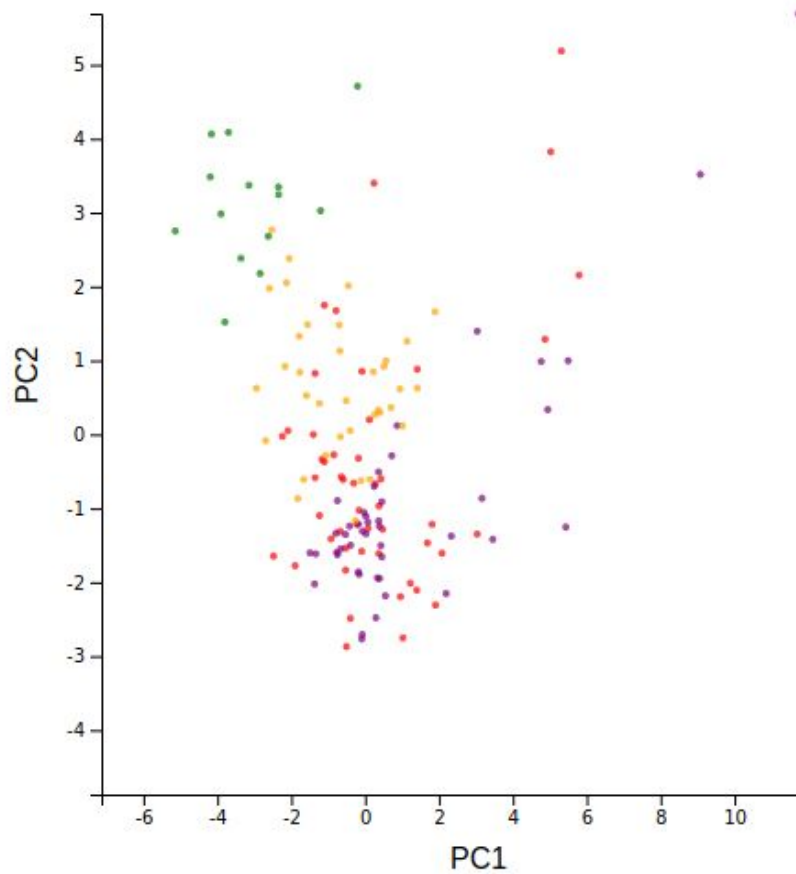


2. Scatter Plot random sampling.





3. Scatter Plot Stratified sampling.

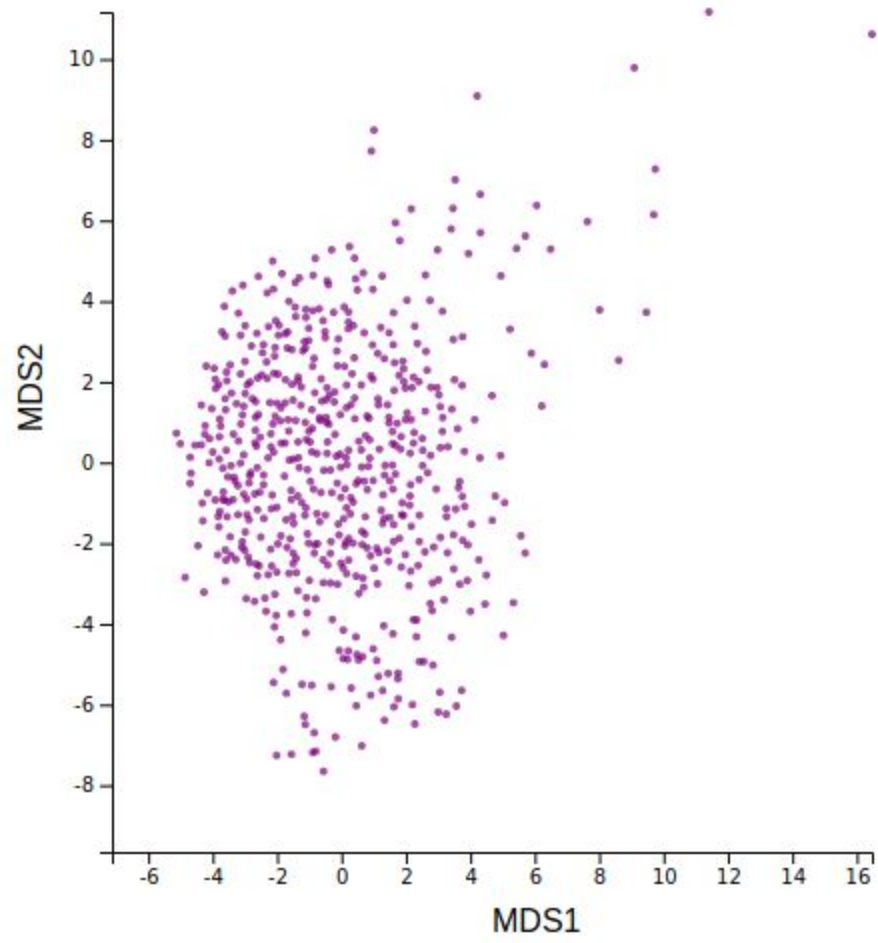


**Observation:** We can see the data is indeed clustered and