

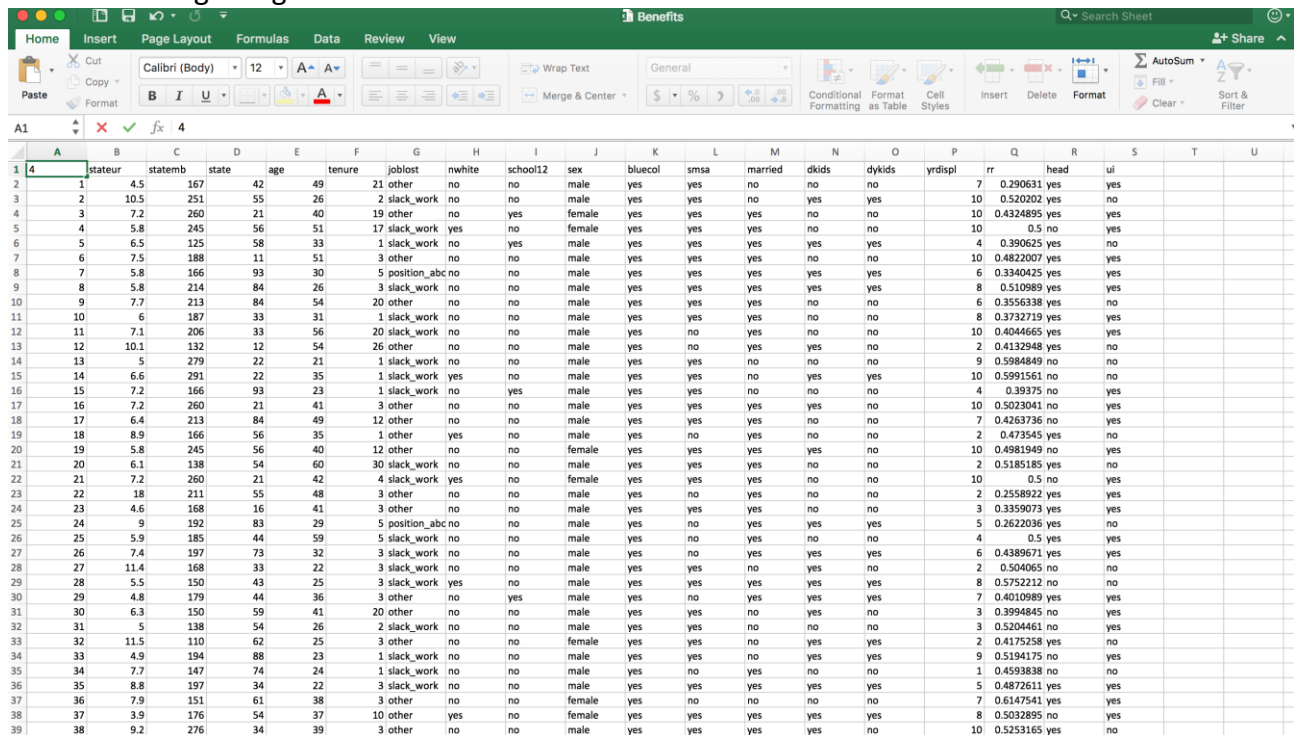
LAB 2

DATASET

Chosen Dataset: **Benefits, Unemployment of Blue Collar Workers**

Link to Dataset: <https://vincentarelbundock.github.io/Rdatasets/csv/Ecdat/benefits.csv>

Dataset contains 4877 rows & 18 dimensions. Some of them are numerical. I have changed the remaining categorical attributes into numerical attributes as well.



	A	B	C	D	E	F	G	H	I	J	K	L	M	N	O	P	Q	R	S	T	U
1	4	stature	stature	state	age	tenure	joblost	nwhite	school12	sex	bluecol	smsa	married	dkids	dykids	yrdispl	rr	head	ui		
2	1	4.5	167	42	49	21	other	no	no	male	yes	yes	no	no	no	7	0.290631	yes	yes		
3	2	10.5	251	55	26	2	slack_work	no	no	male	yes	yes	no	yes	yes	10	0.520202	yes	no		
4	3	7.2	260	21	40	19	other	no	yes	female	yes	yes	yes	no	no	10	0.4324895	yes	yes		
5	4	5.8	245	56	51	17	slack_work	yes	no	female	yes	yes	yes	no	no	10	0.5	no	yes		
6	5	6.5	125	58	33	1	slack_work	no	yes	male	yes	yes	yes	yes	yes	4	0.390625	yes	no		
7	6	7.5	188	11	51	3	other	no	no	male	yes	yes	yes	no	no	10	0.4822007	yes	yes		
8	7	5.8	166	93	30	5	position_abc	no	no	male	yes	yes	yes	yes	yes	6	0.3340425	yes	yes		
9	8	5.8	214	84	26	3	slack_work	no	no	male	yes	yes	yes	yes	yes	8	0.510989	yes	yes		
10	9	7.7	213	84	54	20	other	no	no	male	yes	yes	yes	no	no	6	0.3556338	yes	no		
11	10	6	187	33	31	1	slack_work	no	no	male	yes	yes	yes	no	no	8	0.3732719	yes	yes		
12	11	7.1	206	33	56	20	slack_work	no	no	male	yes	no	yes	no	no	10	0.4044665	yes	yes		
13	12	10.1	132	12	54	26	other	no	no	male	yes	no	yes	yes	no	2	0.4132948	yes	no		
14	13	5	279	22	21	1	slack_work	no	no	male	yes	yes	no	no	no	9	0.5984849	no	no		
15	14	6.6	291	22	35	1	slack_work	yes	no	male	yes	yes	no	yes	yes	10	0.5991561	no	no		
16	15	7.2	166	93	23	1	slack_work	no	yes	male	yes	yes	no	no	no	4	0.39375	no	yes		
17	16	7.2	260	21	41	3	other	no	no	male	yes	yes	yes	yes	no	10	0.5023041	no	yes		
18	17	6.4	213	84	49	12	other	no	no	male	yes	yes	yes	no	no	7	0.4263736	no	yes		
19	18	8.9	166	56	35	1	other	yes	no	male	yes	no	yes	no	no	2	0.473545	yes	no		
20	19	5.8	245	56	40	12	other	no	no	female	yes	yes	yes	yes	no	10	0.4981949	no	yes		
21	20	6.1	138	54	60	30	slack_work	no	no	male	yes	yes	yes	no	no	2	0.5185185	yes	no		
22	21	7.2	260	21	42	4	slack_work	yes	no	female	yes	yes	yes	no	no	10	0.5	no	yes		
23	22	18	211	55	48	3	other	no	no	male	yes	no	yes	no	no	2	0.2558922	yes	yes		
24	23	4.6	168	16	41	3	other	no	no	male	yes	yes	yes	no	no	3	0.3359073	yes	yes		
25	24	9	192	83	29	5	position_abc	no	no	male	yes	no	yes	yes	yes	5	0.2622036	yes	no		
26	25	5.9	185	44	59	5	slack_work	no	no	male	yes	no	yes	no	no	4	0.5	yes	yes		
27	26	7.4	197	73	32	3	slack_work	no	no	male	yes	no	yes	yes	yes	6	0.4389671	yes	yes		
28	27	11.4	168	33	22	3	slack_work	no	no	male	yes	yes	no	yes	no	2	0.504065	no	no		
29	28	5.5	150	43	25	3	slack_work	yes	no	male	yes	yes	yes	yes	yes	8	0.5752212	no	no		
30	29	4.8	179	44	36	3	other	no	yes	male	yes	no	yes	yes	yes	7	0.4010989	yes	yes		
31	30	6.3	150	59	41	20	other	no	no	male	yes	yes	no	yes	no	3	0.3994845	no	yes		
32	31	5	138	54	26	2	slack_work	no	no	male	yes	yes	no	no	no	3	0.5204461	no	yes		
33	32	11.5	110	62	25	3	other	no	no	female	yes	yes	no	yes	yes	2	0.4175258	yes	no		
34	33	4.9	194	88	23	1	slack_work	no	no	male	yes	yes	no	yes	yes	9	0.5194175	no	yes		
35	34	7.7	147	74	24	1	slack_work	no	no	male	yes	no	yes	no	no	1	0.4593838	no	no		
36	35	8.8	197	34	22	3	slack_work	no	no	male	yes	yes	yes	yes	yes	5	0.4872611	yes	yes		
37	36	7.9	151	61	38	3	other	no	no	female	yes	no	no	no	no	7	0.6147541	yes	yes		
38	37	3.9	176	54	37	10	other	yes	no	female	yes	yes	yes	yes	yes	8	0.5032895	no	yes		
39	38	9.2	276	34	39	3	other	no	no	male	yes	yes	yes	yes	no	10	0.5253165	yes	no		

