# Link for video: https://youtu.be/GeUfMmVI-90

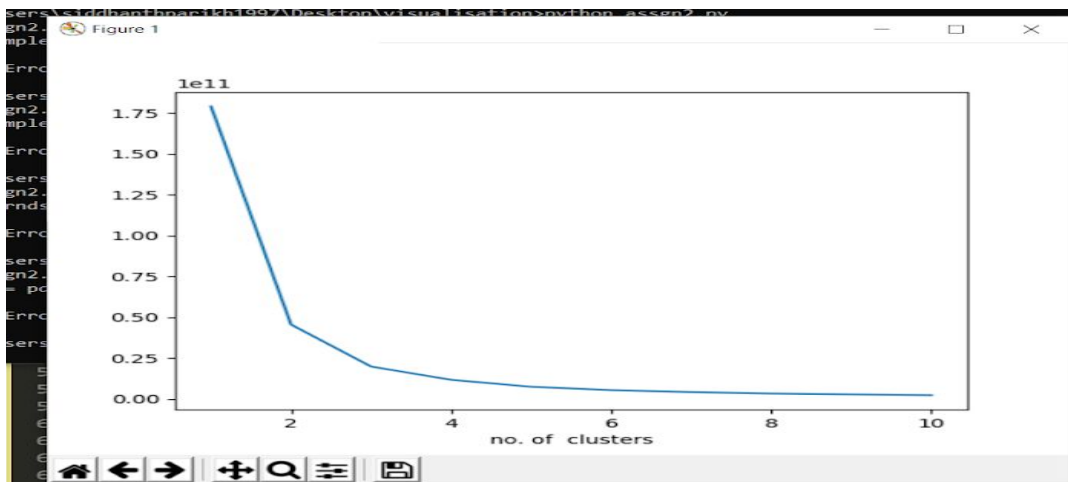# Visualization HW2 Report

1. Chosen Dataset:
   Video Games Sales which lists all the various names and platforms a game is released on along with the user ratings, user scores, critic ratings, critic scores, Age rating and the developer of the game.
2. Task 1:
   - 1a and b: This task involved performing random and stratified sampling on the dataset
   - For stratified sampling we first need to apply kmeans clustering algorithm on the data to obtain clusters and then extract samples from individual clusters.
   - The code for this is snapshotted and pasted below:-

```python
19
20  def find_opt_k(data, max_k):
21      wcss = []
22      for i in range(1, max_k + 1):
23          kmeans = KMeans(n_clusters = i, init = 'k-means++', max_iter = 300, n_init = 10, random_state = 0)
24          kmeans.fit(data)
25          wcss.append(kmeans.inertia_)
26      plt.plot(range(1, max_k + 1), wcss)
27      plt.xlabel('no. of  clusters')
28      plt.show()
29
30  def stratified_sampling(data, no_clusters, frac):
31      kmeans1 = KMeans(n_clusters = no_clusters, random_state = 0)
32      result = kmeans1.fit(data)
33      data['titles'] = kmeans1.labels_
34      stratified_rows = []
35      for i in range(no_clusters):
36          length_of_labels = (int)(len(data[data['titles'] == i])*frac)
37          index_of_labels = list(data[data['titles'] == i].index)
38          rnd_sample = random.sample(index_of_labels, length_of_labels)
39          stratified_rows.append(data.loc[rnd_sample])
40
41      stratified_sample = pd.concat(stratified_rows)
42      del stratified_sample['titles']
43      return stratified_sample
44
```

   - The elbow curve looked something like this :-

Task 2:

2a) Dimension reduction

- Intrinsic Dimensionality of data is calculated using the PCA plot for both random and adaptive data samples. The Scree plot visualization is produced and intrinsic dimensionality is marked. It comes out to be 3 in both the cases as shown below the code snippet

```
61      return jsonify({ key :exp_var })
62
63   def perform_pca(data):
64       pca1 =  PCA()
65       data = pca1.fit_transform(data)
66       exp_var = pca1.explained_variance_
67       return exp_var
68
69   @app.route('/pca_org', methods = ['GET', 'POST'])
70   def pca_org():
71       # sample = stratified_sampling(dataset, 3, 0.25)
72       pca1 = PCA()
73       pca_data = pca1.fit_transform(dataset)
74       return jsonify({"key":pca_data})
75
76   @app.route('/pca_rnd', methods = ['GET', 'POST'])
77   def pca_rnd():
78       sample = random.sample(dataset, 0.25)
79       pca1 = PCA()
80       pca_data = pca1.fit_transform(sample)
81       return jsonify({"key":pca_data})
```