looks similar to original data, more clarity about it comes from the fact that the top Pca loadings of the original and stratified data are same.

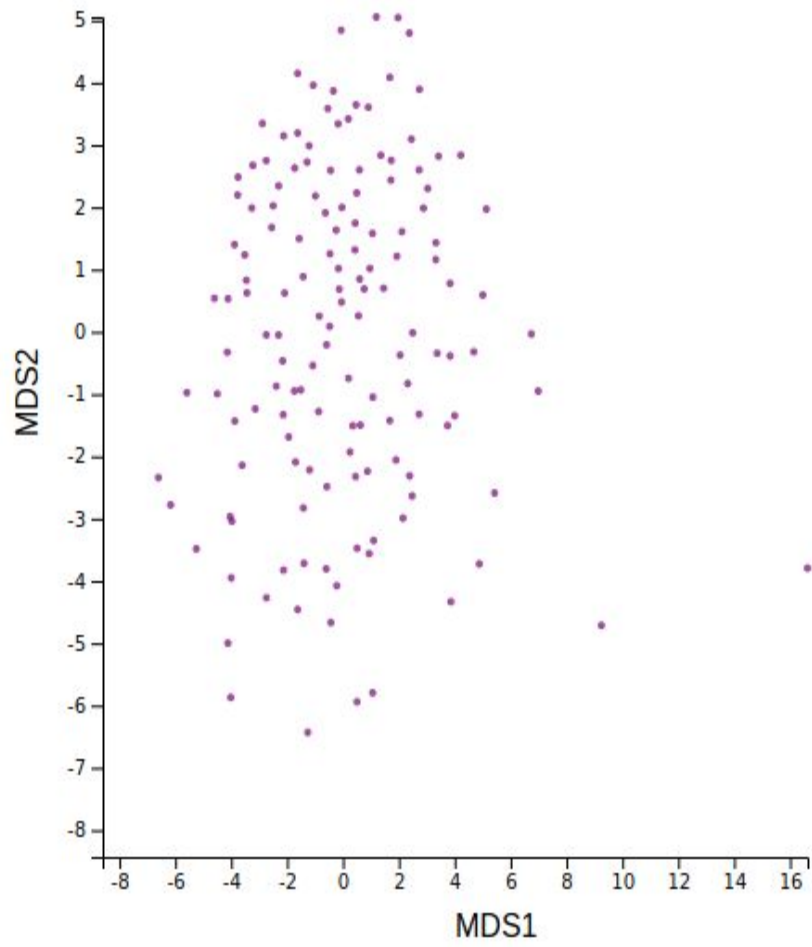
MDS dimension reduction

Euclidean:

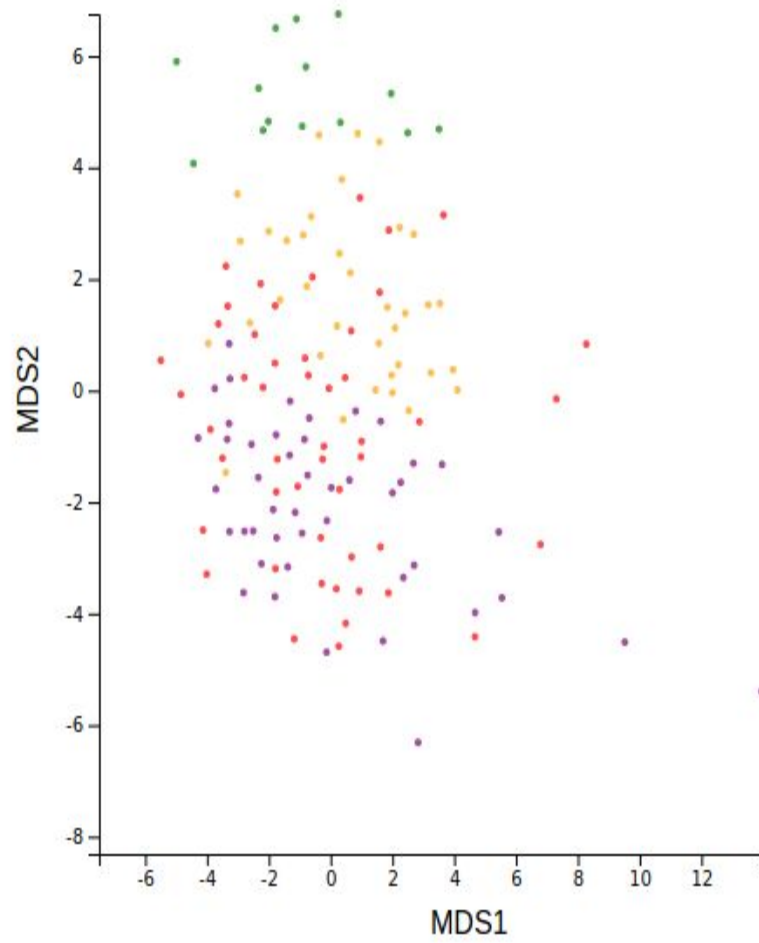
1. Scatter plot Original Data.