I have used client server system: python for processing (server), D3 (v3) for VIS (client).

TASKS AND IMPLEMENTATION

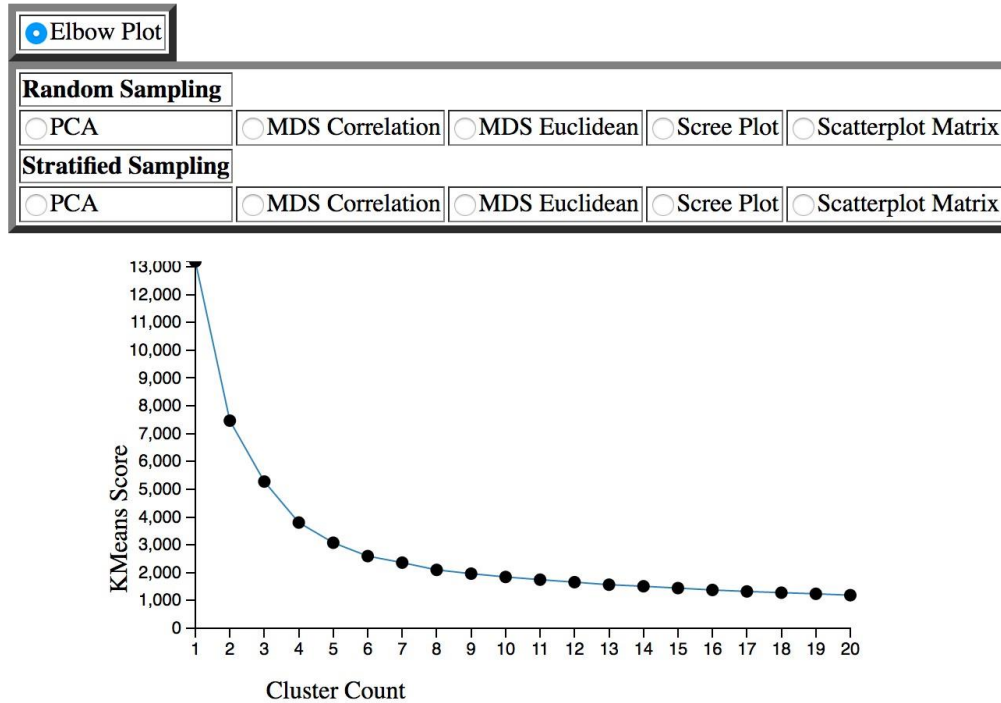
Task1: data clustering and decimation (30 points)

I used Python's numpy' random choice function to randomly select 20% for Random Sampling. For Stratified Sampling, I used K Means Clustering & constructed an elbow function to obtain the appropriate number of clusters, which turned out to be 4 and then selected 20% from each of the four clusters.

K means Elbow Plot

localhost:8000/Visualization-Project2-master/assignment2.html

Visualization Assignment 2



Code Snippets

```
def adaptive_sampling(data_frame, cluster_count, fraction):
    k_means = Kcluster.KMeans(n_clusters=cluster_count)
    print('k_means before', k_means)
    k_means.fit(data_frame)
    print('k_means after', k_means)
    data_frame['label'] = k_means.labels_
    print('label', k_means.labels_)
    adaptiveSampleRows = []
    for i in range(cluster_count):
        adaptiveSampleRows.append(data_frame.ix[random.sample(data_frame[data_frame['label'] == i].index, (int)
            (len(data_frame[data_frame['label'] == i])*fraction))])

    adaptiveSample = pd.concat(adaptiveSampleRows)
    del adaptiveSample['label']

    return adaptiveSample

def random_sampling(data_frame, fraction):
    rows = random.sample(data_frame.index, (int)(len(data_frame)*fraction))
    return data_frame.ix[rows]
```

```

def find_Elbow_Point(data_frame):
    sse = []
    maxK = 18
    minInertia = 10000000000000
    minK = 0
    kArr = list(xrange(21))
    kArr.pop(0)
    inertiaArr= []
    print kArr

    for k in range(1,21):
        # sse[k] = 0
        k_means = Kcluster.KMeans(n_clusters=k)
        k_means.fit(data_frame)
        clusters = k_means.labels_
        # print ('clusters', clusters)
        inertias = k_means.inertia_
        inertiaArr.append(inertias)
        if(minInertia>inertias):
            minInertia = inertias
            minK = k

        print ('inertia',inertias)
        # print('inertia',inertias/k)
    print ("minK ",minInertia,minK)

    # plt.plot( kArr, inertiaArr,c='blue',label="test1")
    # plt.xticks(np.arange(1, 21, 1.0))
    # plt.show()
    elbowDF = pd.DataFrame({'KMeans_Score' : inertiaArr})
    elbowDF['Cluster_Count'] = kArr
    elbowDF.to_csv('elbow.csv', sep=',')#adaptivescreeplt

```

Task 2: dimension reduction (use decimated data) (30 points)

Scree plots for both the random samples and stratified samples were obtained as shown below. The intrinsic dimensionality of the data is all points obtained where the Eigen Values were > 1 i.e. 4.

Scree Plots for Random and Stratified Samples

Now, we obtain the three highest attributes with highest PCA loadings and saved them for producing a scatter plot matrix later.

