# EECS E6893 Big Data Analytic HW4

Chong Hu ch3467

November 15, 2019

**Problem 1.** Bar Chart

1. Answer these questions in simple words.

1.1 SVG Coordinate Space basically works in the way of Mathematical / Graph Coordinate Space, except for:

SVG Coordinate space has x=0 and y=0 coordinates fall on the top left, while Mathematical / Graph Coordinate Space has x=0 and y=0 fall on the bottom left.

SVG Coordinate space has the Y coordinate growing from top to bottom.

1.2 `.enter` identifies any DOM elements that need to be added when the joined array is longer than the selection. `.exit` returns an exit selection which consists of the elements that need to be removed from the DOM.

1.3 `transform` in SVG is used to transform a single SVG shape element or group of SVG elements. SVG transform supports `translate, scale, rotate` and `skew`. `translation` moves all the points of an element in the same direction and by the same amount. It preserves parallelism, angles and distances.

1.4 Anonymous function in d3.js is used to easily manipulate data and data could be applied in DOM. In this case, the answer is
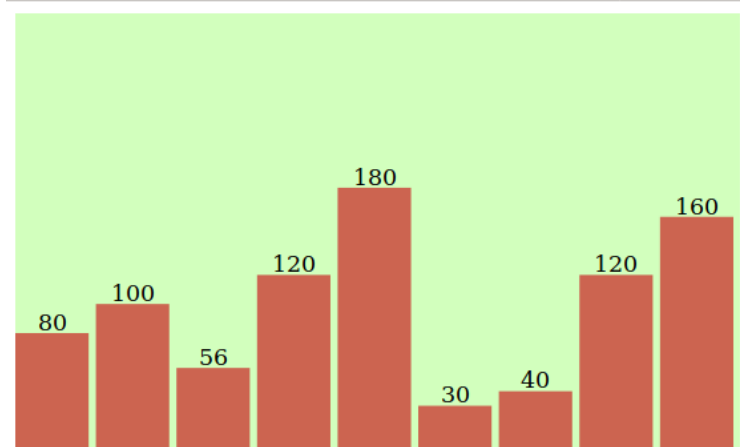
```
[5,6,7,8,9]
```

2. bar-chart



Figure 1: Bar Chart plot using SVG.

Code: Here I write the javascript in a single file (`part1.js`), separated with the structure file `part1.html`

`part1.html`

```html
<!DOCTYPE html>
<html lang="en">
<head>
    <meta charset="UTF-8">
    <title>Homework 4 Question 1</title>
</head>

<body>
<style>body {
}

.bar-chart {
    background-color: #D2FFBD;
}</style>
<svg class="bar-chart"></svg>

<script src="https://d3js.org/d3.v5.min.js"></script>
<script src="part1.js"></script>

</body>
</html>
```

part1.js

```
1   var data = [80, 100, 56, 120, 180, 30, 40, 120, 160];
2   var svgWidth = 500, svgHeight = 300;
3   // The required padding between bars is 5px.
4   // The label must locate 2px above the middle of each bar.
5
6   var svg = d3.select('svg')
7       .attr("width", svgWidth)
8       .attr("height", svgHeight);
9
10  const barPadding = 5;
11  const barWidth = svgWidth / data.length - barPadding;
12
13  function translateBarHelper(d, i) {
14      return "translate(" + (barWidth + barPadding) * i + ","
15          + (svgHeight - d) + ")";
16  }
17
18  function translateTextHelper(d, i) {
19      return "translate(" + ((barWidth + barPadding) * i + barWidth / 2) + ","
20          + (svgHeight - d - 2) + ")";
21  }
22
23  var barChart = svg.selectAll("rect")
24      .data(data)
25      .enter();
26
27  barChart.append("rect")
28      .attr("class", "bar")
29      .attr("height", function (d) {
30          return d;
31      })
32      .attr("width", barWidth)
33      .attr("transform", translateBarHelper)
34      .attr("fill", "#CC6450");
35
36  barChart.append("text")
37      .text(function (d) {
38          return d;
39      }).attr("transform", translateTextHelper)
40      .style("text-anchor", "middle");
```

**Problem 2.** Dashboard

1. Data processing.

   Result preview:



Figure 2: Preview of the table in the BigQuery

Code of data processing: Here I directly use SQL in BigQuery to process the data. Here is the SQL query.

```
1  CREATE VIEW  IF NOT EXISTS twitter_analysis.data as
2  SELECT time, count as data
3  FROM `hardy-symbol-252200.twitter_analysis.wordcount`
4  WHERE word="data";
5
6  CREATE VIEW IF NOT EXISTS twitter_analysis.ai as
7  SELECT time, count as ai
8  FROM `hardy-symbol-252200.twitter_analysis.wordcount`
9  WHERE word="ai";
10
11 CREATE VIEW  IF NOT EXISTS twitter_analysis.good as
12 SELECT time, count as good
13 FROM `hardy-symbol-252200.twitter_analysis.wordcount`
14 WHERE word="good";
15
16 CREATE VIEW IF NOT EXISTS twitter_analysis.movie as
17 SELECT time, count as movie
18 FROM `hardy-symbol-252200.twitter_analysis.wordcount`
19 WHERE word="movie";
```

```
20
21   CREATE VIEW  IF NOT EXISTS twitter_analysis.spark as
22   SELECT time, count as spark
23   FROM `hardy-symbol-252200.twitter_analysis.wordcount`
24   WHERE word="spark";
25
26   CREATE TABLE IF NOT EXISTS twitter_analysis.rstcnt AS
27   (SELECT COALESCE(t1.time, t2.time) as time, IFNULL(data, 0) as data,
28   IFNULL(ai, 0) as ai, IFNULL(good, 0) as good, IFNULL(movie, 0) as movie,
29   IFNULL(spark, 0) as spark
30   FROM
31   (SELECT COALESCE(t1.time, t2.time) as time, IFNULL(data, 0) as data,
32   IFNULL(ai, 0) as ai, IFNULL(good, 0) as good, IFNULL(movie, 0) as movie
33   FROM
34   (SELECT COALESCE(t1.time, t2.time) as time, IFNULL(data, 0) as data,
35   IFNULL(ai, 0) as ai
36   FROM twitter_analysis.data t1
37   FULL OUTER JOIN
38   twitter_analysis.ai t2
39   ON t1.time = t2.time)  t1
40   FULL OUTER JOIN
41   (SELECT COALESCE(t1.time, t2.time) as time, IFNULL(good, 0) as good,
42   IFNULL(movie, 0) as movie
43   FROM twitter_analysis.good t1
44   FULL OUTER JOIN
45   twitter_analysis.movie t2
46   ON t1.time = t2.time) t2
47   ON t1.time = t2.time) t1
48   FULL OUTER JOIN
49   twitter_analysis.spark t2
50   ON t1.time = t2.time);
```

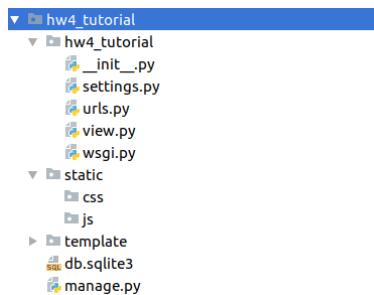2. Create Django project.

Screenshot of Directory Structure:



Figure 3: Screenshot of Directory Structure
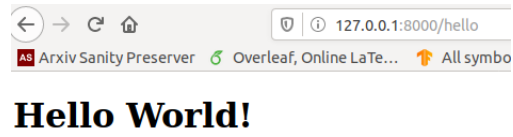
Screenshot to show `helloworld` page:



Figure 4: Screenshot of `helloworld` page
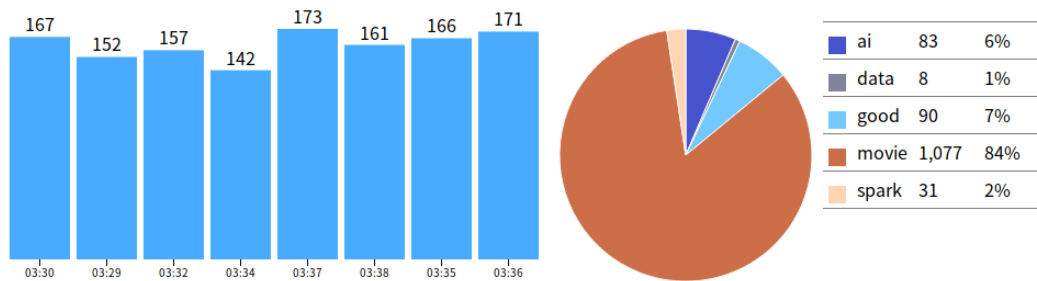
3. Finish the code. Output result:
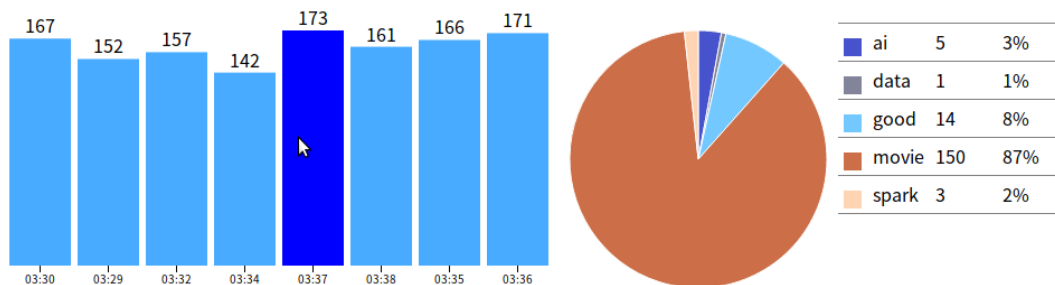


Figure 5: Screenshot of Dashboard.



Figure 6: Screenshot of Dashboard when mouse on the bar.
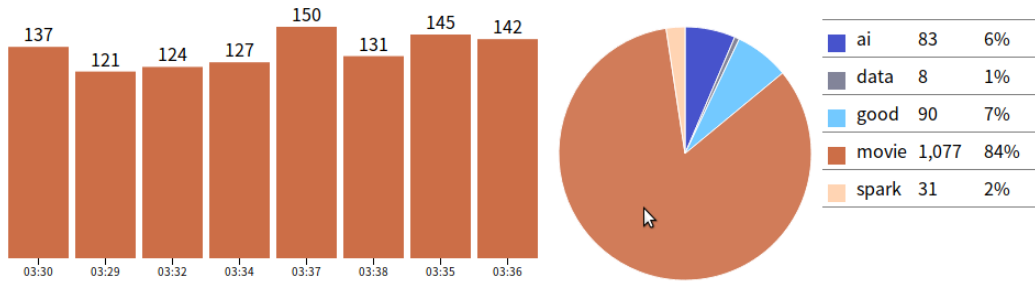
## Question 2 - Dashboard



Figure 7: Screenshot of Dashboard when mouse on the Pie chart.

Code:

Inside `view.py`:

```
1  credentials = service_account.Credentials.from_service_account_file(
2      '/home/huchong/Downloads/hardy-symbol-252200-25dbece318bf.json')
3
4
5  def dashboard(request):
6      pandas_gbq.context.credentials = credentials
7      pandas_gbq.context.project = "hardy-symbol-252200"
8
9      SQL = "SELECT time, ai, data, good, movie, spark " \
10         "FROM `hardy-symbol-252200.twitter_analysis.rstcnt` " \
11         "LIMIT 8"
12     df = pandas_gbq.read_gbq(SQL)
13     df_list = df.to_dict('records')
14
15     data_list = []
16     for df_row in df_list:
17         data_row = dict()
18         data_row["Time"] = df_row["time"].strftime(format="%H:%M")
19         df_row = dict(df_row)
20         df_row.pop("time")
21         data_row["count"] = df_row
22         data_list.append(data_row)
23
24     data = dict()
25     data["data"] = data_list
26
```

```python
27        '''
28            TODO: Finish the SQL to query the data, it should be limited to 8 rows.
29            Then process them to format below:
30            Format of data:
31            {'data': [{'Time': hour:min, 'count': {'ai': xxx, 'data': xxx, 'good': xxx, 'movie': xxx,
32                      {'Time': hour:min, 'count': {'ai': xxx, 'data': xxx, 'good': xxx, 'movie': xxx,
33                      ...
34                      ]
35            }
36        '''
37
38        return render(request, 'dashboard.html', data)
```