2. Scatter Plot random data.

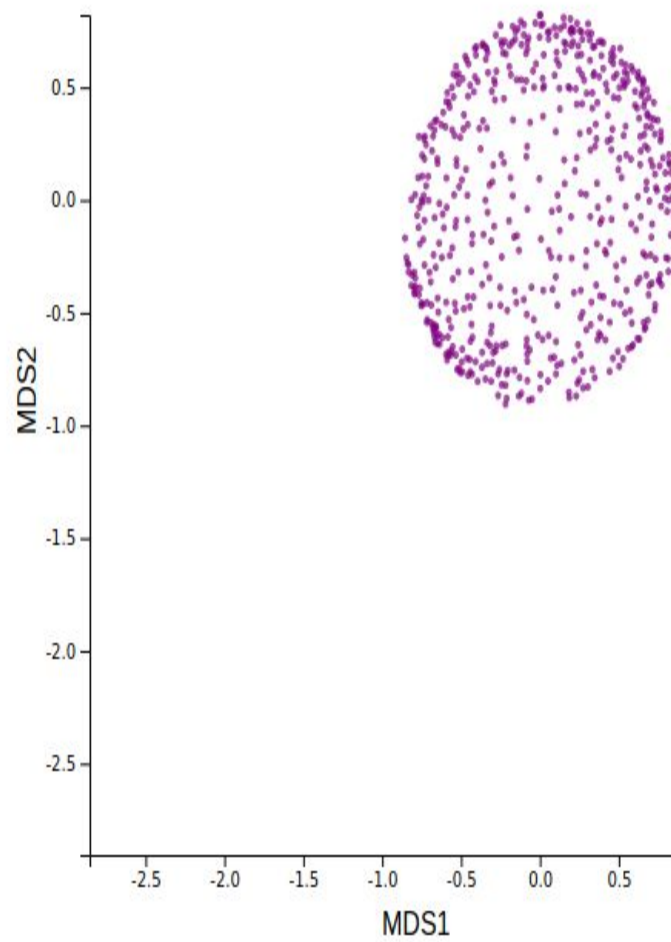


3. Scatter plot Stratified data.

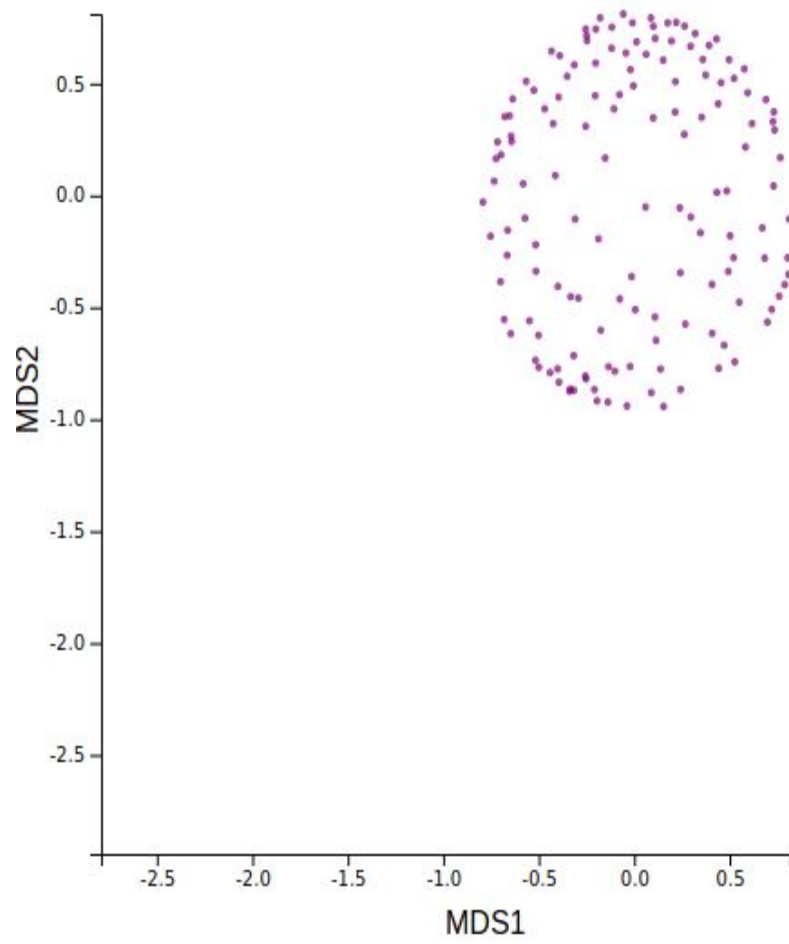


MDS Correlation

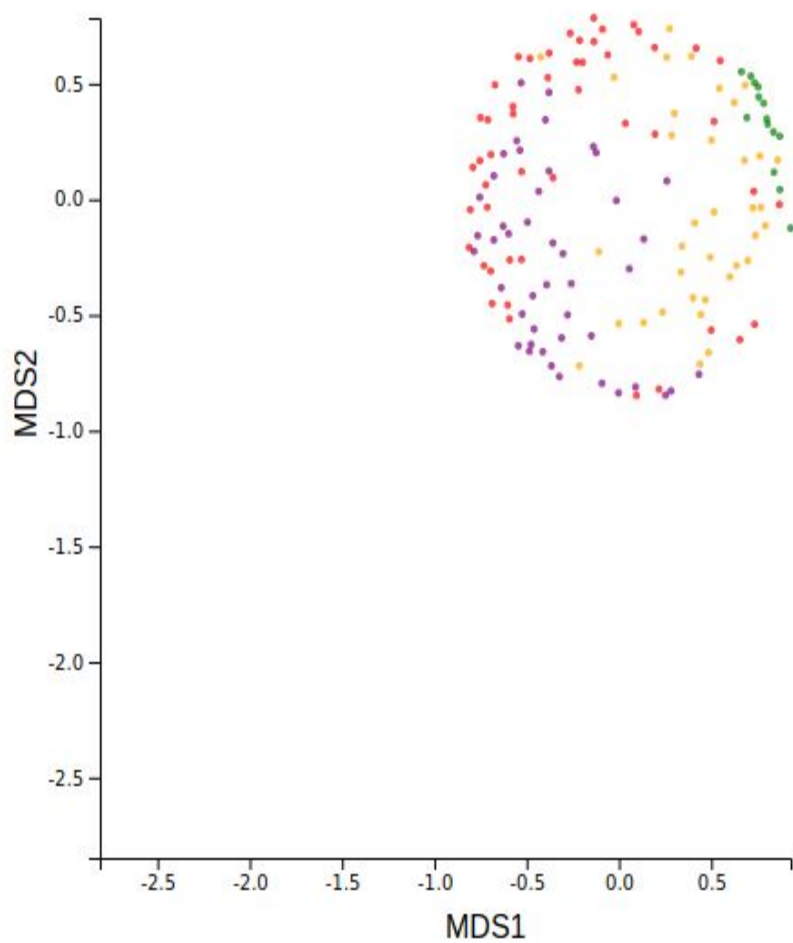
1. Scatter plot original data.