Loadings for Random Samples

← → ↻ ⓘ localhost:8000/Visualization-Project2-master/assignment2.html

Visualization Assignment 2

☐ Elbow Plot

Random Sampling

☐ PCA

☐ MDS Correlation

☐ MDS Euclidean

☒ Scree Plot

☐ Scatterplot Matrix

Stratified Sampling

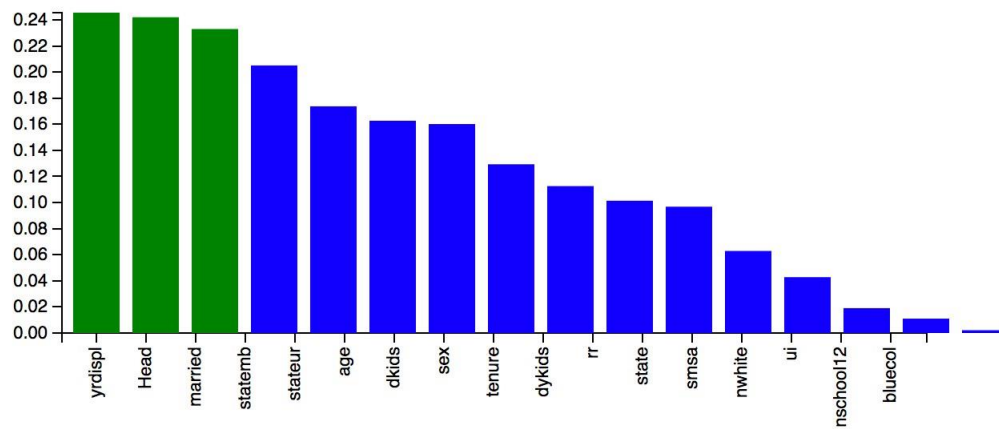
☐ PCA

☐ MDS Correlation

☐ MDS Euclidean

☐ Scree Plot

☐ Scatterplot Matrix



Loadings for Stratified Samples

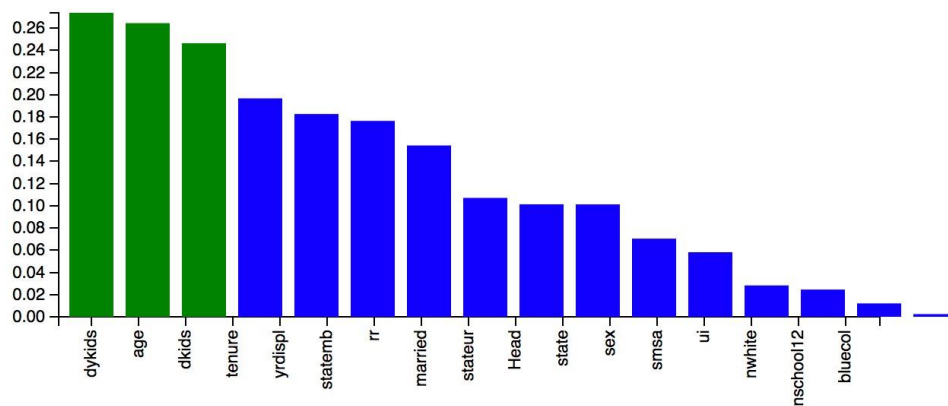
← → ↻ localhost:8000/Visualization-Project2-master/assignment2.html

Visualization Assignment 2

☐ Elbow Plot

Random Sampling
☐ PCA ☐ MDS Correlation ☐ MDS Euclidean ☐ Scree Plot ☐ Scatterplot Matrix

Stratified Sampling
☐ PCA ☐ MDS Correlation ☐ MDS Euclidean ☒ Scree Plot ☐ Scatterplot Matrix

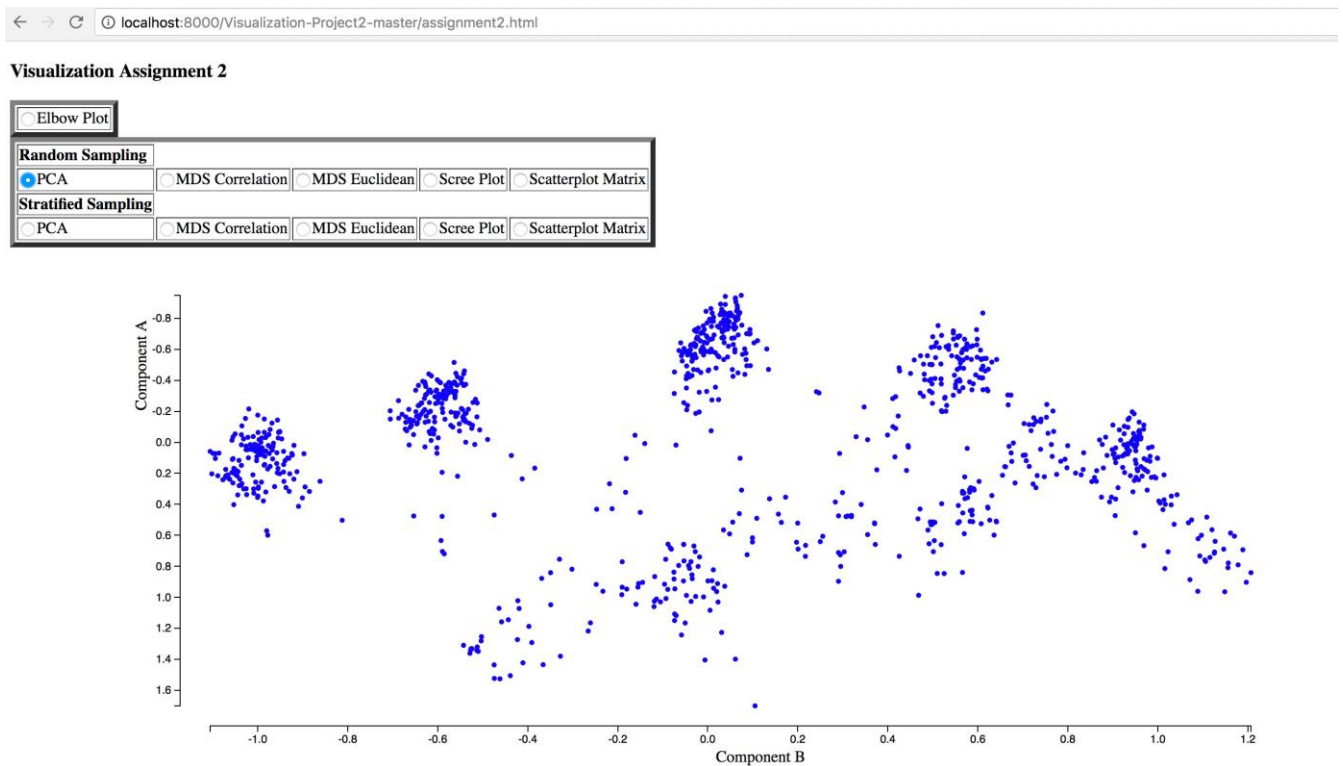


To obtain the three highest attributes with maximum PCA loading, data is first reduced to only the intrinsic dimensions and then I found the sum of squared loading value for all the attributes. Attributes with the highest sum of this squared value were considered as the three attributes with maximum PCA loadings.