2c) top 3 highest loadings obtained:-

```python
def get_squared_loadings(dataframe, intrinsic):
    std_input = StandardScaler().fit_transform(dataframe)
    pca = PCA(n_components=intrinsic)
    pca.fit_transform(std_input)
    loadings = pca.components_
    # print("loadings shape ")
    # print(pd.DataFrame(loadings).shape)
    squared_loadings = []
    a = np.array(loadings)
    a = a.transpose()
    for i in range(len(a)):
        squared_loadings.append(np.sum(np.square(a[i])))
    # print(len(squared_loadings))
    # save squared_loadings in csv
    df_attributes = pd.DataFrame(pd.DataFrame(dataframe).columns)
    df_attributes.columns = ["attributes"]
    df_sqL = pd.DataFrame(squared_loadings)
    df_sqL.columns = ["squared_loadings"]
    sample = df_attributes.join([df_sqL])
    sample = sample.sort_values(["squared_loadings"], ascending=[False])
    sample.to_csv("./data/squared_loadings.csv", sep=',')
    return sample

def getTop3attributes(squared_loadings):
    top3 = squared_loadings.head(n = 3)
    return top3['attributes'].values.tolist()
    # arr = np.array(squared_loadings)
```

3a) This task involved getting the top 2 pca vectors and plotting a scatterplot, the code snippet for which is posted below:

```javascript
317
318    // scatterplot
319    function scatterplotfunc(){
320      $.post('/scatterplot', {'data': 'recvd'}, function(data_post){
321        data_post = data_post["key"];
322        scatterplot(data_post);
323      } );
324    }
325    // scatterplotfunc()
326    function scatterplot(data){
327      var margin = {top: 10, right: 30, bottom: 30, left: 60},
328        width = 460 - margin.left - margin.right,
329        height = 400 - margin.top - margin.bottom;
330
331      // append the svg object to the body of the page
332      var svg = d3.select("#my_dataviz")
333        .append("svg")
334          .attr("width", width + margin.left + margin.right)
335          .attr("height", height + margin.top + margin.bottom)
336        .append("g")
337          .attr("transform",
338                "translate(" + margin.left + "," + margin.top + ")");
339
340
341      // Add X axis
342      var x = d3.scaleLinear()
343        .domain([0, d3.max(data, function(d){return d["axis1"]})])
344        .range([ 0, width ]);
345      svg.append("g")
346        .attr("transform", "translate(0," + height + ")")
347        .call(d3.axisBottom(x));
348
349      // Add Y axis
```

3c) This task involved plotting a scatterplotmatrix for the top 3 highest loadings and the code snippet for this is attached bellow

```
04
05    traits.forEach(function(trait) {
06      domainByTrait[trait] = d3.extent(data, function(d) { return d[trait]; });
07    });
08
09    xAxis.tickSize(size * n);
10    yAxis.tickSize(-size * n);
11
12    var svg = d3.select("body").append("svg")
13        .attr("width", size * n + padding)
14        .attr("height", size * n + padding)
15      .append("g")
16        .attr("transform", "translate(" + padding + "," + padding / 2 + ")");
17
18    svg.selectAll(".x.axis")
19        .data(traits)
20      .enter().append("g")
21        .attr("class", "x axis")
22        .attr("transform", function(d, i) { return "translate(" + (n - i - 1) * size + ",0)"; })
23        .each(function(d) {
24          x.domain(domainByTrait[d]).nice();
25          d3.select(this).call(xAxis);
26        });
27
28    svg.selectAll(".y.axis")
29        .data(traits)
30      .enter().append("g")
31        .attr("class", "y axis")
32        .attr("transform", function(d, i) { return "translate(0," + i * size + ")"; })
33        .each(function(d) { y.domain(domainByTrait[d]); d3.select(this).call(yAxis); });
34
35    var cell = svg.selectAll(".cell")
36        .data(cross(traits, traits))
37      .enter().append("g")
38        .attr("class", "cell")
39        .attr("transform", function(d) { return "translate(" + (n - d.i - 1) * size + "," + d.j * size + ")"; })
40        .each(plot);
```