2. Scatter plot random data.



3. Scatter plot stratified data.



Code snippets:

```
#mds Euclidian
def mdsEScatter(sample):

    mat = sample.values
    mat = mat.astype(float)
    mat = StandardScaler().fit_transform(mat)
    mds = manifold.MDS(n_components=2, dissimilarity='precomputed')
    similarity = pairwise_distances(mat, metric='euclidean')
    X = mds.fit_transform(similarity)
    return X
```



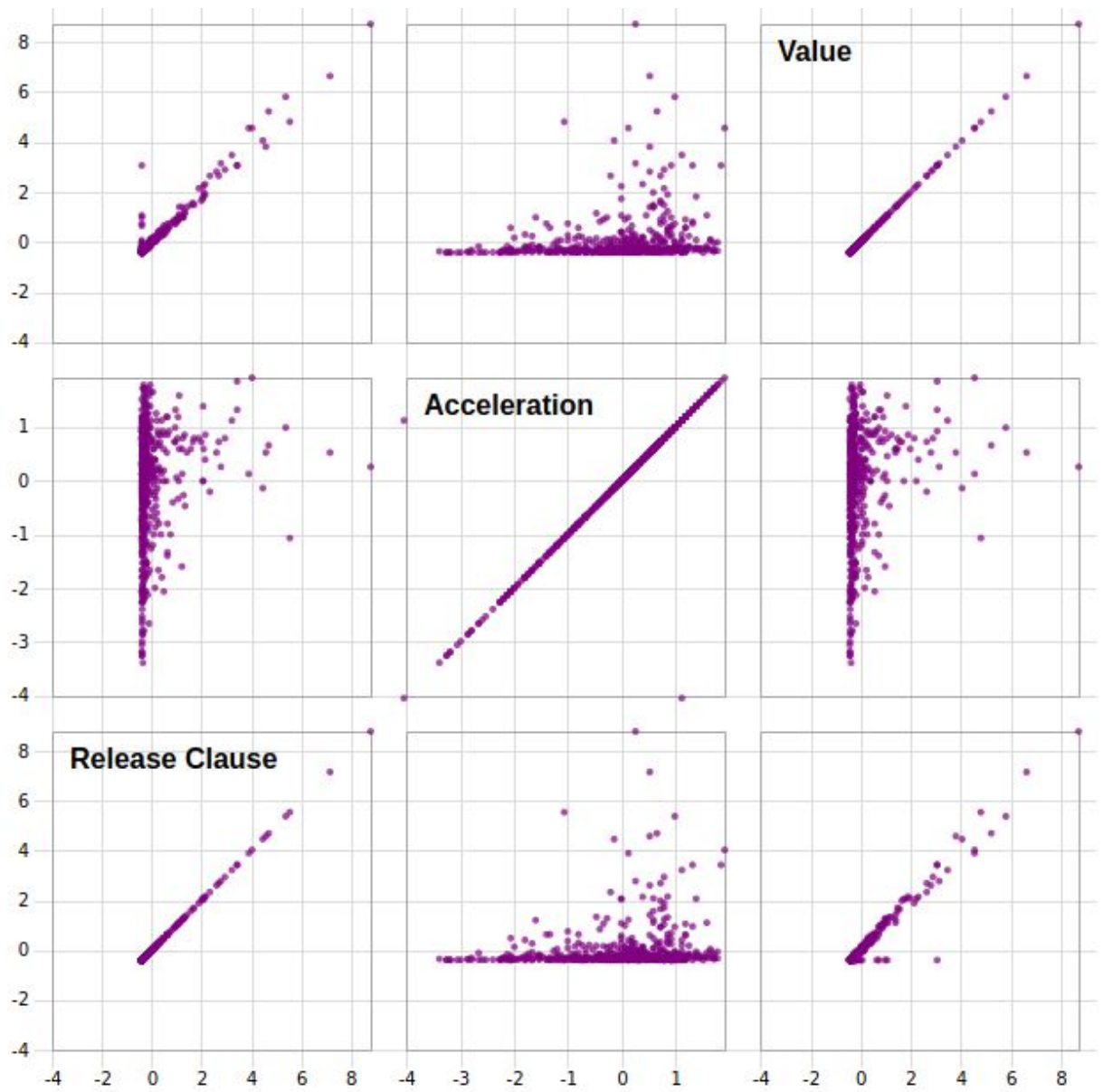
```
#mds correlation
def mdsCScatter(sample):

    mat = sample.values
    mat = mat.astype(float)
    mat = StandardScaler().fit_transform(mat)
    mds = manifold.MDS(n_components=2, dissimilarity='precomputed')
    similarity = pairwise_distances(mat, metric='correlation')
    X = mds.fit_transform(similarity)
    return X
```