Task 3: visualization (use dimension reduced data) (40 points)

Visualize data projected into the top two PCA vectors via 2D scatterplot

PCA (Random Sampling)

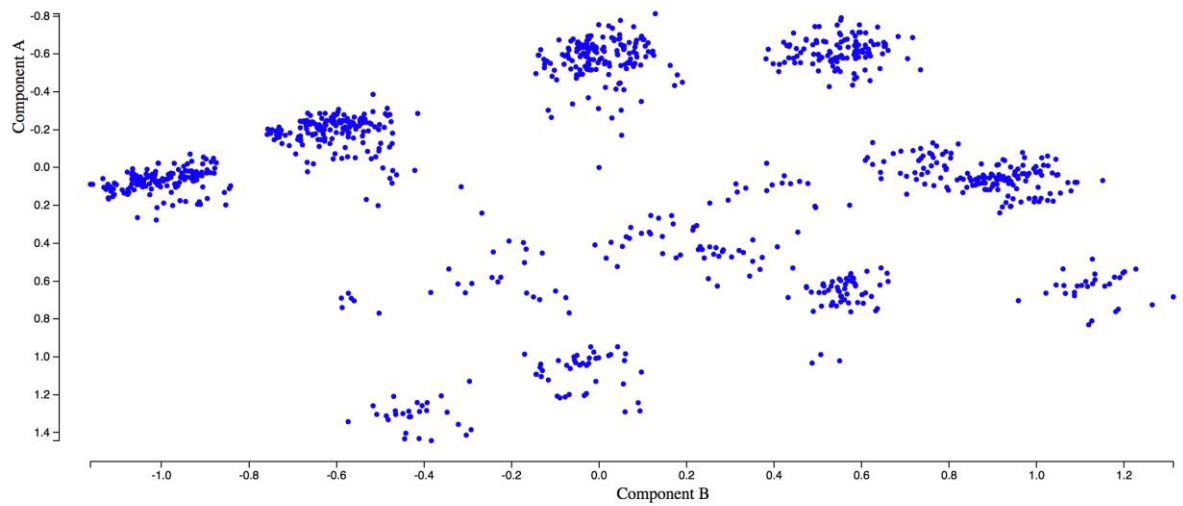


PCA (Stratified Sampling)

localhost:8000/Visualization-Project2-master/assignment2.html

Visualization Assignment 2

<input type="radio"/> Elbow Plot				
Random Sampling				
<input type="radio"/> PCA	<input type="radio"/> MDS Correlation	<input type="radio"/> MDS Euclidean	<input type="radio"/> Scree Plot	<input type="radio"/> Scatterplot Matrix
Stratified Sampling				
<input checked="" type="radio"/> PCA	<input type="radio"/> MDS Correlation	<input type="radio"/> MDS Euclidean	<input type="radio"/> Scree Plot	<input type="radio"/> Scatterplot Matrix



