

REPORT

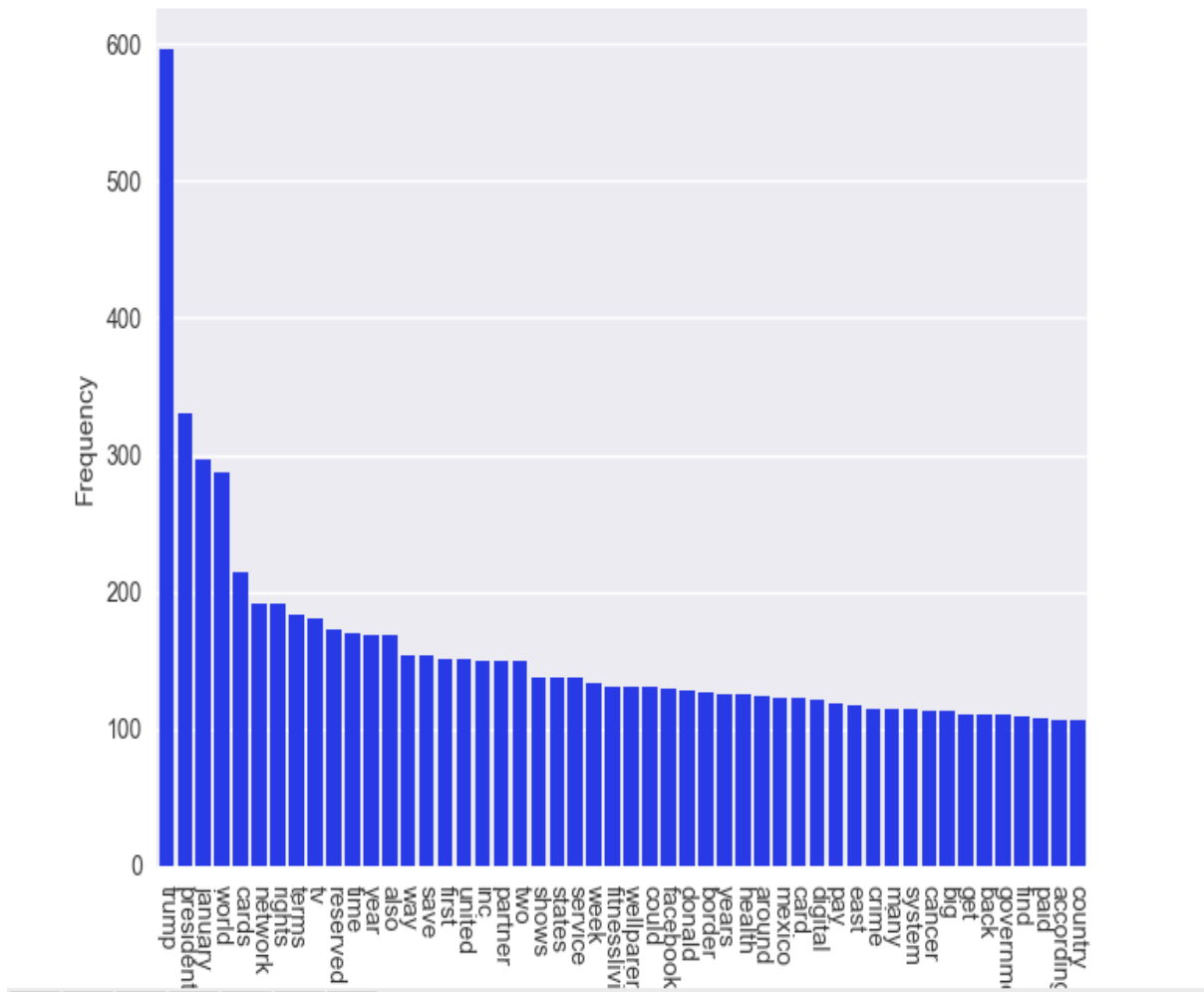
Manish Jaisinghani

1.Objective: Demonstrate Kmeans algorithm using different distance metrics. The metrics involve Original KMeans, KMeans with Cosine distance, KMeans with Euclidean distance and KMeanse with Jaccard distance.

2.Data Collection: The data collection process is completed in two steps as described below:

- a. Website List: I collected the website list by using online website scrapper <http://www.xsitemap.com>. The scrapper parses all the sub URLs from a given URL and provides a list of the parsed URLs. The URL list has been divided into 4 categories and 25 URLs have been taken from each category:
 - Politics – 25 articles
 - Health – 25 articles
 - World – 25 articles
 - Technology – 25 articles
- b. Website Scrapper: I have written HTML website scrapper in python using beautiful soup and NLTK as packages to parse the HTML packages and remove words from the collected list that do not add to meaning of the article, for example “a, the, he, she, these etc.” A complete list of the words which have been remove or not taken into consideration has been attached by the file name “stopwords”.

Top 50 Word Frequency



3.Data Processing: The data processing is done in two parts as described below:

- a. Website parsing: A list of URLs is provided in the file by name “**website_list**”. The python script parses through each website one by one. The script makes use of beautiful soup, parses the website and pulls out text data eliminating all the symbols, punctuations, styles and scripts. Once the data is collected the data is recorded in a list data structure and then stored making sure that there are no duplicate words in the list. The same procedure is repeated for every URL. While parsing the text of current article frequency of words occurring in the article is also recorded. The final step is just writing this data to a CSV file(“**data.csv**”). In the csv file each article is represented as a row and

each word is recorded in a column and the final matrix consists of the frequency of words.

- b. **Data manipulation:** For Jaccard distance metric we convert the frequency matrix to similarity matrix which only reports presence or absence of a word in an article. We also normalize the data by feature subset selection.

4. Module specifications: The algorithms have been implemented using python 3 and below listed libraries have been used:

- **Pandas** – Deals with csv data and helps in manipulation.
- **Numpy** – Deals with data in the form of arrays and assists in manipulation
- **Matplotlib** – Helps to visualize the data.
- **Seaborn** – Another library which helps in data visualization, the only difference here is seaborn is vast than matplotlib.
- Used KMeans algorithm provided in lecture notes and implemented all the distance metrics in the same algorithm.

5. Terminology & Explanation:

- *Euclidean Distance:* The **Euclidean distance** between points **p** and **q** is the distance between two points and is given by:

$$d(\mathbf{p}, \mathbf{q}) = d(\mathbf{q}, \mathbf{p}) = \sqrt{\sum_{i=1}^n (q_i - p_i)^2}.$$

- *Cosine Distance:* Cosine similarity is a measure of similarity between two non-zero vectors of an inner product space that measures the cosine of the angle between them.

$$\text{similarity} = \cos(\theta) = \frac{\mathbf{A} \cdot \mathbf{B}}{\|\mathbf{A}\| \|\mathbf{B}\|}$$

$$\text{Cosine Distance} = 1 - \text{Similarity}$$

- *Jaccard Distance:* Given two objects, *A* and *B*, each with *n* binary attributes, the Jaccard coefficient is a useful measure of

the overlap that A and B share with their attributes. Each attribute of A and B can either be 0 or 1. The total number of each combination of attributes for both A and B are specified as follows:

- represents the total number of attributes where A and B both have a value of 1.
- represents the total number of attributes where the attribute of A is 0 and the attribute of B is 1.
- represents the total number of attributes where the attribute of A is 1 and the attribute of B is 0.
- represents the total number of attributes where A and B both have a value of 0.