Inside `dashboard.js`

```javascript
1      // Define the color to change if your mouse move on the bar
2      var barColor = '#49abff';
3
4      // Choose color for each word:
5      function segColor(c) {
6          cmap = {
7              ai: "#4753CC",
8              data: "#828499",
9              good: "#73C9FF",
10             movie: "#CC6E47",
11             spark: '#FFD4B3'
12         };
13         /* TO FINISH */
14         return cmap[c];
15     }
```

```javascript
1      // compute total for each state.
2      fData.forEach(function (d) {
3          d.total = NaN;
4          /* TO FINISH */
5          d.total = Object.keys(d.count)
6              .reduce((sum, key) => sum + d.count[key], 0);
7      });
```

```javascript
1      bars.append("rect")
2      /* TO FINISH */
```

```
3        .attr("x", function (d, i) {
4            return x(d[0]);
5        })
6        /* TO FINISH */
7        .attr("y", function (d) {
8            return y(d[1]);
9        })
10       .attr("width", x.rangeBand())
11       .attr("height", function (d) {
12           return hGDim.h - y(d[1]);
13       })
14       ...
```

```
1        //Create the frequency labels ABOVE the rectangles.
2        bars.append("text").text(function (d) {
3            return d3.format(",")(d[1])
4        })
5            .attr("x", function (d) {
6                /* TO FINISH */
7                return x(d[0]) + x.rangeBand() / 2;
8            })
9            /* TO FINISH */
10           .attr("y", function (d) {
11               return y(d[1]) - 5;
12           })
13           .attr("text-anchor", "middle");
```

```
1        // transition the height and color of rectangles.
2        bars.select("rect").transition().duration(500)
3        /* TO FINISH */
4            .attr("y", function (d) {
5                return y(d[1]);
6            })
7            ...
```

```
1        // calculate total count by segment for all state.
2        var tF = ['ai', 'data', 'good', 'movie', 'spark'].map(function (d) {
3            return {
4                type: d, count: d3.sum(fData.map(function (t) {
5                    /* TO FINISH */
6                    return t.count[d];
```

```
7              }))
8          };
9      });
```