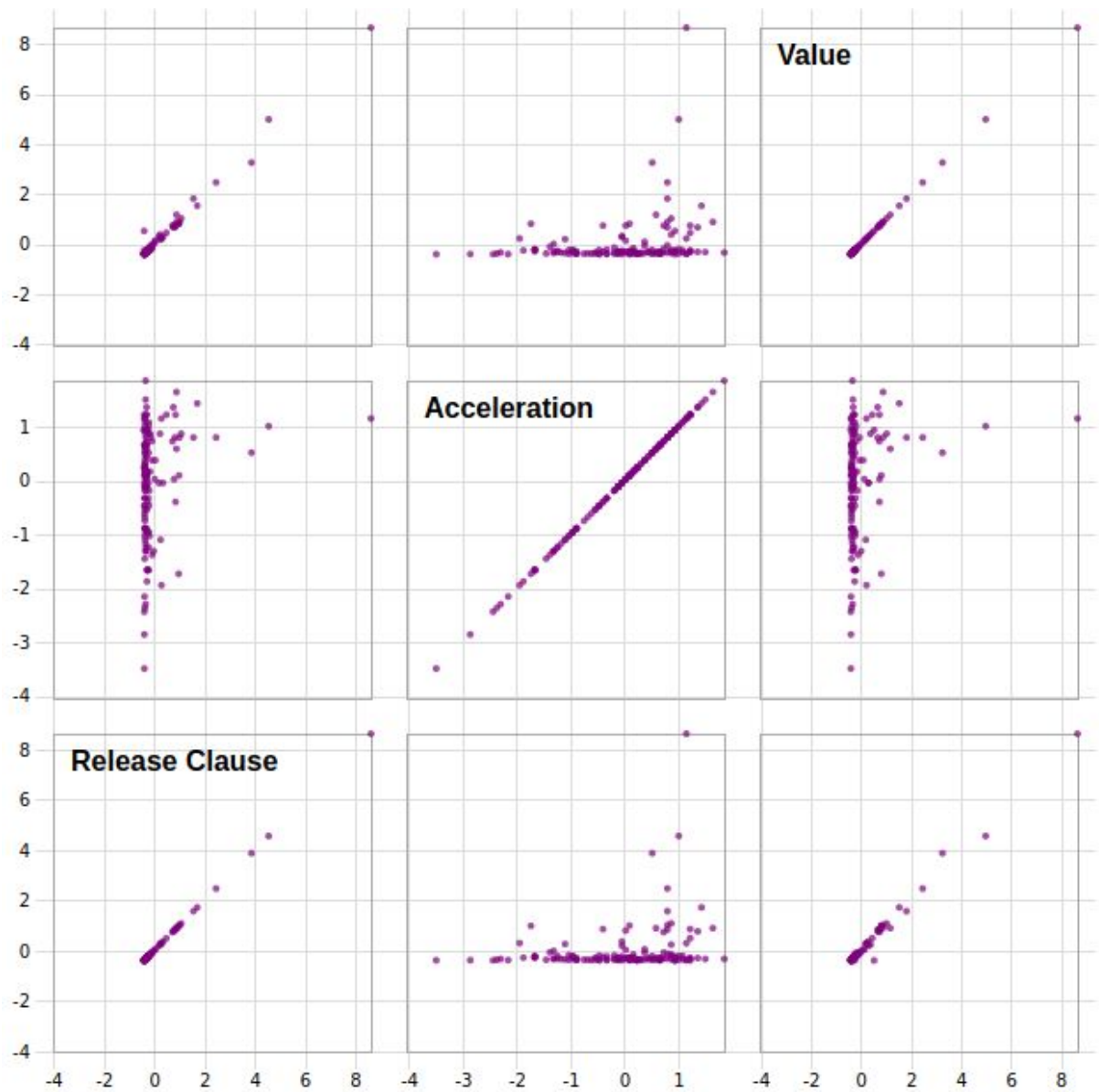
For pca the same techniques as in part 2 is used.