The Jaccard distance, d_J , is given as

$$d_J = \frac{M_{01} + M_{10}}{M_{01} + M_{10} + M_{11}} = 1 - J.$$

- *sum of squared errors of prediction (SSE)*: is the sum of the squares of residuals (deviations predicted from actual empirical values of data). It is a measure of the discrepancy between the data and an estimation model.

$$SSE = \sum_{i=1}^n (y_i - (\alpha + \beta x_i))^2$$

6. Implementation details: Below are the details on how the algorithms have been written and pre-requisites for successfully running the same:-

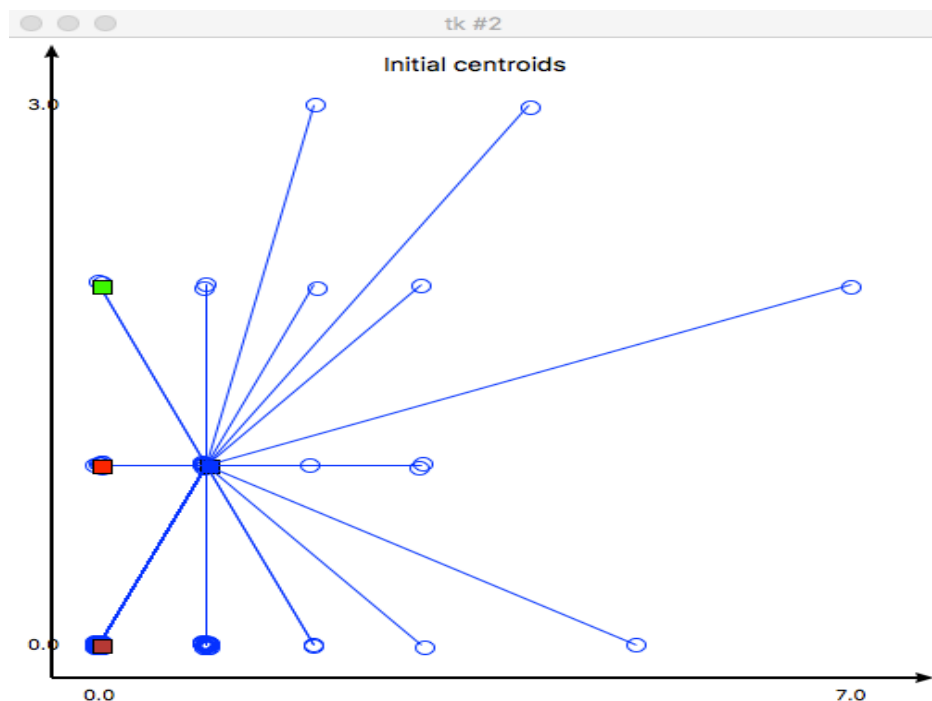
- **kmeans_Cosine.py** - The file contains the implementation of KMeans algorithm using Cosine distance metric.
- **kmeans_Euclidean.py** – The file contains the implementation of KMeans algorithm using Euclidean distance metric.
- **Kmeans_Jaccard.py** – The file contains the implementation of KMeans algorithm using Jaccard distance metric.
- **kmeans_Original.py** – The file contains the implementation of KMeans algorithm.
- **Jaccard_Data.csv** – This file contains the data with respect to Jaccard distance. A separate data file is used for jaccard distance because jaccard reports presence or absence rather than frequency of words.

- **Data_KMeans.csv** – This file contains word frequency metrics which is utilized by all the other distance metrics except jaccard similarity.
- **Constraints** – The system should be installed with all python modules listed in “Module Specifications” section.
- The python script should be executed using python 2.0.

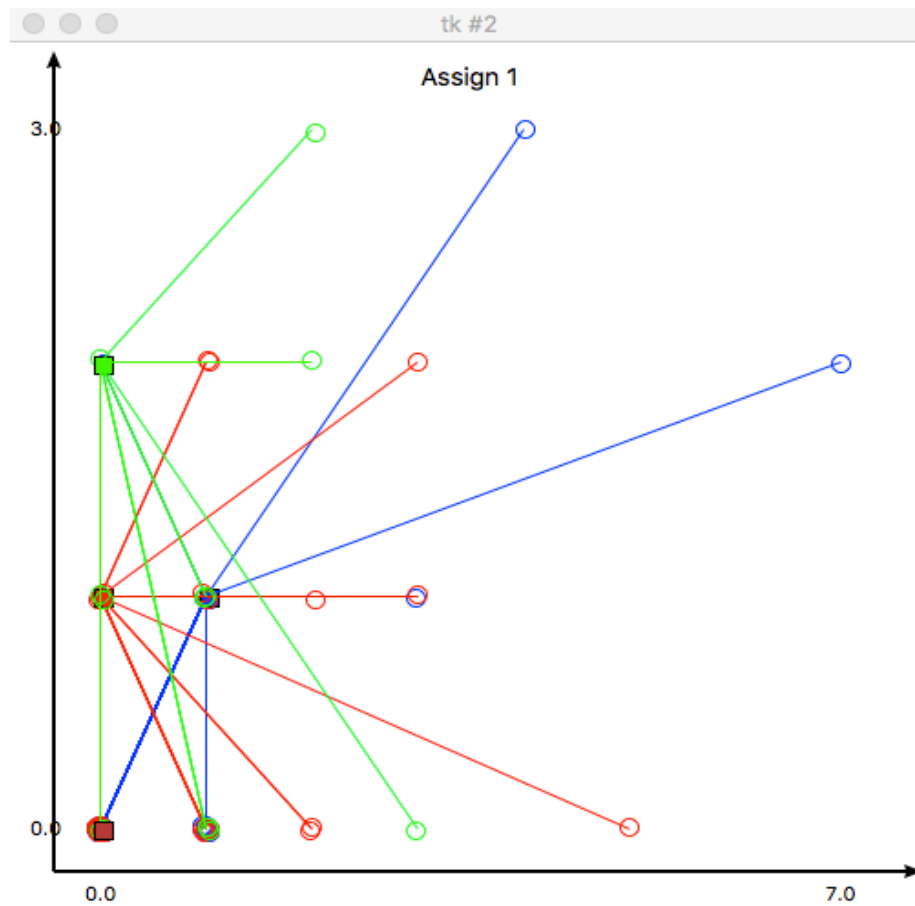
7. Data Analysis: The data analysis section is divided in 5 subsections:

- KMeans algorithm with Eculidean distance.
 - KMeans algorithm with Cosine distance.
 - KMeans algorithm with Jaccard distance.
 - KMeans algorithm original.
- a. KMeans algorithm: The original KMeans algorithm converges in 6 iterations.
- Below is the graphical representation of cluster formation:

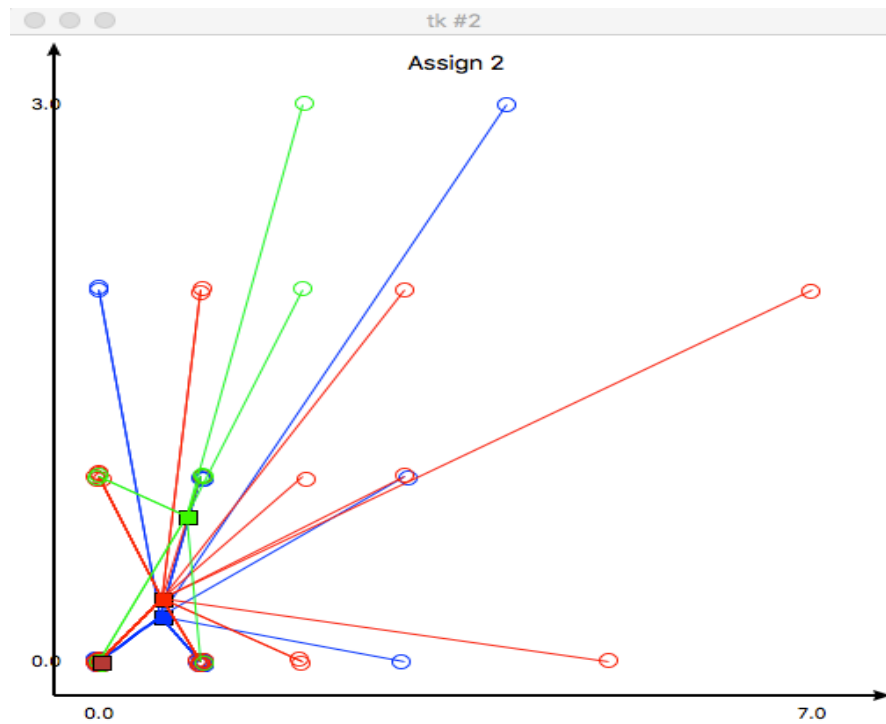
- Initial Centroids:



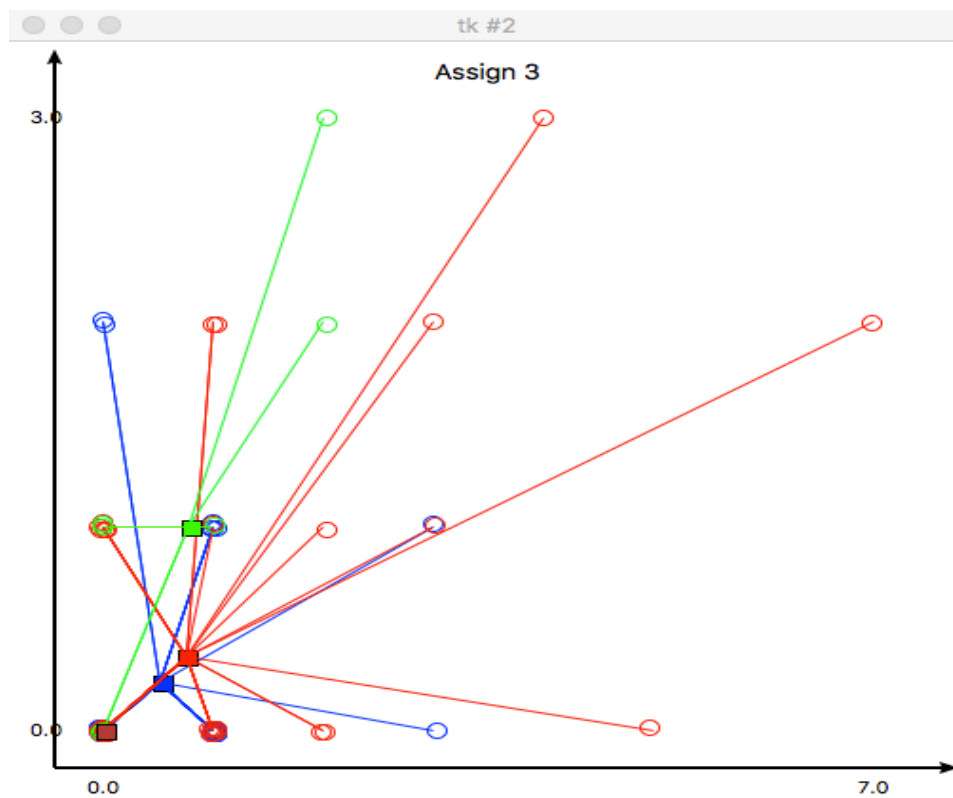
- First iteration depicting assignment of data points to centroids:



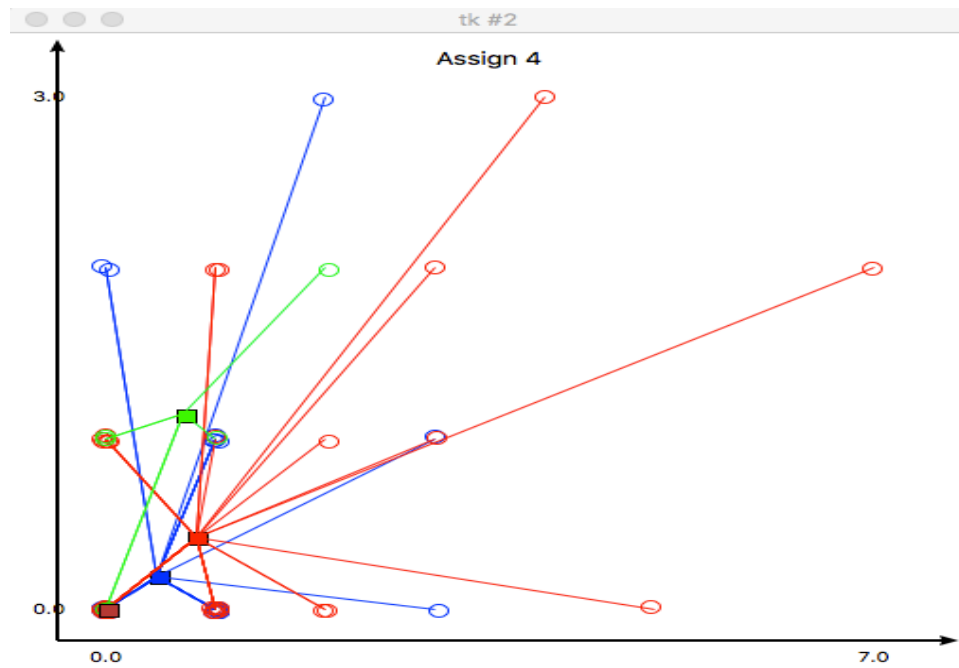
- Second iteration depicting assignment of data points to centroids:



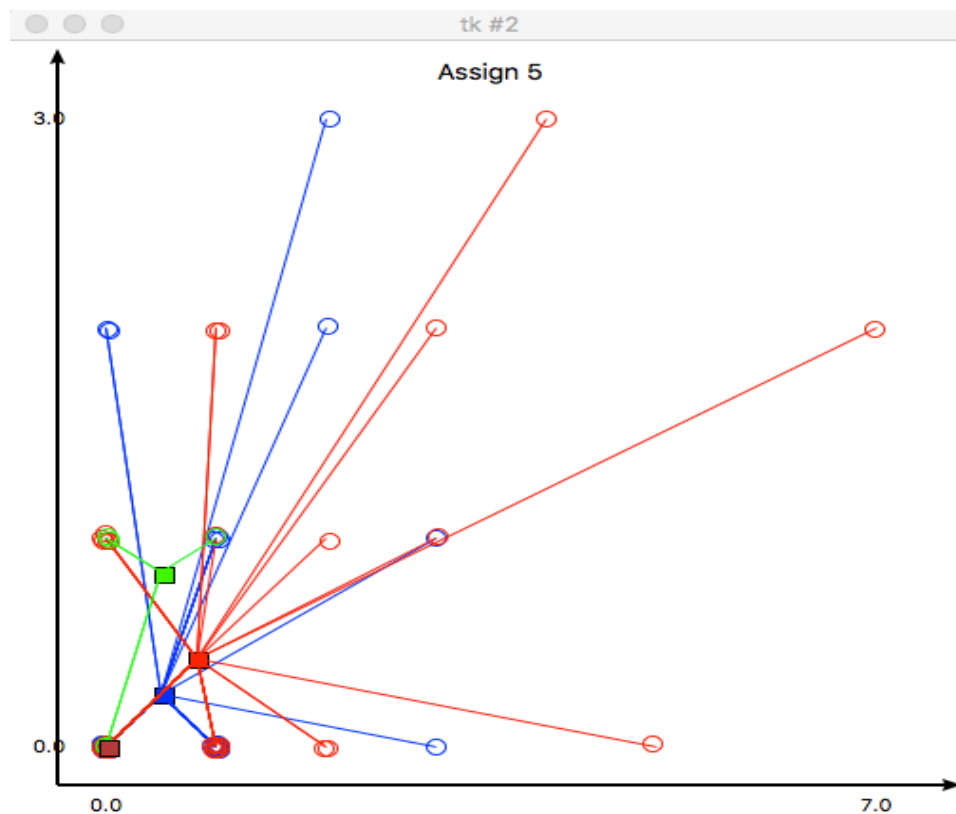
- Third iteration depicting assignment of data points to centroids:



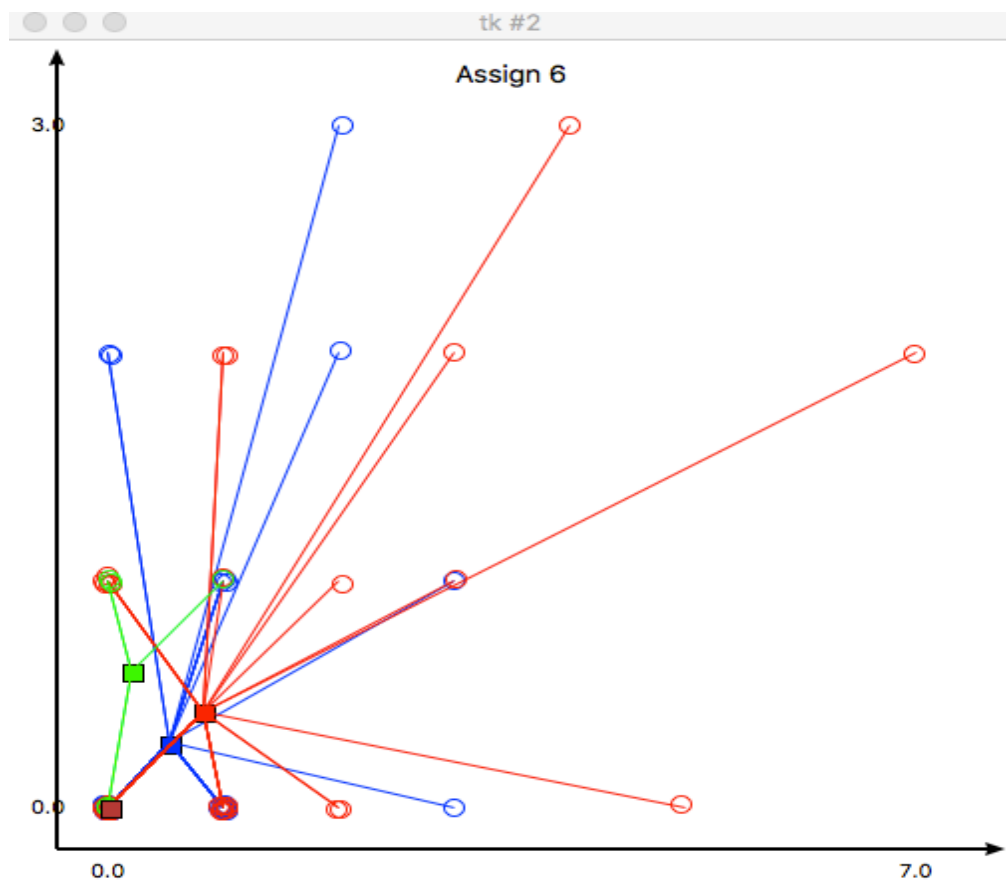
- Fourth iteration depicting assignment of data points to centroids:



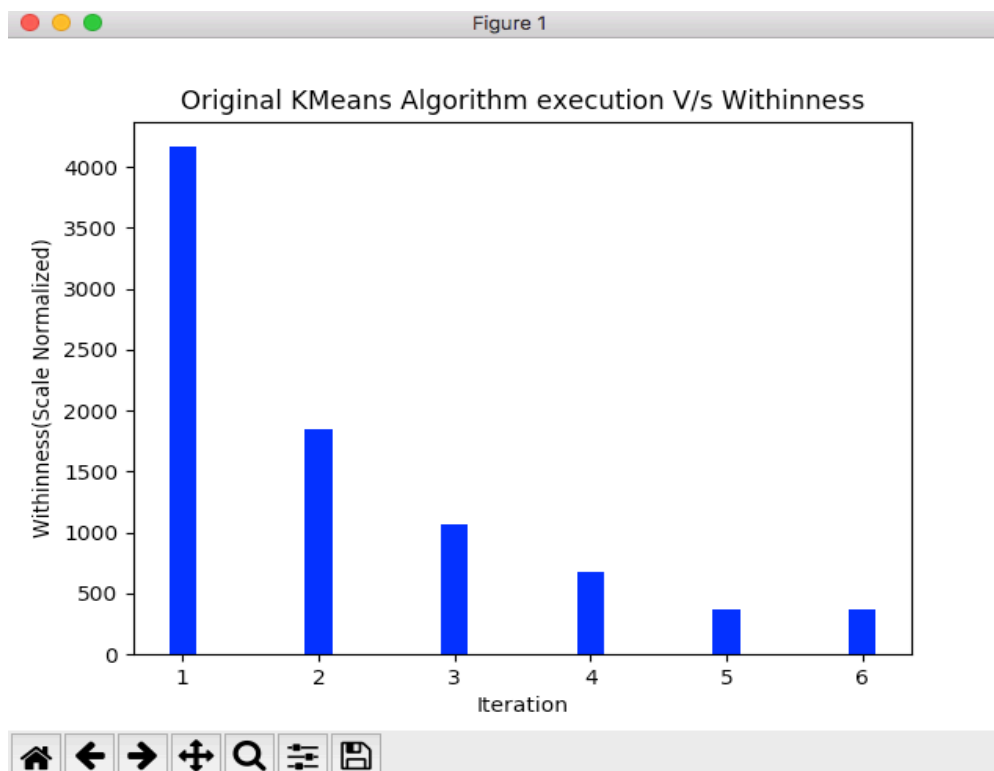
- Fifth iteration depicting assignment of data points to centroids:



- Sixth iteration depicting assignment of data points to centroids:



- The **withinness calculation** for Original KMeans is displayed below:



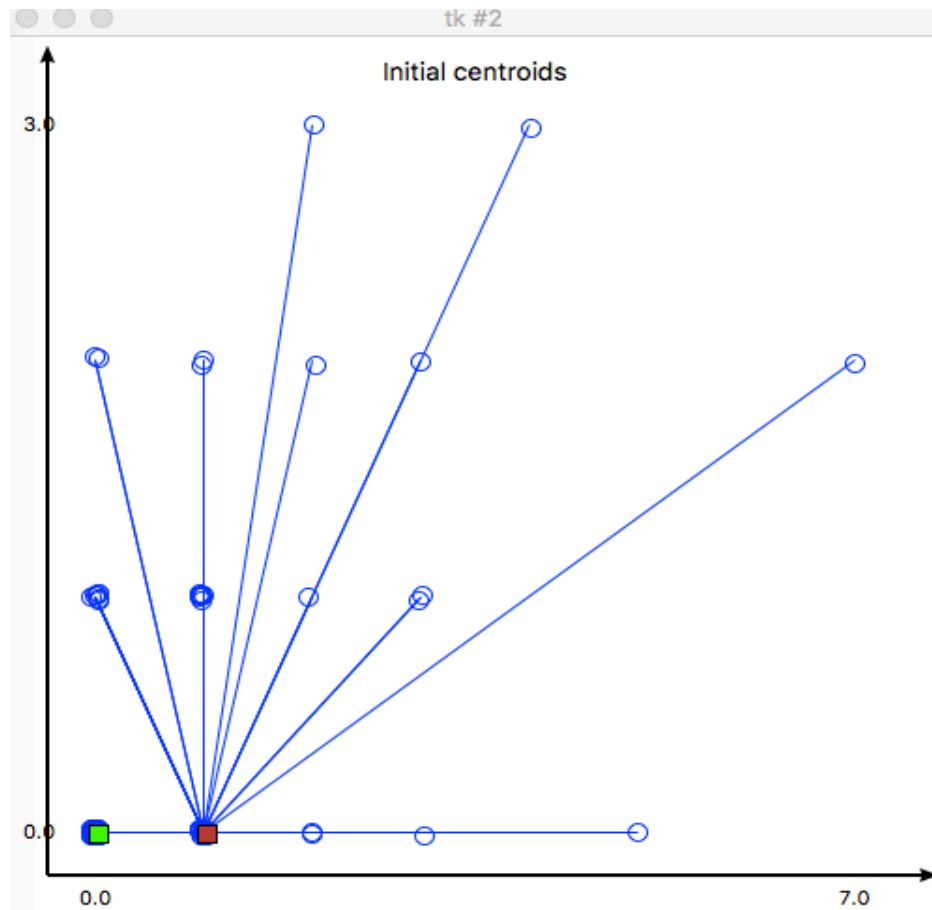
- Clearly visible from the above graph the withinness of cluster formation decreases along the successive assignments.
- Centroid allocation: Below are the centroids allocated once the algorithm converges

```
centroid0    0.53 0.29 3.63 0.02 0.27 0.27 0.43 0.37 0.31 0.43 0.49 0.45 0.02 0.00 0.24 0.63 0.94 0.51 0.69 1.73 0.82 0.67 0.73 0.63 0.63 0.14 0.82 0.08 0.
22 0.08 0.43 0.00 1.51 0.39 0.98 0.35 1.10 0.14 0.14 0.45 0.02 0.10 0.31 0.67 0.82 0.33 0.59 0.37 0.47 0.18 0.61 0.61 0.69 0.35 1.27 1.10 1.80 0.18 0.71 0.69
0.78 0.39 0.39 1.98 0.80 0.41 0.41 0.35 0.12 0.29 0.71 0.20 0.02 0.43 0.90 1.27 0.24 0.78 0.33 0.39 0.24 0.35 0.31 0.27 2.39 0.10 0.02 0.31 0.00 0.20 0.16 0.0
2 0.08 0.06 0.27 0.59 0.71 0.16 0.08 0.12 0.29 0.27 0.12 0.10 0.06 0.47 0.53 0.08 0.10 1.04 0.73 1.45 1.00 0.55 0.49 0.53 1.12 2.00 2.41 0.80 0.73 0.45 0.63 1
.73 0.63 2.69 2.47 2.51 0.02 0.37 7.84 3.59 0.37 0.14 0.08 0.35 0.16 0.12 0.98 0.27 0.37 0.29 0.94 0.22 0.10 0.08 0.16 0.00 0.27 0.00 0.76 0.31 2.29 0.16 0.69
0.37 0.18 0.90 0.14 0.31 0.55 0.04 0.14 0.14 0.35 0.71 0.08 0.04 0.33 0.43 0.27 0.06 0.00 0.43 0.24 0.04 0.49 0.27 0.16 0.14 0.53 0.55 0.31 0.29 1.06 1.37 0.
37 1.08 1.27 0.51 0.27 0.31 0.04 0.10 0.33 0.47 0.12 0.45 0.16 0.16 0.63 0.37 0.63 0.29 0.33 0.84 0.27 0.80 0.51 0.08 0.80 0.61 1.08 0.82 0.57 0.24 0.20 0.20
0.10 0.33 0.76 0.61 0.08 0.22 0.41 0.10 0.73
centroid1    0.82 0.42 1.96 0.78 2.56 2.56 3.47 1.09 2.04 2.09 2.11 2.53 0.20 0.13 0.64 1.00 1.69 1.00 1.51 4.16 1.44 0.80 0.93 1.13 0.96 0.40 1.33 0.33 0.
80 1.22 0.98 0.33 0.58 0.53 1.09 0.98 2.11 0.04 0.42 0.44 0.38 0.16 0.67 1.78 0.51 0.27 0.80 0.27 0.31 0.87 0.71 0.27 1.44 0.20 1.53 0.82 1.47 0.58 0.36 0.62
2.04 0.33 0.58 1.02 1.22 0.22 0.56 1.18 0.36 0.44 0.42 0.04 0.09 0.60 0.67 0.22 0.44 1.64 0.09 0.33 0.53 0.67 0.27 0.71 0.24 0.71 0.22 0.40 0.27 0.11 0.40 2.4
9 0.49 0.33 0.13 1.02 0.62 0.60 0.53 0.62 0.53 2.22 0.04 0.13 0.04 0.96 0.42 0.09 0.04 1.02 1.16 1.02 2.11 1.22 0.98 1.62 1.07 0.98 1.91 1.76 1.20 1.07 1.40 2
.11 1.69 1.16 1.00 1.18 0.11 0.09 0.60 2.22 0.33 0.56 0.04 0.13 0.02 0.36 0.20 0.02 0.71 0.20 0.09 0.09 0.13 0.47 0.31 0.00 0.33 0.00 0.18 0.47 0.33 0.44 0.13
0.51 0.07 0.18 0.38 0.02 0.13 0.00 0.38 0.00 0.02 0.09 0.00 0.04 0.00 0.04 0.07 0.31 0.11 0.02 0.02 0.04 0.44 0.71 0.24 0.13 0.07 0.58 0.51 0.38 1.02 1.42 0.
16 0.76 2.11 0.56 0.84 0.16 0.22 0.49 0.09 0.11 0.36 0.27 0.09 0.27 0.73 0.22 0.64 0.09 1.42 1.20 0.53 0.18 0.40 0.00 0.84 1.40 1.80 1.82 0.42 0.51 0.18 0.16
0.38 0.69 0.29 0.18 0.18 0.58 0.44 0.02 0.49
centroid2    0.20 0.60 13.00 0.00 0.00 0.00 0.00 0.00 1.20 0.00 0.00 0.00 0.00 0.00 0.20 1.00 1.00 1.00 1.20 2.40 1.60 3.00 2.20 2.40 2.40 0.20 0.40 0.60 1
.00 1.00 0.80 0.20 0.20 0.00 1.00 1.00 3.20 0.00 0.60 2.20 1.40 0.60 1.40 3.40 3.60 0.20 0.00 4.80 2.40 0.20 0.00 6.80 2.20 0.60 3.80 1.40 3.20 0.60 2.20 0.20
4.20 0.80 0.20 1.00 2.40 1.40 0.00 0.60 0.00 0.60 1.40 0.00 0.60 1.40 4.00 6.40 0.20 2.40 1.40 0.20 0.00 1.60 0.80 1.00 1.60 0.20 0.00 0.20 0.00 0.00 0.00 0.
00 0.00 0.00 0.60 1.80 0.60 0.40 0.20 0.00 0.20 0.80 0.60 1.00 0.00 1.40 0.60 0.60 0.80 0.80 6.00 0.80 1.60 1.00 1.20 0.80 0.80 1.20 1.60 0.80 2.00 0.80 3.60
2.00 1.40 1.40 1.00 1.20 6.20 4.40 37.00 4.00 18.00 1.60 1.40 20.60 2.20 0.60 14.20 15.20 1.80 2.20 8.80 2.60 1.40 2.60 2.00 0.20 1.20 1.20 2.80 1.20 1.40 1.6
0 0.20 1.20 0.40 1.40 0.60 0.20 2.00 1.00 0.60 0.40 1.00 0.40 0.80 0.40 1.00 2.20 0.80 0.60 0.60 0.80 1.40 1.40 2.60 3.80 0.60 1.40 2.60 3.00 1.80 0.40 2.40 2
.00 0.40 0.80 2.00 1.40 1.00 2.20 0.80 1.80 0.20 0.80 0.20 1.20 0.40 2.20 2.00 1.00 1.00 1.00 0.40 3.00 0.20 0.20 1.60 1.00 2.00 1.40 3.00 2.80 0.40 1.40 1.20
1.00 0.20 0.60 0.60 0.40 1.20 1.00 0.40 0.20 1.00
centroid3    0.00 0.00 0.00 2.00 3.00 3.00 3.00 1.00 2.00 2.00 2.00 2.00 0.00 0.00 0.00 1.00 2.00 1.00 1.00 3.00 2.00 0.00 2.00 2.00 2.00 0.00 4.00 0.00 0.
00 0.00 1.00 0.00 1.00 1.00 1.00 30.00 3.00 0.00 0.00 0.00 0.00 0.00 1.00 0.00 0.00 3.00 0.00 0.00 0.00 0.00 0.00 0.00 0.00 0.00 0.00 0.00 0.00 0.00 4.00 1.00
3.00 0.00 22.00 1.00 0.00 0.00 0.00 1.00 0.00 0.00 0.00 0.00 0.00 1.00 0.00 0.00 2.00 1.00 0.00 3.00 2.00 0.00 0.00 1.00 1.00 1.00 0.00 0.00 0.00 0.00 0.00 0.00 0
.00 0.00 0.00 0.00 0.00 0.00 0.00 1.00 1.00 1.00 8.00 0.00 0.00 0.00 0.00 0.00 0.00 0.00 1.00 1.00 1.00 2.00 1.00 1.00 1.00 1.00 2.00 2.00 1.00 1.00 1.00
2.00 1.00 1.00 1.00 1.00 0.00 0.00 0.00 1.00 0.00 1.00 0.00 0.00 0.00 0.00 0.00 0.00 0.00 0.00 0.00 0.00 0.00 0.00 0.00 0.00 0.00 0.00 0.00 0.00 0.00 0.00 0.00 0.
00 0.00 1.00 0.00 0.00 0.00 0.00 0.00 0.00 0.00 0.00 0.00 0.00 0.00 0.00 0.00 0.00 0.00 0.00 0.00 0.00 0.00 0.00 0.00 0.00 0.00 0.00 0.00 0.00 0.00 0.00 0.00 0.00 0.
00 0.00 1.00 2.00 0.00 0.00 0.00 0.00 1.00 0.00 0.00 0.00 0.00 0.00 0.00 0.00 0.00 0.00 0.00 0.00 0.00 0.00 0.00 0.00 0.00 0.00 0.00 0.00 0.00 0.00 0.00 0.00 0.00 0.00
0.00 0.00 0.00 0.00 0.00 0.00 0.00 0.00 0.00
```

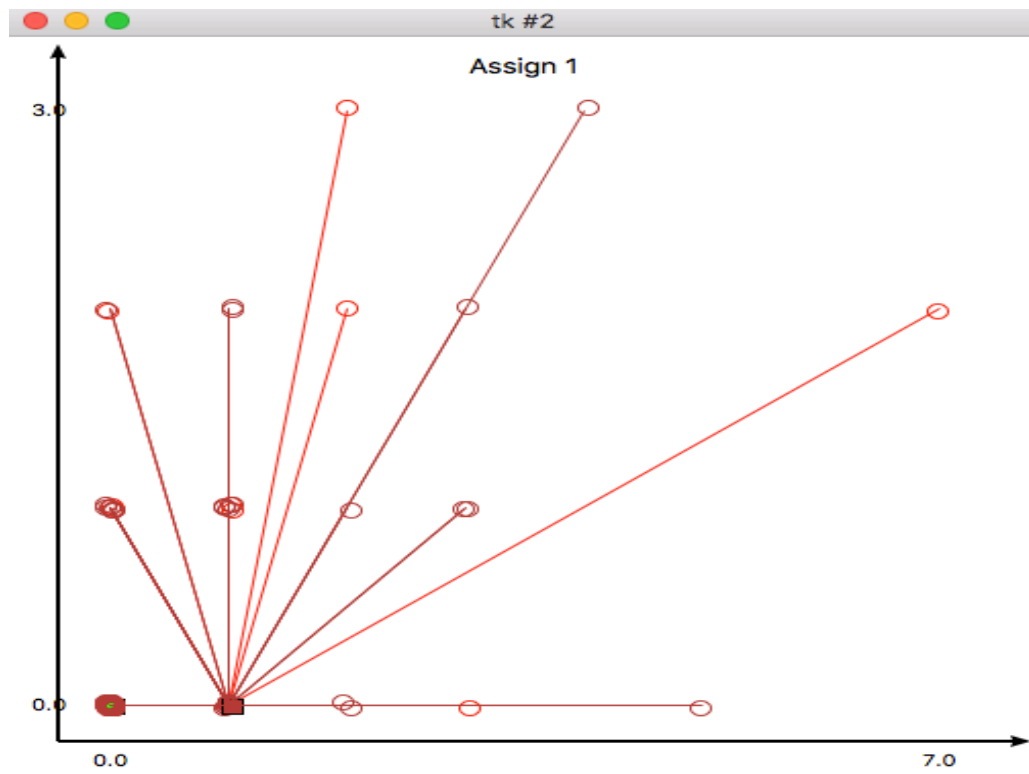
b. KMeans algorithm with Euclidean: The Euclidean distance is used as distance metric with KMeans algorithm. Using euclidean as the distance metric KMeans algorithm converges in 8 iterations.