## Problem 3. Connection

1. Data processing.

   Here I process the data in local and save required nodes and edges into `nodes.csv` and `edges.csv` separately, and then upload to `BigQuery`. Here is the python script to deal with data, including both nodes and edges data.

```python
1  import csv
2  import os
3
4  from graphframes import *
5  from pyspark import SQLContext
6  from pyspark import SparkConf, SparkContext
7
8
9  def getData(sc, filename):
10     """
11     Load data from raw text file into RDD and transform.
12     Hint: transfromation you will use: map(<lambda function>).
13     Args:
14         sc (SparkContext): spark context.
15         filename (string): hw2.txt cloud storage URI.
16     Returns:
17         RDD: RDD list of tuple of (<User>, [friend1, friend2, ... ]),
18         each user and a list of user's friends
19     """
20     # read text file into RDD
21     data = sc.textFile(filename)
22     data = data.map(lambda line: line.split("\t")).map(
23         lambda line: (int(line[0]), [int(x) for x in line[1].split(",")] if len(
24             line[1]) else []))
25     return data
26
27
28  def get_vertices(data, sqlcontext):
29     """
30     get vertices
31     :param data: RDD list of tuple of (<User>, [friend1, friend2, ... ]),
32         each user and a list of user's friends
33     :param sqlcontext: SQLContext
```

```python
34          :return: dataframe
35          """
36          vertices = data.map(lambda line: (line[0],))
37
38          return sqlcontext.createDataFrame(vertices, schema=["id"])
39
40
41  def get_edges(data, sqlcontext):
42          """
43          get edges
44          :param data: RDD list of tuple of (<User>, [friend1, friend2, ... ]),
45              each user and a list of user's friends
46          :param sqlcontext: SQLContext
47          :return:
48          """
49
50          def map_friends(line):
51              """
52              map function to construct edge between friends
53              construct a pair of ((friend1, friend2) -> common friends list)
54              if two friends are already direct friends, then common friends list
55              is empty.
56              :param line: tuple of (<User>, [friend1, friend2, ... ]),
57                          each user and a list of user's friends
58              :return: friend pair
59              """
60              user = line[0]
61              friends = line[1]
62              for i in range(len(friends)):
63                  yield (user, friends[i])
64
65          edges = data.flatMap(map_friends)
66          return sqlcontext.createDataFrame(edges, schema=["src", "dst"])
67
68
69  def save_nodes(nodes):
70          """
71          save node list to csv
72          :param nodes: list that contain nodes in the cluster of 25 users
73          :return:
74          """
75          with open("nodes.csv", "w") as csv_file:
76              csv_writer = csv.writer(csv_file)
77              csv_writer.writerow(["node"])
78              csv_writer.writerows(nodes)
```

```python
79

80

81   def connected_components(graph):
82       """
83       run connected components on graph
84       :param graph: Graph contains vertices and edges
85       :return:
86       """
87       print("connected components")
88       result = graph.connectedComponents(algorithm="graphx")

89

90       # How many clusters / connected components in total for this dataset
91       cluster_num = result.select("component").distinct().count()
92       print("clusters amount: ", cluster_num)
93       print()

94

95       # How many users in the top 10 clusters?
96       print("number of users in top 10 cluster")
97       res1 = result.groupBy("component").count().orderBy('count',
98                                                          ascending=False)
99       res2 = res1.head(10)
100      total = 0
101      for row in res2:
102          total += row["count"]
103          print("cluster id:\t%d\tnumber of users:\t%d" % (
104              row["component"], row["count"]))
105      print("Total number of users in top 10 cluster:\t", total)
106      print()

107

108      # What are the user ids for the cluster which has 25 users?
109      print("user ids for the cluster which has 25 users")
110      cluster_id = res1.where(res1["count"] == 25).select("component").collect()
111      cluster_id = [row["component"] for row in cluster_id]
112      user_list = result.where(result["component"].isin(cluster_id)).select(
113          "id").collect()
114      user_ls = [row["id"] for row in user_list]
115      user_ls.sort()
116      save_nodes([[node] for node in user_ls])
117      print(user_ls)
118      print()

119

120      # get edges for 25 nodes
121      df_edges = graph.edges.filter(
122          graph.edges.dst.isin(user_ls) & graph.edges.src.isin(user_ls))
123      df_edges = df_edges.rdd.map(
```

```python
124             lambda x: (user_ls.index(x[0]), user_ls.index(x[1]))).toDF(
125             ["source", "target"])
126         # write edges to csv
127         df_edges.toPandas().to_csv("edges.csv", header=True, index=False)
128         return


131 def main():
132     # Configure Spark
133     if not os.path.isdir("checkpoints"):
134         os.mkdir("checkpoints")
135     conf = SparkConf().setMaster('local').setAppName('connected components')
136     sc = SparkContext(conf=conf)
137     sqlcontext = SQLContext(sc)
138     SparkContext.setCheckpointDir(sc, "checkpoints")

140     # The directory for the file
141     filename = "q1.txt"

143     # Get data in proper format
144     data = getData(sc, filename)
145     edges = get_edges(data, sqlcontext)
146     vertices = get_vertices(data, sqlcontext)
147     graph = GraphFrame(vertices, edges)
148     connected_components(graph=graph)


151 if __name__ == '__main__':
152     main()
```

Previews in the `BigQuery`.



Figure 9: Screenshot of preview of `nodes` in `BigQuery`



Figure 10: Screenshot of preview of `edges` in `BigQuery`

2. Finish the code.

   Inside `view.py`: Notice that possible duplicates of edges' data are already removed in the data processing step and the amount of edges is correct.

```python
credentials = service_account.Credentials.from_service_account_file(
    '/home/huchong/Downloads/hardy-symbol-252200-25dbece318bf.json')


def connection(request):
    pandas_gbq.context.credentials = credentials
    pandas_gbq.context.project = "hardy-symbol-252200"
    SQL1 = 'SELECT node ' \
           'FROM `hardy-symbol-252200.graph.nodes`'
    df1 = pandas_gbq.read_gbq(SQL1)

    SQL2 = 'SELECT source, target ' \
           'FROM `hardy-symbol-252200.graph.edges`'
    df2 = pandas_gbq.read_gbq(SQL2)

    data = {
        'n': list(df1.T.to_dict().values()),
```

```
18          'e': list(df2.T.to_dict().values())
19      }
20
21      '''
22          TODO: Finish the SQL to query the data, it should be limited to 8 rows.
23          Then process them to format below:
24          Format of data:
25          {
26          'n': [{'node': 18233},{'node': 18234},...]
27          'e': [{'source':0, 'target':0},{'source':0, 'target':1},... ]
28          }
29      '''
30
31      return render(request, 'connection.html', data)
```

