# Assignment 4
## on
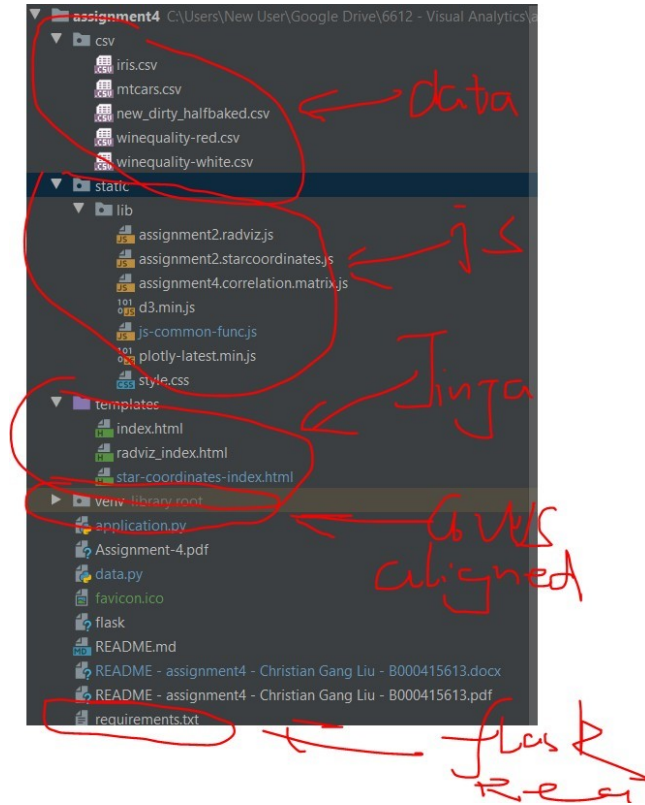# CSCI 6610 Visual Analytics

Christian Gang Liu (B000415613)

November 29, 2019

**Abstract**

*Important Notes: because I am principle programmer of our "outliers" group project implementation, so you might find similar coding on our project as well.*

1. Project technical used: General speaking, the project is based on JS and Python 3.4+ as required by assignment 4 instruction.

   - Flask: A python framework used for quick REST services establishment;

   - D3: A Javascript library used for generating the interactive ML graphic charts;

   - Plotly: A popular javascript library which based on D3 provides the complicated chart view (here I only used it for show the tabular table view on the page, this requirement is not asked by assignment, I used it only to demonstrate how data comes from backend framework (Flask);

   - Others: JQuery / Bootstrap (only for UI layout and event triggering)

   - IDE: PyCharm;

   - AWS (extra work): public cloud deployment:$http://assignment4-dev.us-east-2.elasticbeanstalk.com/$
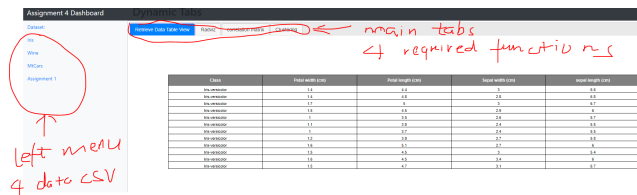
2. Project structure:



As you can see above, I constructed project structure to comply with the AWS elastic sever standardization, in order to deploy on AWS for public access.

I utilized FLASK framework to quickly establish this lightweight project, like what I did for my group project as well $http://flask-env.zsewpnnzda.us-east-2.elasticbeanstalk.com/$.

In order to simplify the assignment 4 review work. I planned to deploy it on AWS for quickly reviewing. Anyhow you can always open my project source folder to execute command: python application.py to exam it on local.