Visualize data via MDS (Euclidian & correlation distance) in 2D scatterplots

Euclidean

localhost:8000/Visualization-Project2-master/assignment2.html

Visualization Assignment 2

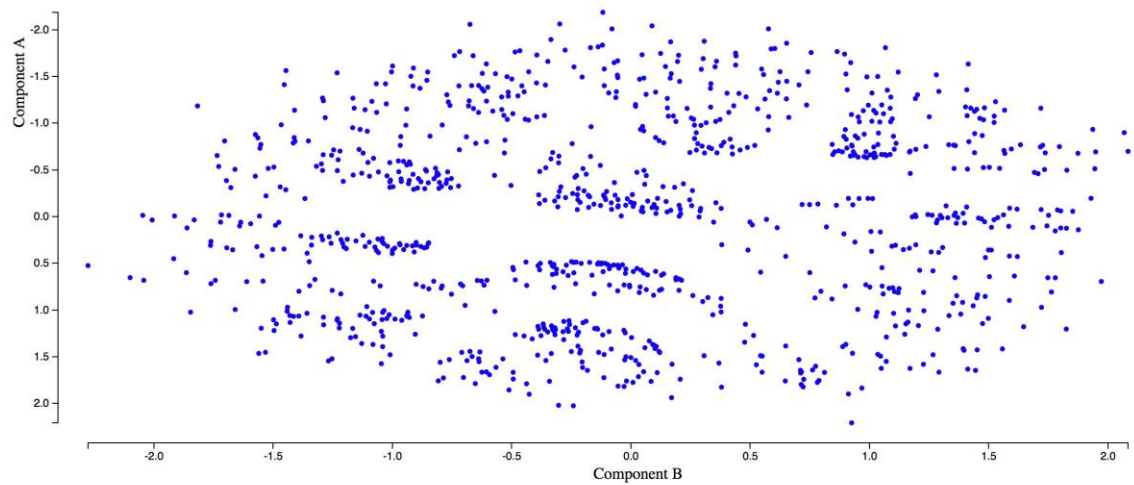
☐ Elbow Plot

Random Sampling

☐ PCA ☐ MDS Correlation ☒ MDS Euclidean ☐ Scree Plot ☐ Scatterplot Matrix

Stratified Sampling

☐ PCA ☐ MDS Correlation ☐ MDS Euclidean ☐ Scree Plot ☐ Scatterplot Matrix



Visualization Assignment 2

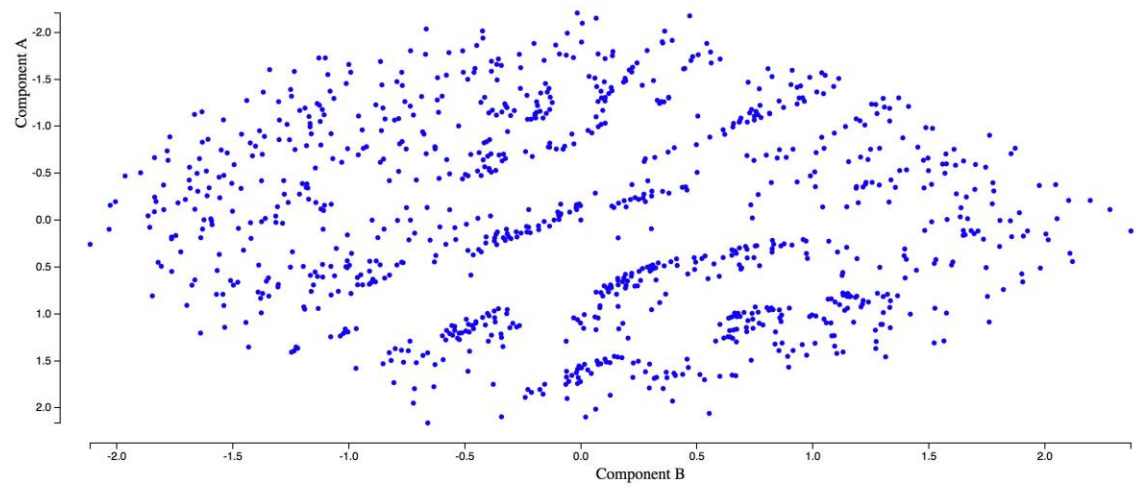
☐ Elbow Plot

Random Sampling

☐ PCA☐ MDS Correlation☐ MDS Euclidean☐ Scree Plot☐ Scatterplot Matrix

Stratified Sampling

☐ PCA☐ MDS Correlation☒ MDS Euclidean☐ Scree Plot☐ Scatterplot Matrix



Correlation

localhost:8000/Visualization-Project2-master/assignment2.html

Visualization Assignment 2

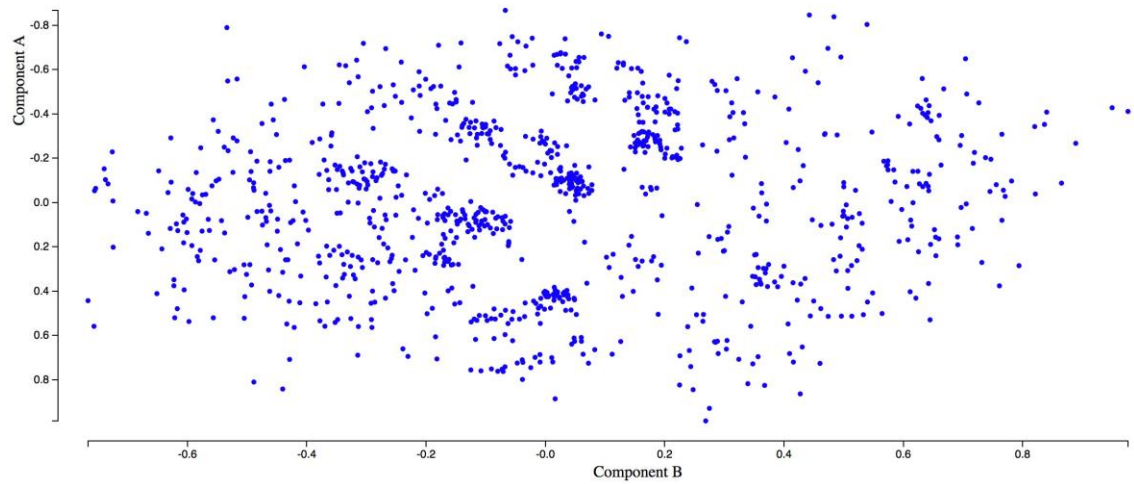
☐ Elbow Plot

Random Sampling

☐ PCA☒ MDS Correlation☐ MDS Euclidean☐ Scree Plot☐ Scatterplot Matrix

Stratified Sampling

☐ PCA☐ MDS Correlation☐ MDS Euclidean☐ Scree Plot☐ Scatterplot Matrix



localhost:8000/Visualization-Project2-master/assignment2.html

Visualization Assignment 2

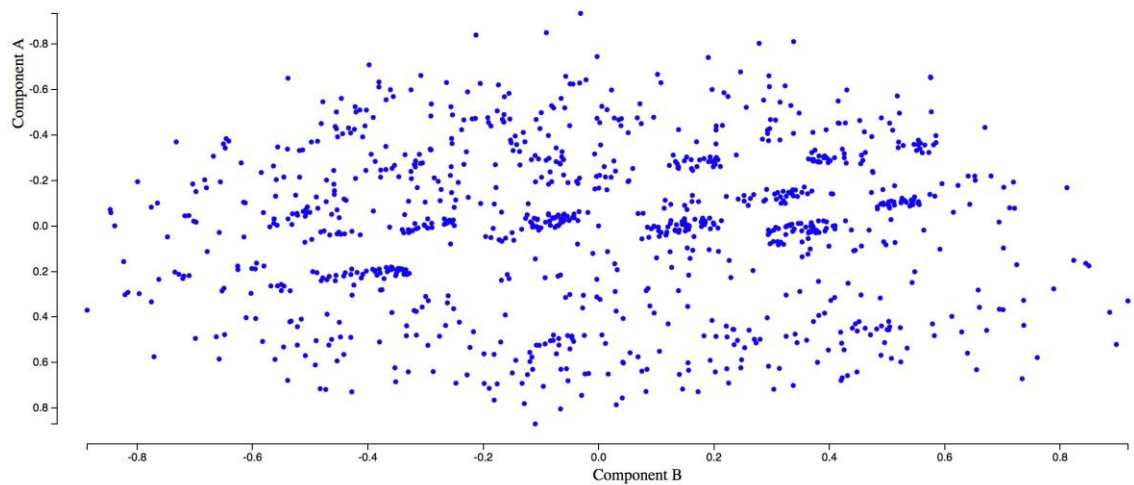
☐ Elbow Plot

Random Sampling

☐ PCA☐ MDS Correlation☐ MDS Euclidean☐ Scree Plot☐ Scatterplot Matrix

Stratified Sampling

☐ PCA☒ MDS Correlation☐ MDS Euclidean☐ Scree Plot☐ Scatterplot Matrix

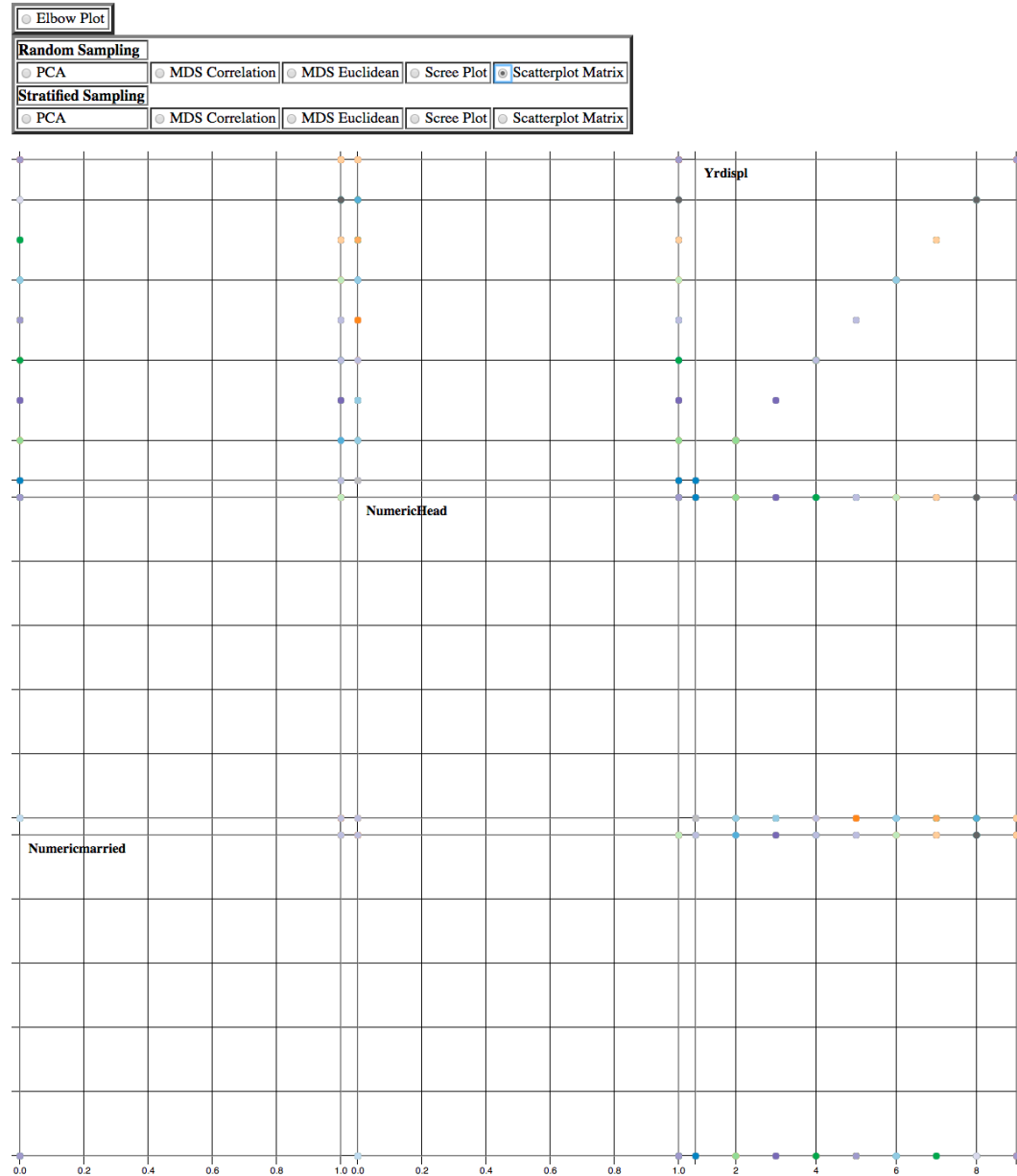


Visualize scatterplot matrix of the three highest PCA loaded attributes

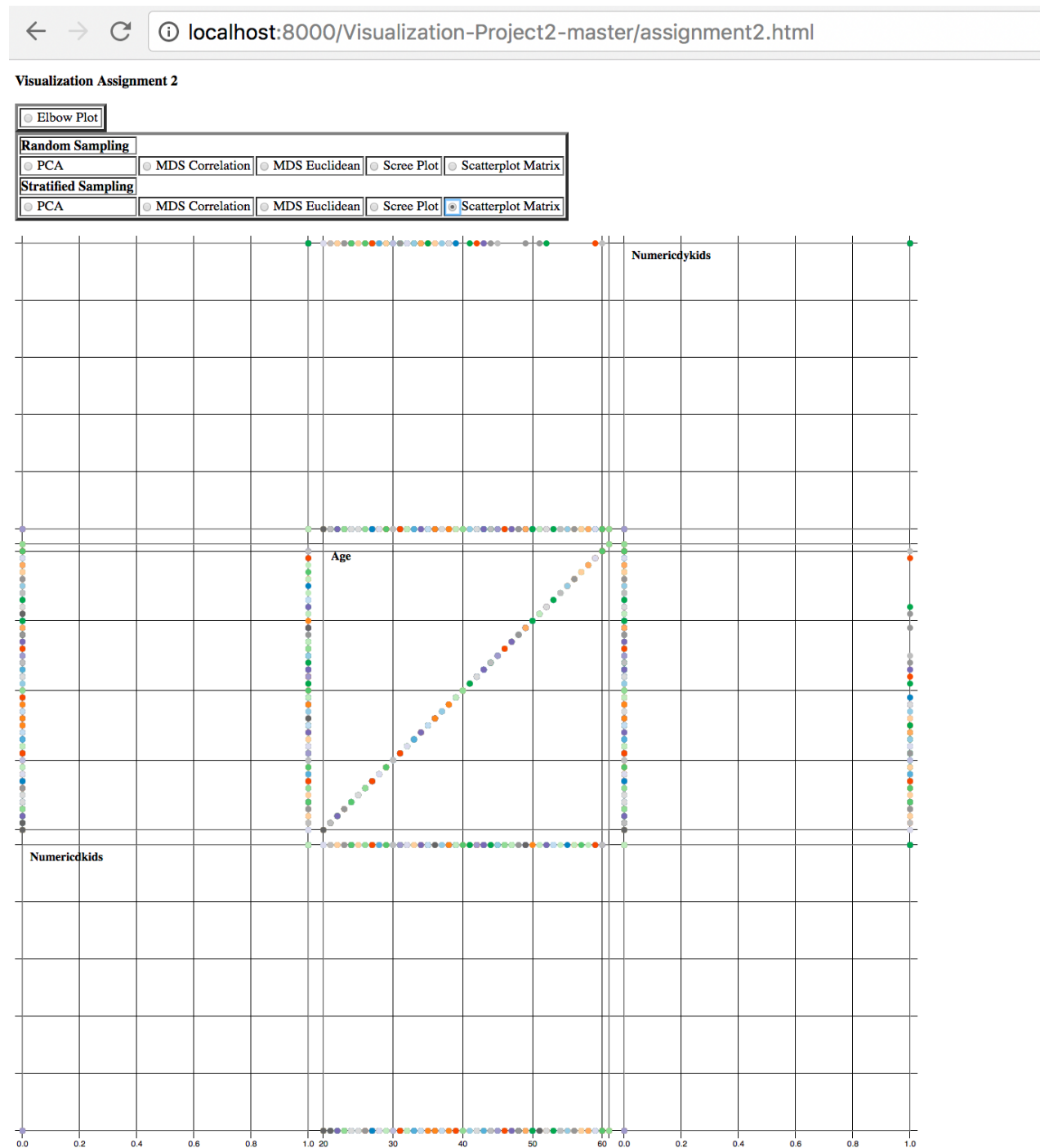
Random Sampling

← → ↻ ⓘ localhost:8000/Visualization-Project2-master/assignment2.html

Visualization Assignment 2



Stratified Sampling



CODE SNIPPETS

elbowplot.js

```
24 // Define the line
25 var valueline = d3.svg.line()
26   .y(function(d) { return y(d.KMeans_Score/1000); })
27   .x(function(d) { return x(d.Cluster_Count); });
28
29 // Get the data