Inside `connection.js`:

```
1      var svg = d3.select("body")
2          .append("svg")
3          /* TO FINISH */
4          .attr("height", height)
5          /* TO FINISH */
6          .attr("width", width);
```

```
1      var svg_edges = svg.selectAll("line")
2      /* TO FINISH */
3          .data(edges)
4          .enter()
5          /* TO FINISH */
6          .append("line")
7          .style("stroke", "#ccc")
8          .style("stroke-width", 1);
```

```
1      var svg_nodes = svg.selectAll("circle")
2      /* TO FINISH */
3          .data(nodes)
4          .enter()
5          /* TO FINISH */
6          .append("circle")
7          .attr("r", 20)
8          /* TO FINISH */
```

```
9          .style("fill", function (d) {
10             return color(d.index);
11         })
12         .call(force.drag);
```

```
1      var svg_texts = svg.selectAll("text")
2          ...
3          /* TO FINISH */
4          .text(function (d) {
5             return d.node;
6          });
```

```
1      force.on("tick", function () {
2      /* TO FINISH */
3      svg_edges.attr("x1", function (d) {
4          return d.source.x;
5      })
6      /* TO FINISH */
7          .attr("y1", function (d) {
8             return d.source.y;
9          })
10         /* TO FINISH */
11         .attr("x2", function (d) {
12            return d.target.x;
13         })
14         /* TO FINISH */
15         .attr("y2", function (d) {
16            return d.target.y;
17            });
```
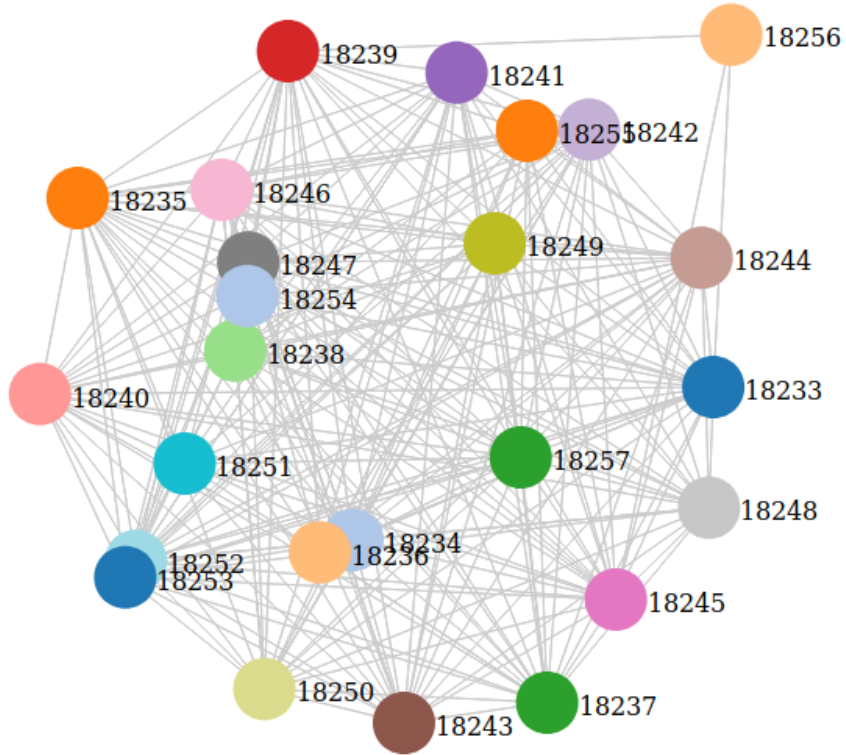
Result:



Figure 10: Screenshot of Connection output result.