3. Execution: the execution instruction is also mentioned in the README file under project folder.

4. General UI introduction (the details will be described specifically to each of questions of assignment 4):

- Left Menu: they are four different tables which are required by assignment 4:

  - Iris
  - Wine
  - Mtcars
  - Assignment 1 dataset (bonus)

- Main Menu:

  - Tabular view on Html page (extra work): it shows tabular view directly from JSON response of Flask REST services:
  - Radviz: The assignment 2 functionality for each of tables;
  - Correlation Matrix: Heatmap of each table;
  - Clustering: there will be a comparison between original labels (classification) and clustering by K-means, as well as a measurement score showing on the page.
  - Notes: the Bonus requirement of assignment 1 dataset will share the same functionalities.

## 0.1 Assignment 4 Qeustions And Answers

**The specific answers to the assignment 4 questions:**

**1** *[10 Marks] Create a backend that will provide the data and metadata that can be used to display the visualization a. Use an HTTP request to retrieve data from the back-end and use it to generate the visualization on the front-end. b. Tip: You can return, along with the data, some metadata like column names or other information that could be useful to handle and/or display the data in the front-end.*

**answer:** Front-end ( JS retrievals JSON data, sends to html page):



Back-end Flask WSGI endpoint (REST services) – Code snippets:

```python
import data
from flask import Flask, jsonify, render_template, request


application = Flask(__name__)
assignment4Data = data.getAssignment4Data()


@application.route('/')
def hello_world():...


'''...'''
@application.route('/read_iris')
def read_iris():...


@application.route('/read_iris_correlation')
def read_iris_correlation():...


@application.route('/read_iris_clustering')
def read_iris_clustering():...


@application.route('/read_wine')
def read_wine():...


@application.route('/read_wine_correlation')
def read_wine_correlation():...
```

```python
@application.route('/read_car_clustering')
def read_car_clustering():...


@application.route('/read_assignment1')
def read_assignment1():...


@application.route('/read_assignment1_correlation')
def read_assignment1_correlation():...


@application.route('/read_assignment1_clustering')
def read_assignment1_clustering():...


@application.route('/show_table')
def show_table():...


if __name__ == '__main__':
    # Setting debug to True enables debug output. This line should be
    # removed before deploying a production app.
    application.debug = True
    application.run()
```

The returning JSON sample like this: $(http://assignment4-dev.us-east-2.elasticbeanstalk.com/read\_iris)$

"[{\"sepal length (cm)\":5.1,\"Sepal width (cm)\":3.5,\"Petal length (cm)\":1.4,\"Petal width (cm)\":0.2,\"Class\":\"Iris-setosa\"},{\"sepal length (cm)\":4.9,\"Sepal width (cm)\":1.4,\"Petal width (cm)\":0.2,\"Class\":\"Iris-setosa\"},{\"sepal length (cm)\":4.7,\"Sepal width (cm)\":3.2,\"Petal length (cm)\":1.3,\"Petal width (cm)\":0.2,\"Class ...

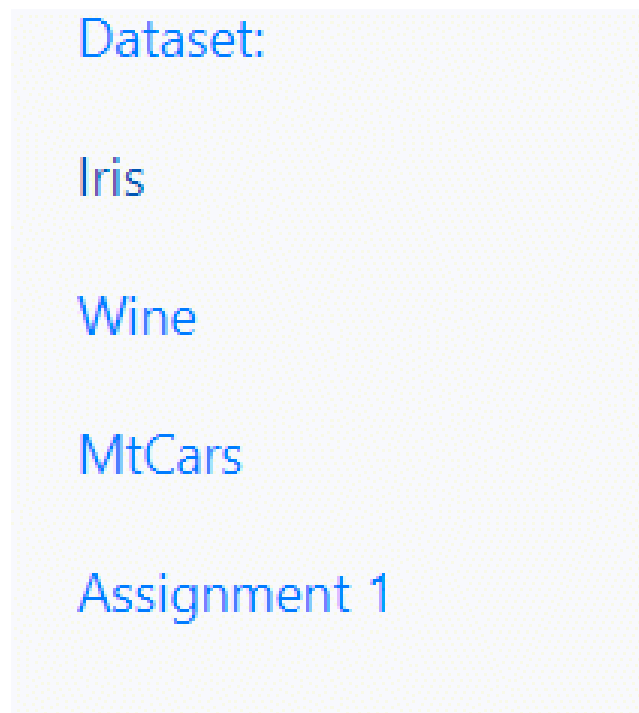The list of REST services showing below (hostname is http://assignment4-dev.us-east- 2.elasticbeanstalk.com):