30 d3.csv("./data/elbow.csv", function(error, data) {
31   data.forEach(function(d) {
32     d.KMeans_Score = +d.KMeans_Score;
33     d.Cluster_Count = +d.Cluster_Count;
34   });
35
36   // Scale the range of the data
37   y.domain([0,d3.max(data, function(d) { return d.KMeans_Score/1000; })]);
38   x.domain(d3.extent(data, function(d) { return d.Cluster_Count; }));
39
40   // Add the valueline path.
41   svg.append("path")
42     .attr("class", "line")
43     .attr("d", valueline(data))
44     .attr("transform", "translate(120,0)");
45
46   svg.selectAll("dot")
47     .data(data)
48     .enter().append("circle")
49     .attr("transform", "translate(120,0)")
50     .attr("r", 4)
51     .attr("cx", function(d) { return x(d.Cluster_Count); })
52     .attr("cy", function(d) { return y(d.KMeans_Score/1000); });
53
54   // Add the X Axis
55   svg.append("g")
56     .attr("class", "x axis")
57     .attr("transform", "translate(120," + height + ")")
58     .call(xAxis);
59
60   // Add the Y Axis
61   svg.append("g")
62     .attr("class", "y axis")
63     .attr("transform", "translate(120,0)")
64     .call(yAxis);
65
66   // Add the text label for the Y axis
67   svg.append("text")
68     .attr("transform", "rotate(-90)")
69     .attr("y", left_pad-140)
70     .attr("x", height-400)
71     .attr("dy", "1em")
72     .style("text-anchor", "middle")
73     .text("KMeans Score");
74
75   svg.append("text")
76     .attr("transform", "translate(" + (width / 2) + " ," + (height + margin.bottom-5) + ")")
77     .style("text-anchor", "middle")
78     .text("Cluster Count");
79
80   svg.append("circle")
81     .data(data)
82     .enter()
83     .append("circle")
84     .attr("r", 3)
85     .attr("cx", function(d){
86       return x(d.Cluster_Count);
87     })
88     .attr("cy", function(d){
89       return y(d.KMeans_Score/1000);
90     })
91   });
```