Scatter matrix:

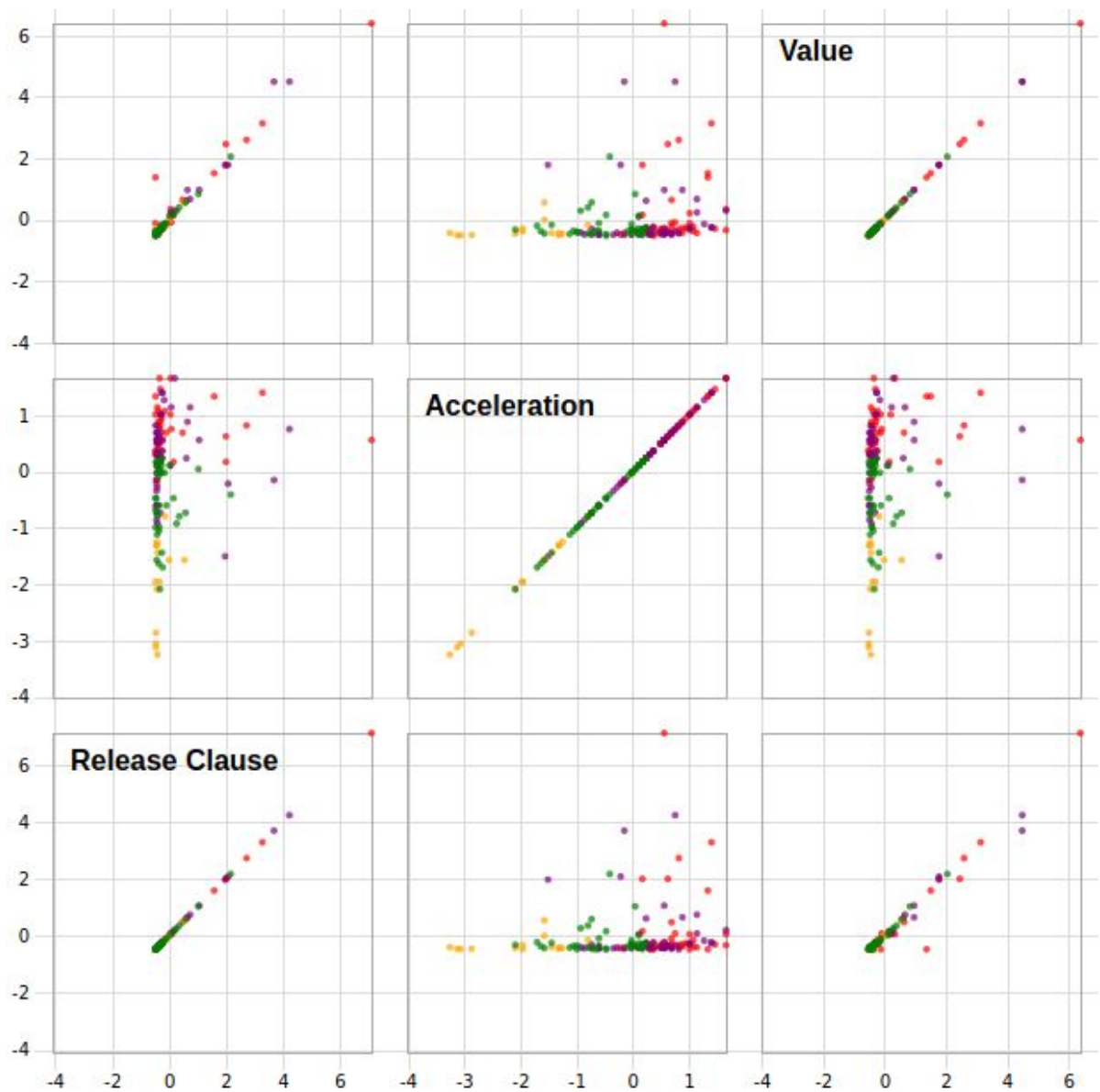
1. Scatter matrix original data.



2. Scatter matrix Random data.



3. Scatter Matrix stratified data.



**Observation :** A positive correlation exist between 'Release Clause' and 'Value'.