- application.route('/')
- application.route('/read_iris')
- application.route('/read_iris_correlation')
- application.route('/read_iris_clustering')
- application.route('/read_wine')
- application.route('/read_wine_correlation')

- application.route('/read_wine_clustering')
- application.route('/read_car')
- application.route('/read_car_correlation')
- application.route('/read_car_clustering')
- application.route('/read_assignment1')
- application.route('/read_assignment1_correlation')
- application.route('/read_assignment1_clustering')
- application.route('/show_table')

**2** *[20 Marks] Add an option on the interface to choose a different dataset (iris or winequality) a. The backend will return the new dataset*

**Answer:**

Front-end: when you click any option in left panel, the back-end will return corresponding response:



**Back-end returning:**

- IRIS: http://assignment4-dev.us-east-2.elasticbeanstalk.com/read_iris
- WINE: http://assignment4-dev.us-east-2.elasticbeanstalk.com/read_wine
- MTCARS: http://assignment4-dev.us-east-2.elasticbeanstalk.com/read_car

- ASSIGNMENT1: http://assignment4-dev.us-east-2.elasticbeanstalk.com/read_assignment1

*b. You can keep a state on the front-end and send it on every request to identify the current dataset in use.*
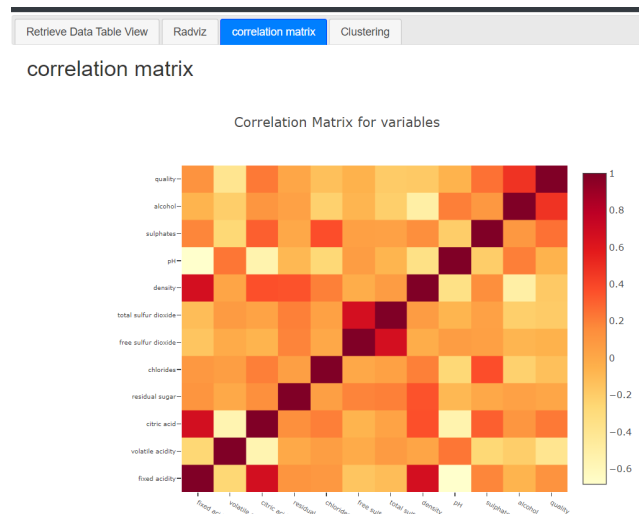
I am using the request parameter $/?data = [dataset\_name]$ to identify current dataset using:

Like $/?data = [iris\|wine\|assignment1\|car]$

**3** *3. [20 Marks] When hovering an instance of a given cluster, show (as a tooltip or in other available space) the correlation matrix for instances of that cluster. a. The correlation matrix should be calculated and returned by the back-end.*

**Answer:**

Front-end:



correlation matrix

Correlation Matrix for variables



Back-end:

```
@application.route('/read_wine_correlation')
def read_wine_correlation():
return jsonify(data.read_wine_correlation())
```