Scatterplot.js (Random sampling)

```
74 }
75
76 function plot_values(filename) {
77     filename = "../data/" + filename;
78     svg.selectAll("*").remove();
79
80     // Load data
81     d3.csv(filename, function(error, data) {
82         data.forEach(function(d) {
83             d.r1 = +d.r1;
84             d.r2 = +d.r2;
85
86         });
87
88         var xValueR = function(d) { return d.r1;};
89         var yValueR = function(d) { return d.r2;};
90
91         xScale.domain([d3.min(data, xValueR), d3.max(data, xValueR)]);
92         yScale.domain([d3.min(data, yValueR), d3.max(data, yValueR)]);
93
94         svg.append("g")
95             .attr("class", "axis")
96             .attr("transform", "translate(0, +(h-pad-10)+)")
97             .call(xAxis);
98
99         svg.append("g")
100             .attr("class", "axis")
101             .attr("transform", "translate(+(left-pad-pad)+, 0)")
102             .call(yAxis);
103
104         svg.append("text")
105             .attr("transform", "rotate(-90)")
106             .attr("y", left-pad-80)
107             .attr("x", h-600)
108             .attr("dy", "1em")
109             .style("text-anchor", "middle")
110             .text("Component A");
111
112         svg.append("text")
113             .attr("y", h-20)
114             .attr("x", h+250)
115             .attr("dy", "1em")
116             .style("text-anchor", "middle")
117             .text("Component B");
118
119         svg.selectAll("circle")
120             .data(data)
121             .enter()
122             .append("circle")
123             .attr("r", 2.5)
124             .attr("cx", function(d){
125                 return xScale(d.r1);
126             })
127             .attr("cy", function(d){
128                 return yScale(d.r2);
129             })
130             .style("fill", "blue");
131     });
132 }
133
134
135
136
```

Scatter_plot.js (Stratified Sampling)

```

74 }
75
76 function plot_values(filename) {
77
78     filename = "../data/" + filename;
79     svg.selectAll("*").remove();
80
81     // Load data
82     d3.csv(filename, function(error, data) {
83         data.forEach(function(d) {
84             d.r1 = +d.r1;
85             d.r2 = +d.r2;
86
87         });
88
89         var xValueR = function(d) { return d.r1;};
90         var yValueR = function(d) { return d.r2;};
91
92         xScale.domain([d3.min(data, xValueR), d3.max(data, xValueR)]);
93         yScale.domain([d3.min(data, yValueR), d3.max(data, yValueR)]);
94
95
96         svg.append("g")
97             .attr("class", "axis")
98             .attr("transform", "translate(0, " + (h-pad-10) + ")")
99             .call(xAxis);
100
101         svg.append("g")
102             .attr("class", "axis")
103             .attr("transform", "translate(" + (left-pad-pad) + ", 0)")
104             .call(yAxis);
105
106         svg.append("text")
107             .attr("transform", "rotate(-90)")
108             .attr("y", left-pad-80)
109             .attr("x", h-600)
110             .attr("dy", "1em")
111             .style("text-anchor", "middle")
112             .text("Component A");
113
114         svg.append("text")
115             .attr("y", h-20)
116             .attr("x", h+250)
117             .attr("dy", "1em")
118             .style("text-anchor", "middle")
119             .text("Component B");
120
121         svg.selectAll("circle")
122             .data(data)
123             .enter()
124             .append("circle")
125             .attr("r", 2.5)
126             .attr("cx", function(d){
127                 return xScale(d.r1);
128             })
129             .attr("cy", function(d){
130                 return yScale(d.r2);
131             })
132             .style("fill", "blue");
133     });
134 }
135
136

```

screplot.js


```

1 function scree_plot(filename) {
2
3     filename = "./data/" + filename;
4     svg.selectAll("*").remove();
5
6     // Set the dimensions of the canvas / graph
7     var margin = {top: 30, right: 20, bottom: 50, left: 50},
8         width = 800 - margin.left - margin.right,
9         height = 400 - margin.top - margin.bottom;
10
11     // Set the ranges
12     var x = d3.scale.linear().range([0, width]);
13     var y = d3.scale.linear().range([height, 0]);
14
15     // Define the axes
16     var xAxis = d3.svg.axis().scale(x)
17         .orient("bottom").ticks(18);
18
19     var yAxis = d3.svg.axis().scale(y)
20         .orient("left").ticks(7);
21
22     // Define the line
23     var valueline = d3.svg.line()
24         .y(function(d) { return y(d.eigen_value); })
25         .x(function(d) { return x(d.PCA_components); });
26
27     var valueline2 = d3.svg.line()
28         .y(function(d) { return y(1); })
29         .x(function(d) { return x(d.PCA_components); });
30
31
32     // Get the data
33     d3.csv(filename, function(error, data) {
34         data.forEach(function(d) {
35             d.eigen_value = +d.eigen_values;
36             d.PCA_components = +d.PCA_components;
37         });
38
39         // Scale the range of the data
40         y.domain([0, d3.max(data, function(d) { return d.eigen_value; }) * 1]);
41         x.domain(d3.extent(data, function(d) { return d.PCA_components; }));
42
43         // Add the valueline path.
44         svg.append("path")
45             .attr("class", "line")
46             .attr("d", valueline(data))
47             .attr("transform", "translate(40,0)");
48         // .attr("data-legend", function(d) { return d.name; });
49
50         svg.selectAll("dot")
51             .data(data)
52             .enter().append("circle")
53             .style("fill", function(d, i) { if(d.eigen_value > 1){return "blue";} else {return "green";}})
54             .attr("transform", "translate(40,0)")
55             .attr("r", 4)
56             .attr("cx", function(d) { return x(d.PCA_components); })
57             .attr("cy", function(d) { return y(d.eigen_value); });
58
59         svg.append("path")
60             .attr("class", "line")
61             .attr("transform", "translate(40,0)")
62             .attr("d", valueline2(data))
63             .attr("style", "stroke:red" );
64

```


Scatterplot_matrix.js

```
1 function scatterplot_matrix(filename) {
2   filename = "../data/" + filename;
3   svg.selectAll("*").remove();
4
5   var width = 960,
6       size = 400,
7       padding = 20;
8
9   var x = d3.scale.linear()
10      .range([padding / 2, size - padding / 2]);
11
12   var y = d3.scale.linear()
13      .range([size - padding / 2, padding / 2]);
14
15   var xAxis = d3.svg.axis()
16      .scale(x)
17      .orient("bottom")
18      .ticks(6);
19
20   var yAxis = d3.svg.axis()
21      .scale(y)
22      .orient("left")
23      .ticks(6);
24
25   var color = d3.scale.category20c();
26
27   d3.csv(filename, function(error, data) {
28     var domainByTrait = {},
29         traits = d3.keys(data[0]).filter(function(d) { return d !== "x"; }),
30         n = traits.length;
31
32     traits.forEach(function(trait) {
33       domainByTrait[trait] = d3.extent(data, function(d) { return d[trait]; });
34     });
35
36     xAxis.tickSize(size * n);
37     yAxis.tickSize(-size * n);
38
39     var brush = d3.svg.brush()
40        .x(x)
41        .y(y)
42        .on("brushstart", brushstart)
43        .on("brush", brushmove)
44        .on("brushend", brushend);
45
46     svg.attr("width", size * n + padding)
47        .attr("height", size * n + padding)
48        .append("g")
49        .attr("transform", "translate(" + padding + "," + padding / 2 + ")");
50
```

```

71
72 // Titles for the diagonal.
73 cell.filter(function(d) { return d.i === d.j; }).append("text")
74   .attr("x", padding)
75   .attr("y", padding)
76   .attr("dy", ".71em")
77   .text(function(d) { return d.x; });
78
79 cell.call(brush);
80
81 function plot(p) {
82   var cell = d3.select(this);
83
84   x.domain(domainByTrait[p.x]);
85   y.domain(domainByTrait[p.y]);
86
87   cell.append("rect")
88     .attr("class", "frame")
89     .attr("x", padding / 2)
90     .attr("y", padding / 2)
91     .attr("width", size - padding)
92     .attr("height", size - padding);
93
94   cell.selectAll("circle")
95     .data(data)
96     .enter().append("circle")
97     .attr("cx", function(d) { return x(d[p.x]); })
98     .attr("cy", function(d) { return y(d[p.y]); })
99     .attr("r", 4)
100    .style("fill", function(d,i) { return color(i); });
101  }
102
103  var brushCell;
104
105  // Clear the previously-active brush, if any.
106  function brushstart(p) {
107    if (brushCell !== this) {
108      d3.select(brushCell).call(brush.clear());
109      x.domain(domainByTrait[p.x]);
110      y.domain(domainByTrait[p.y]);
111      brushCell = this;
112    }
113  }
114
115  // Highlight the selected circles.
116  function brushmove(p) {
117    var e = brush.extent();
118    svg.selectAll("circle").classed("hidden", function(d) {
119      return e[0][0] > d[p.x] || d[p.x] > e[1][0]
120        || e[0][1] > d[p.y] || d[p.y] > e[1][1];
121    });
122  }
123
124  // If the brush is empty, select all circles.
125  function brushend() {
126    if (brush.empty()) svg.selectAll(".hidden").classed("hidden", false);
127  }
128
129  function cross(a, b) {
130    var c = [], n = a.length, m = b.length, i, j;
131    for (i = -1; ++i < n; ) for (j = -1; ++j < m; ) c.push({x: a[i], i: i, y: b[j], j: j});
132    return c;
133  }

```

```

def determine_pca(data_frame,type):
    data_frame = data_frame.as_matrix()
    kArr = list(xrange(17))
    # kArr.pop(0)
    pca = PCA(n_components=2)
    print (pca)
    print "DATA",data_frame
    pca.fit_transform(data_frame)
    print "Eigen vectors ",pca.components_

    print "Eigens ",len(pca.components_)

    print "After performing eigen"

    temp = []
    for i in range(0,17):
        temp.append(pca.components_[0][i] * pca.components_[0][i] + pca.components_[1][i] * pca.components_[1][i]
        ])

    print "temp ",temp,len(temp)
    objects = ('stateur','statemb','state','age','tenure','yrdispl','rr','NumericHead','Numericwhite',
        'Numericschool12','Numericsex','Numericbluecol','Numericsmsa','Numericmarried','Numericdkids',
        'Numericdykids','Numericui')
    screeDataFrame = pd.DataFrame({'PCA_components' : objects})
    screeDataFrame['eigen_value'] = temp
    # sample.to_csv(file_name, sep=',')
    if type == 1:
        screeDataFrame.to_csv('randomscreeplt.csv', sep=',')
    else:
        screeDataFrame.to_csv('adaptivescreeplt.csv', sep=',')#adaptivescreeplt
    print "variance ",pca.explained_variance_
    return pd.DataFrame(pca.fit_transform(data_frame))

```

```

def find_MDS(dataframe, type):
    dis_mat = SK_Metrics.pairwise_distances(dataframe, metric = type)
    mds = MDS(n_components=2, dissimilarity='precomputed')
    return pd.DataFrame(mds.fit_transform(dis_mat))

data_directory = "/Users/karanmalhotra/Downloads/"
def create_File(random_sample, adaptive_sample, file_name):
    # print "random ",random_sample
    # print "adaptive ",adaptive_sample
    # random_sample['type'] = 1
    # adaptive_sample['type'] = 2
    random_sample.columns = ["r1","r2"]
    adaptive_sample.columns = ["a1","a2"]

    sample = random_sample.join([adaptive_sample])

    file_name = data_directory + file_name
    sample.to_csv(file_name, sep=',')

def calculate_values(random_sample, adaptive_sample,function,file_name):
    create_File(function(random_sample,1), function(adaptive_sample,2),file_name + ".csv")

def Apply_Normalization(dataFrame):
    # dataFrame = dataFrame.apply(zscore)
    # return dataFrame
    min_max_scaler = preprocessing.StandardScaler()
    np_scaled = min_max_scaler.fit_transform(dataFrame)
    df_normalized = pd.DataFrame(np_scaled)
    return df_normalized

def main():
    benefitDF = pd.read_csv("/Users/karanmalhotra/Downloads/Benefits.csv")
    precision = 3
    benefitDF['rr'] = benefitDF['rr'].round(decimals = precision)
    # print benefitDF['rr']

```

The complete code and the video can be seen in the attachments of submission.

Interesting observations about data

The dataset had 18 attributes, out of which only 4 has eigen value 1 or higher. The intrinsic dimensionality of the data was 4 for a dataset of 18 attributes. The data of some attributes varies largely as compared to data of other attributes, thus it becomes essential to normalize them so that one attribute does not completely overwhelm the other.

Instead of using random and stratified sampling, reservoir or inversion sampling could also be used. Reservoir sampling is more appropriate when we have streaming data. Since we have a fixed dataset, stratified sampling seems to be the most appropriate sampling technique which also gives us a uniform sample data spread across all the clusters. Similarly, instead of using PCA, we could also have used techniques like LDA, GDA, Isomap etc.

REFERENCES

- 1.) Scatter Plot reference:
<https://bl.ocks.org/mbostock/3887118>
- 2.) Scatter Plot Matrix reference
<https://bl.ocks.org/mbostock/3213173>
- 3.) Bar Chart Reference
<https://bl.ocks.org/mbostock/3885304>