So whenever you click "correlation matrix" tab, it will send REST call like(giving an example):

```
def read_wine_correlation()
-> readAssignment4Data:
rawData = getAssignment4Data().wine correlation_result = rawData.corr()
return correlation_result.to_json(orient='records')
```
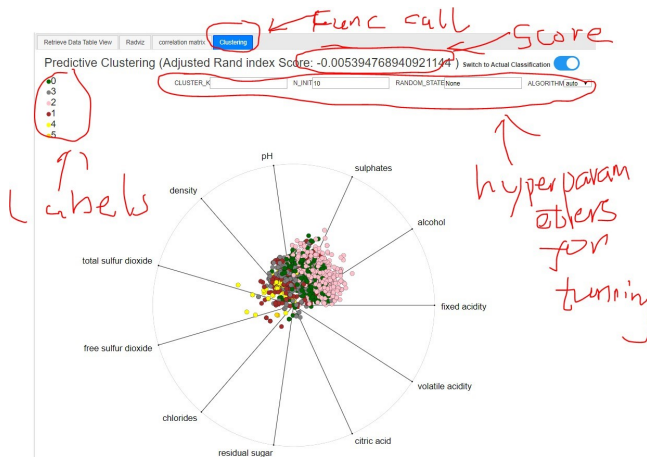
*b. This should be displayed on the front-end through a color matrix. See example of such matrix below (Tip: feel free to use libraries to help you):*

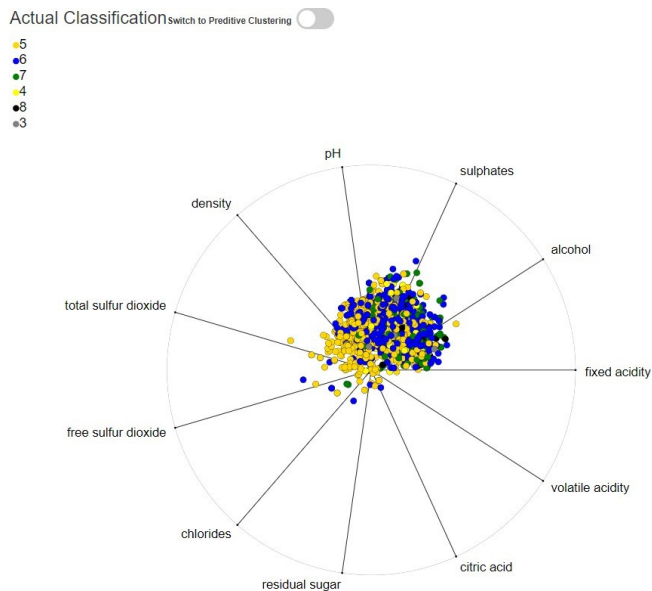So I utilize the corr() function of pandas to generate correlation heatmap

**4** *[40 Marks] Implement a button that requests the backend to clusterize the data using one of: K-Means or DBScan a. You should color the instances using the clustering information*
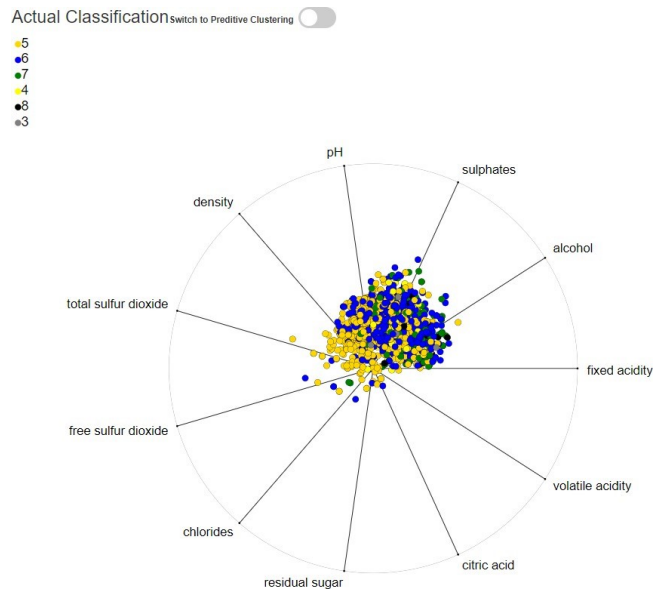
**Answer:**



*b. Add a switch button to choose between the color modes: cluster colors or class-based colors.*

**Answer:**

After we switch back to actual classification:

Actual Classification Switch to Preditive Clustering ⬭

- 5
- 6
- 7
- 4
- 8
- 3



*c. The clusterization should be performed on the same dataset currently seem in the visualization.*
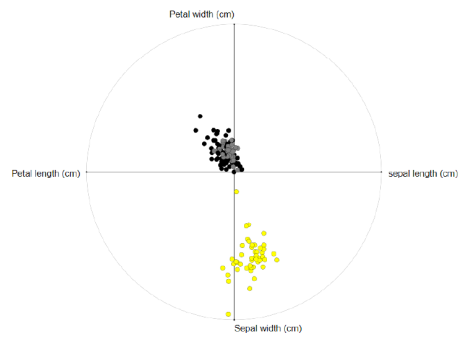
**Answer:**

Since $/?data = [dataset\_name]$ remembers the current dataset. Clustering tab will go further to call

- @application.route('/read_iris_clustering')
- @application.route('/read_wine_clustering')
- @application.route('/read_car_clustering')
- @application.route('/read_assignment1_clustering')

Page will show corresponding result to current dataset:

Predictive Clustering (Adjusted Rand index Score: 0.7302382722834697 ) Switch to Actual Classification

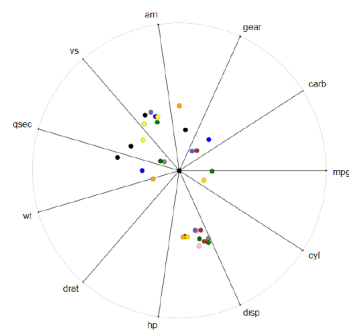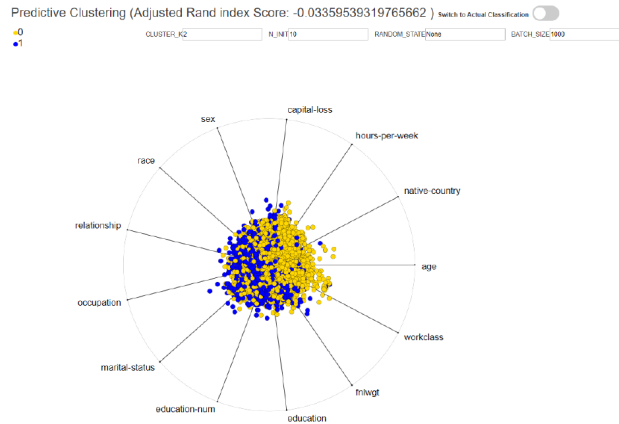Predictive Clustering (Adjusted Rand index Score: 0.006445156285441907 ) Switch to Actual Classification

Predictive Clustering (Adjusted Rand index Score: 1 ) Switch to Actual Classification

Predictive Clustering (Adjusted Rand index Score: -0.03359539319765662 ) Switch to Actual Classification

d. *You just need to implement for one clusterization algorithm.*

**Answer:**

K-MEAN

e. *You may use existing implementations of the clustering algorithms.*

**Answer:**

SKLearn.cluster.kmean

**5** *[10 Marks] Add one (or more) options to configure the parameters of the clustering algorithm a. Clicking the button should make a new clusterization with the new parameters and update the colors on the visualization.*
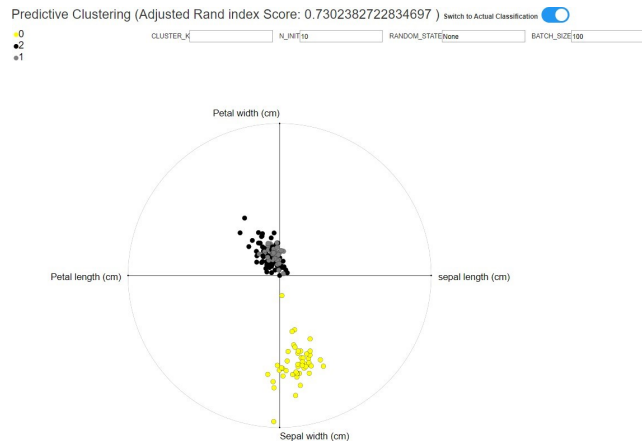
**Answer:**

Here I am going to tune up four major different parameters (Quoted from https://scikit- learn.org/stable/modules/generated/sklearn.cluster.MiniBatchKMeans.html)

- n_clusters : int, optional, default: 8 The number of clusters to form as well as the number of centroids to generate.

- batch_size : int, optional, default: 100 Size of the mini batches.

- random_state : int, RandomState instance or None (default) Determines random number generation for centroid initialization and random reassignment. Use an int to make the randomness deterministic.

- n_init : int, default=3 Number of random initializations that are tried. In contrast to KMeans, the algorithm is only run once, using the best of the n_init initializations as measured by inertia.

The UI toolbar looks like this CLUSTER_K N_INIT 10 RANDOM_STATE None BATCH_SIZE 100
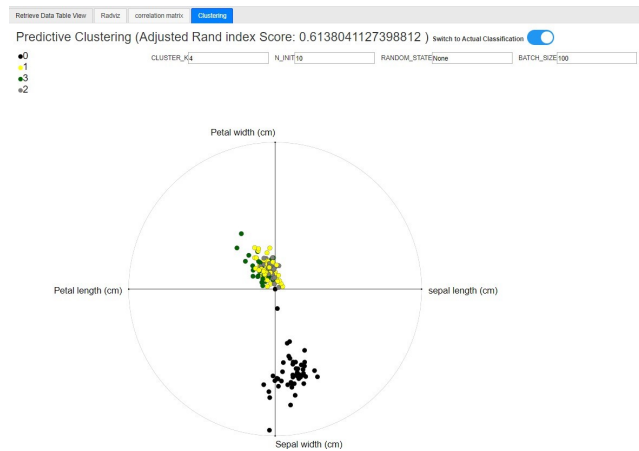
It will allow us to adjust four different major parameters of k-mean , let us see the effects:
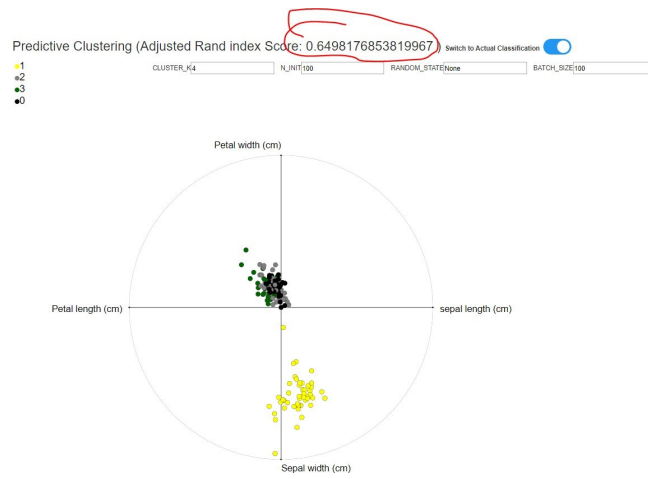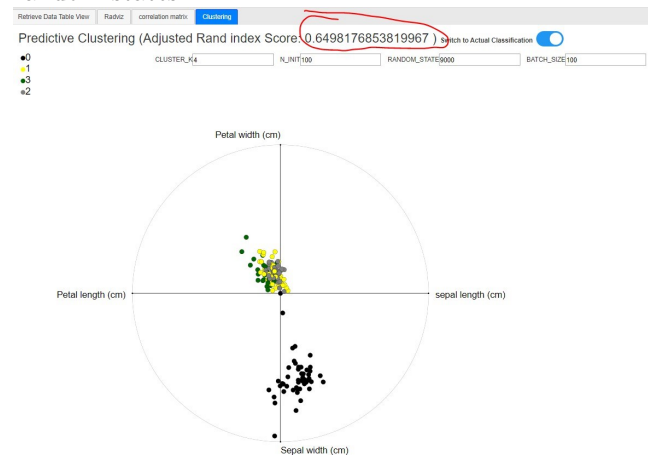
Originally we have:

Predictive Clustering (Adjusted Rand index Score: 0.7302382722834697 ) Switch to Actual Classification



let us see how hyperparameters tunning up:

- K-cluster:



- N_init adjustment:

Predictive Clustering (Adjusted Rand index Score: 0.6498176853819967 ) Switch to Actual Classification

| | CLUSTER_K 4 | N_INIT 100 | RANDOM_STATE None | BATCH_SIZE 100 |

- 1
- 2
- 3
- 0

Petal width (cm)

Petal length (cm)                    sepal length (cm)

Sepal width (cm)

- random_states:

Retrieve Data Table View    Radviz    correlation matrix    Clustering

Predictive Clustering (Adjusted Rand index Score: 0.6498176853819967 ) Switch to Actual Classification

| | CLUSTER_K 4 | N_INIT 100 | RANDOM_STATE 9000 | BATCH_SIZE 100 |

- 0
- 1
- 3
- 2

Petal width (cm)

Petal length (cm)                    sepal length (cm)

Sepal width (cm)

- Batch_size adjustment:

Predictive Clustering (Adjusted Rand index Score: 0.6498176853819967 ) Switch to Actual Classification

**You can see the result will slightly different as we tune the parameters based on the measurement score.**

**6** *[+30 Bonus Mganr4k1s7] 1A2d6d -anCohpritsiotinanonGthaenginLteirufaceBt0o0c0h4o1o5se61to3see the preprocessed dataset generated by your A1 assignment. a. RadViz/StarCoordinates should only show the numerical columns as anchor points*

**Answer:**

Using sklearn normalization library to scale and label the columns as digit numbers

*b. The Categorical columns should be shown as the color of the plot. Make an input box selector in the interface to choose the categorical column to be shown as the color.*
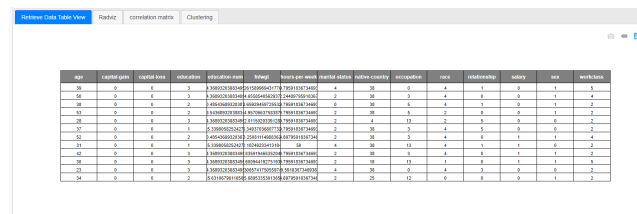
**Answer:**

Make the salary as label: ¡50K = 0, ¿50K = 1

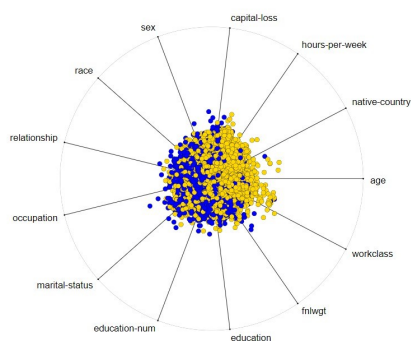Assignment 1 dataset will share same functionalities as three other datasets:

**Correlation:**

correlation matrix



**tabular view:**



**clustering:**



**original Redviz:**

12

Actual Classification Switch to Preditive Clustering ⚪
● 0
● 1

## 0.2 Miscellaneous: Error Handling

I added error message box just in the form to show any error messages, whenever clustering has the problem at the back-end. Eg: when user tries to input random_state = 9000, which is not valid for k-means, then front-page will honestly record the error messages returning from back-end, like this:



'9000' cannot be used to seed a numpy.random.RandomState instance

## 0.3 Referencing

**1** D3: https://d3js.org/

**2** Plotly: https://plot.ly/

**3** Flask:https://realpython.com/tutorials/flask/

**4** Jinja2 Template: https://jinja.palletsprojects.com/en/2.10.x/

**5** Jquery: https://api.jquery.com/

**6** Bootstrap: https://getbootstrap.com/

**7** AWS: https://www.awseducate.com/