- Below is the graphical representation of cluster formations:

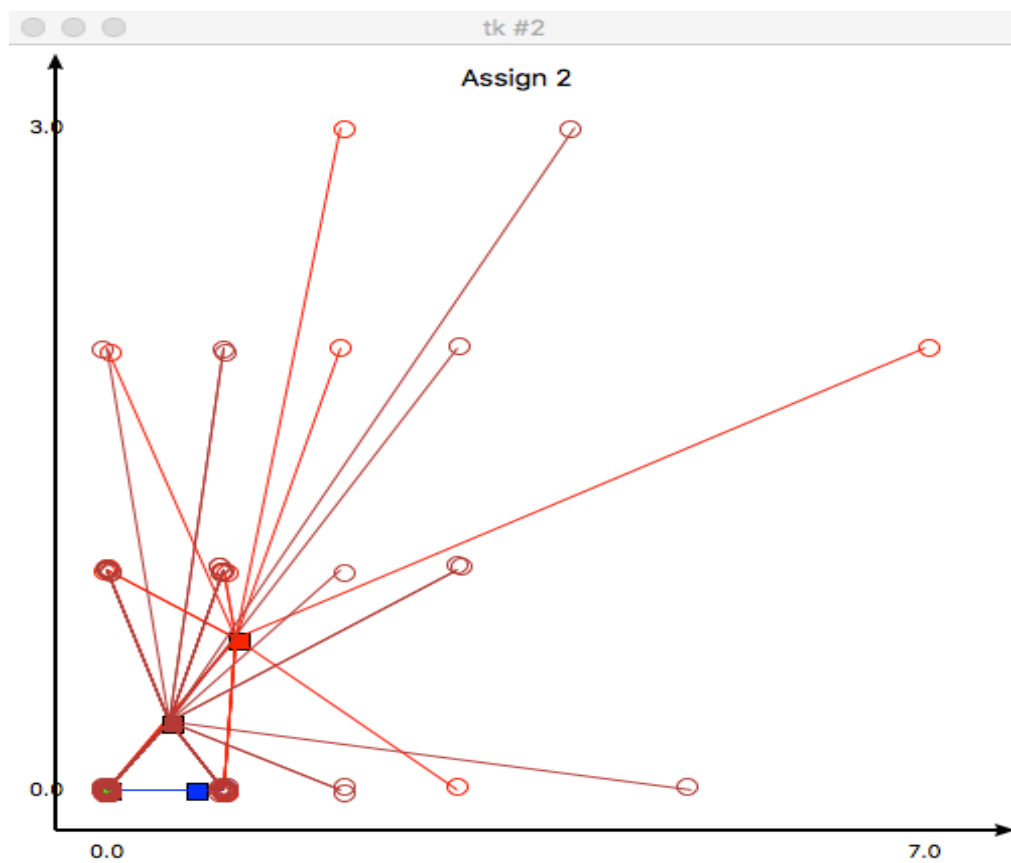
- Initial Centroids:



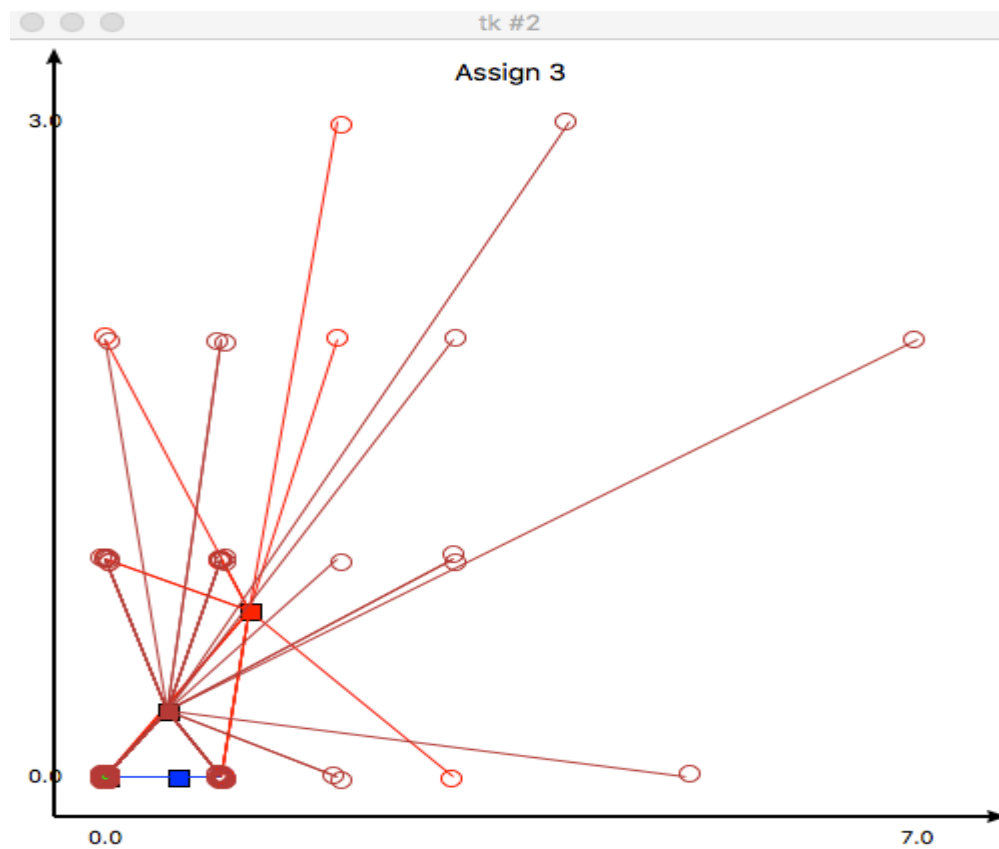
- First iteration depicting assignment of data points to centroids:



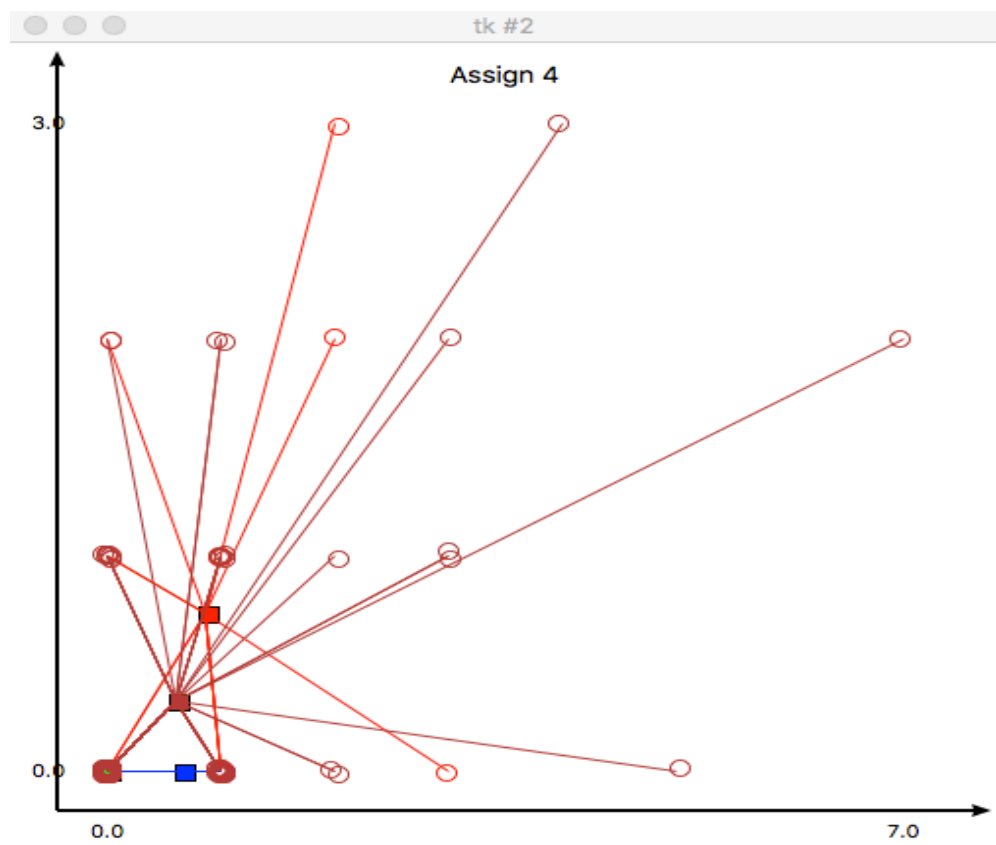
- Second iteration depicting assignment of data points to centroids:



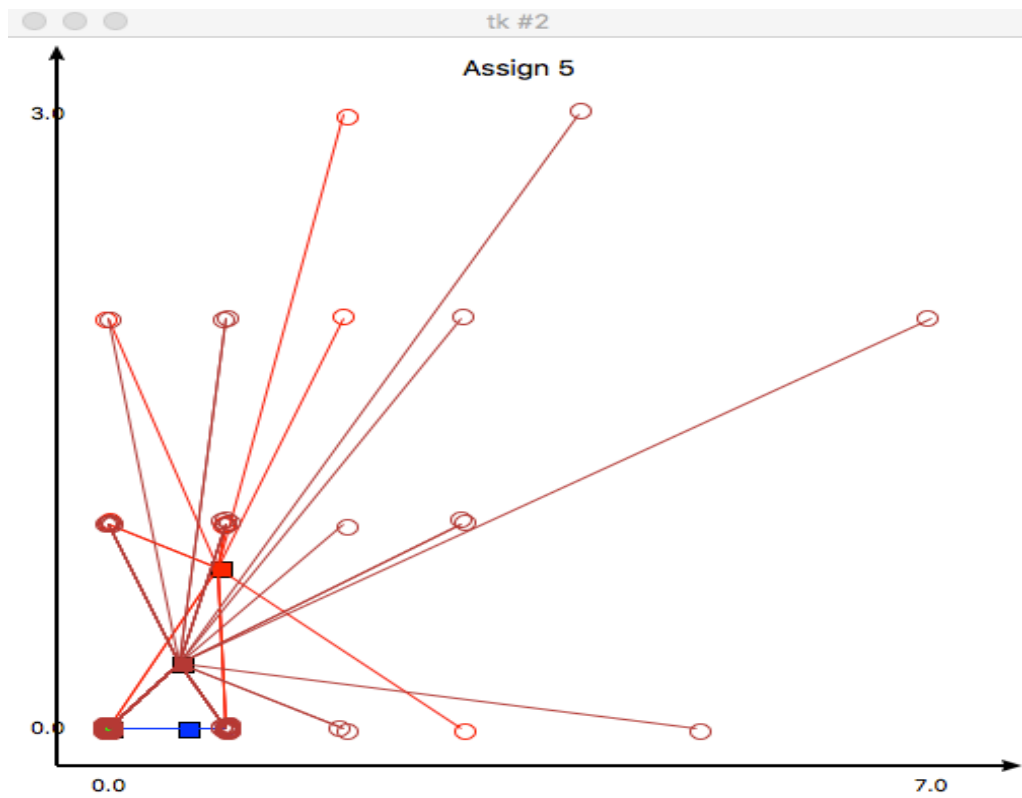
- Third iteration depicting assignment of data points to centroids:



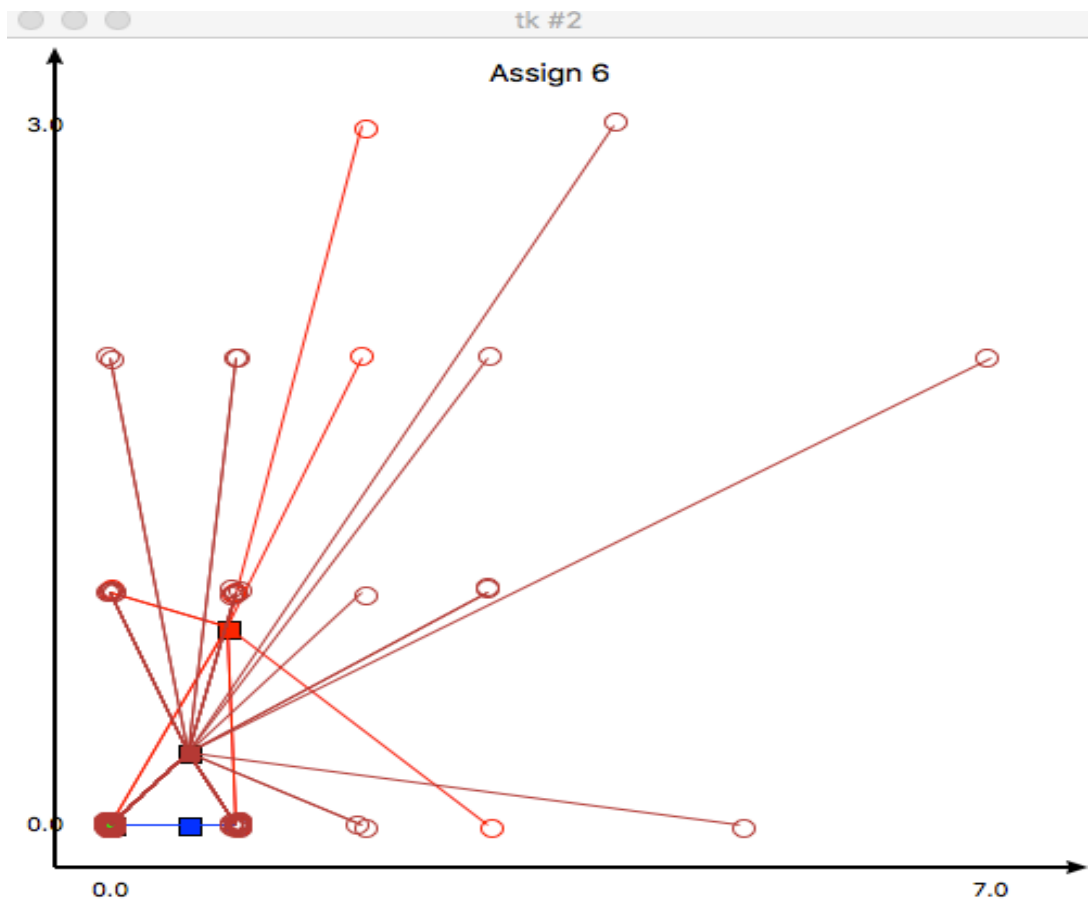
- Fourth iteration depicting assignment of data points to centroids:



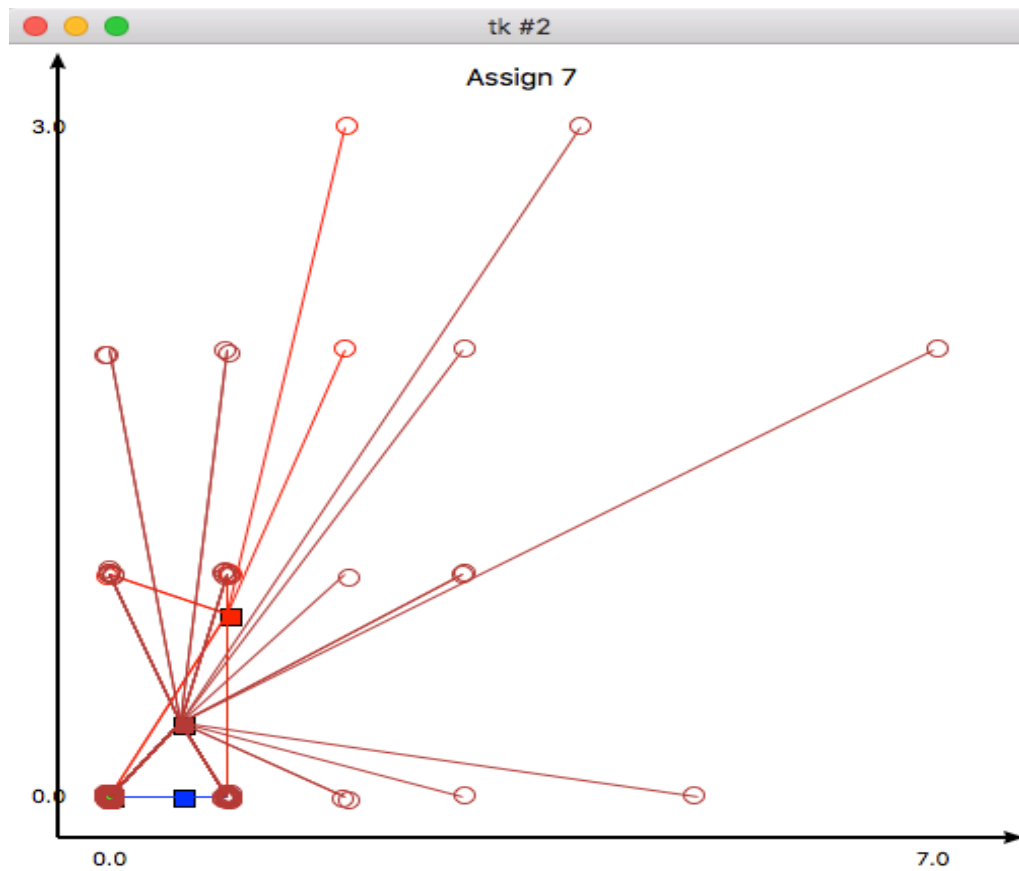
- Fifth iteration depicting assignment of data points to centroids:



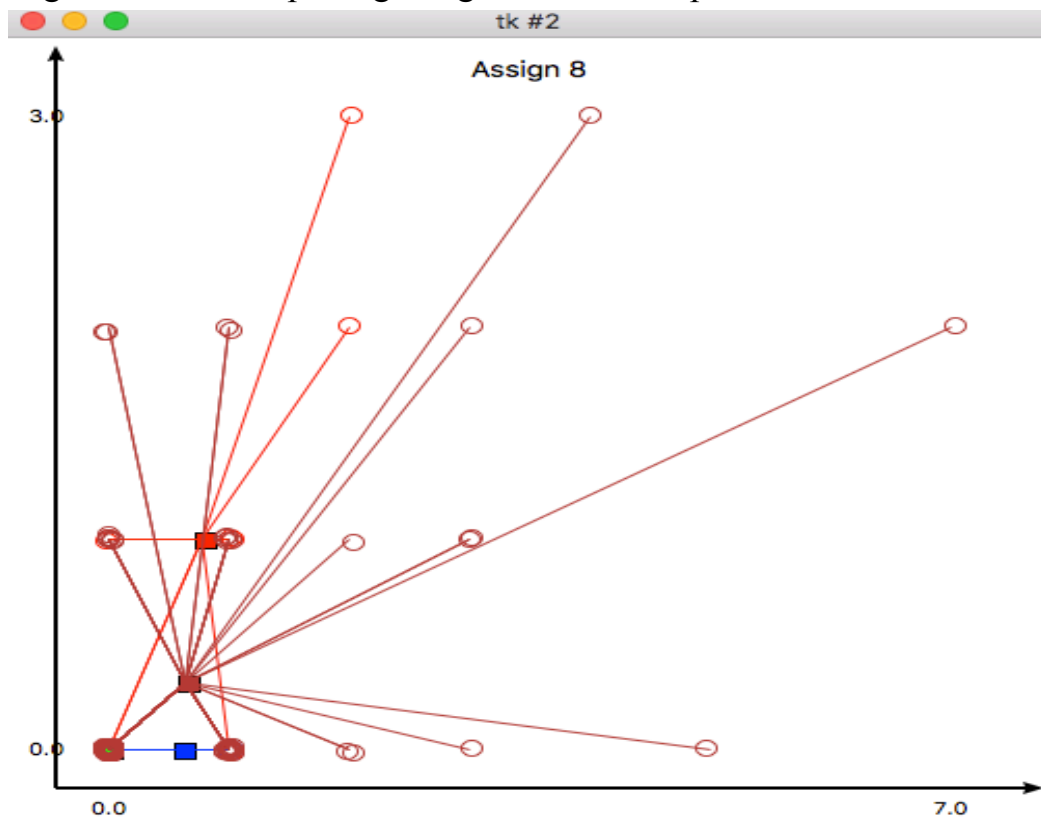
- Sixth iteration depicting assignment of data points to centroids:



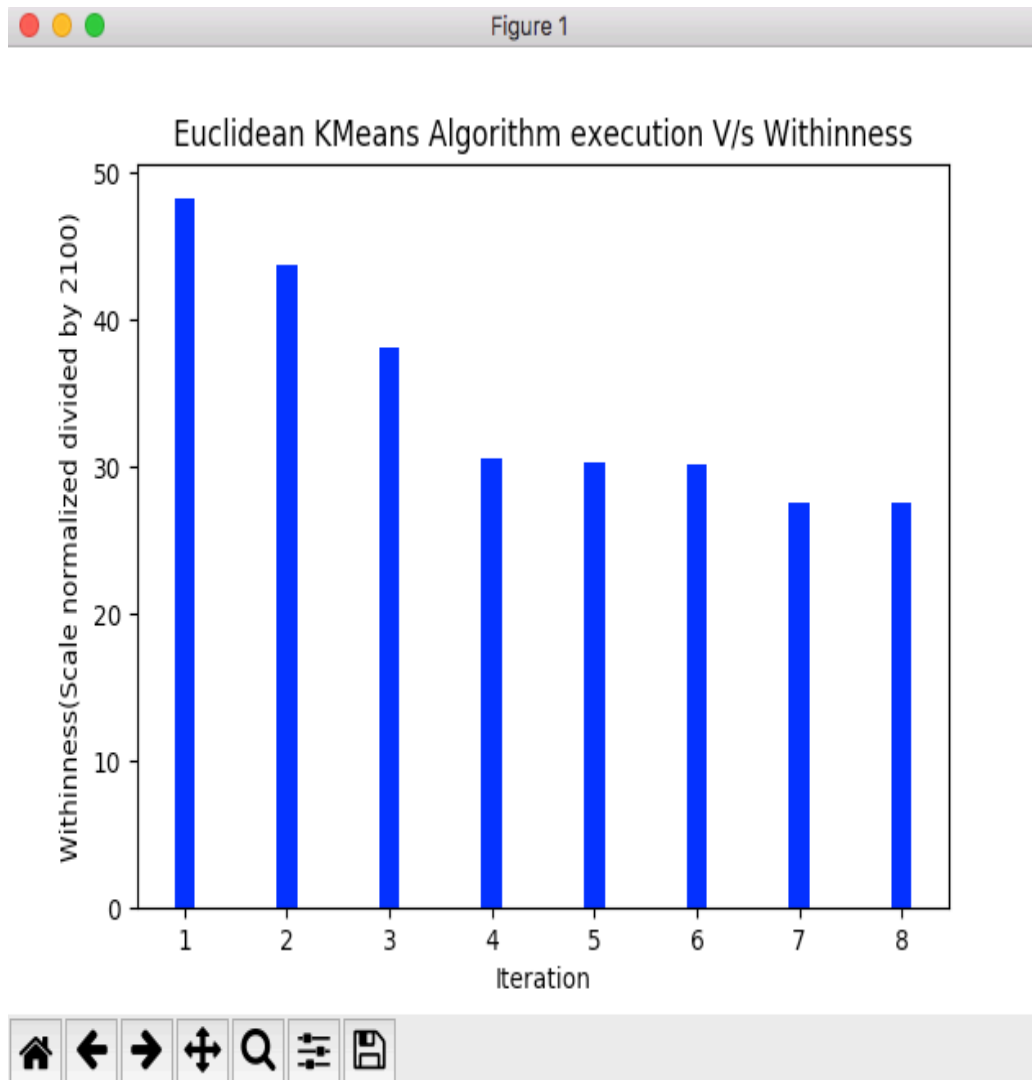
- Seventh iteration depicting assignment of data points to centroids:



- Eighth iteration depicting assignment of data points to centroids:



- The **withinness calculation** for Euclidean KMeans algorithm is displayed below:



- Clearly visible from the above graph the withinness of cluster formation decreases along the successive assignments for Euclidean distance metric as well.

- Centroid allocation: Below are the centroids allocated once the algorithm converges (The text file eucdist.txt contains the centroid allocation information):

```

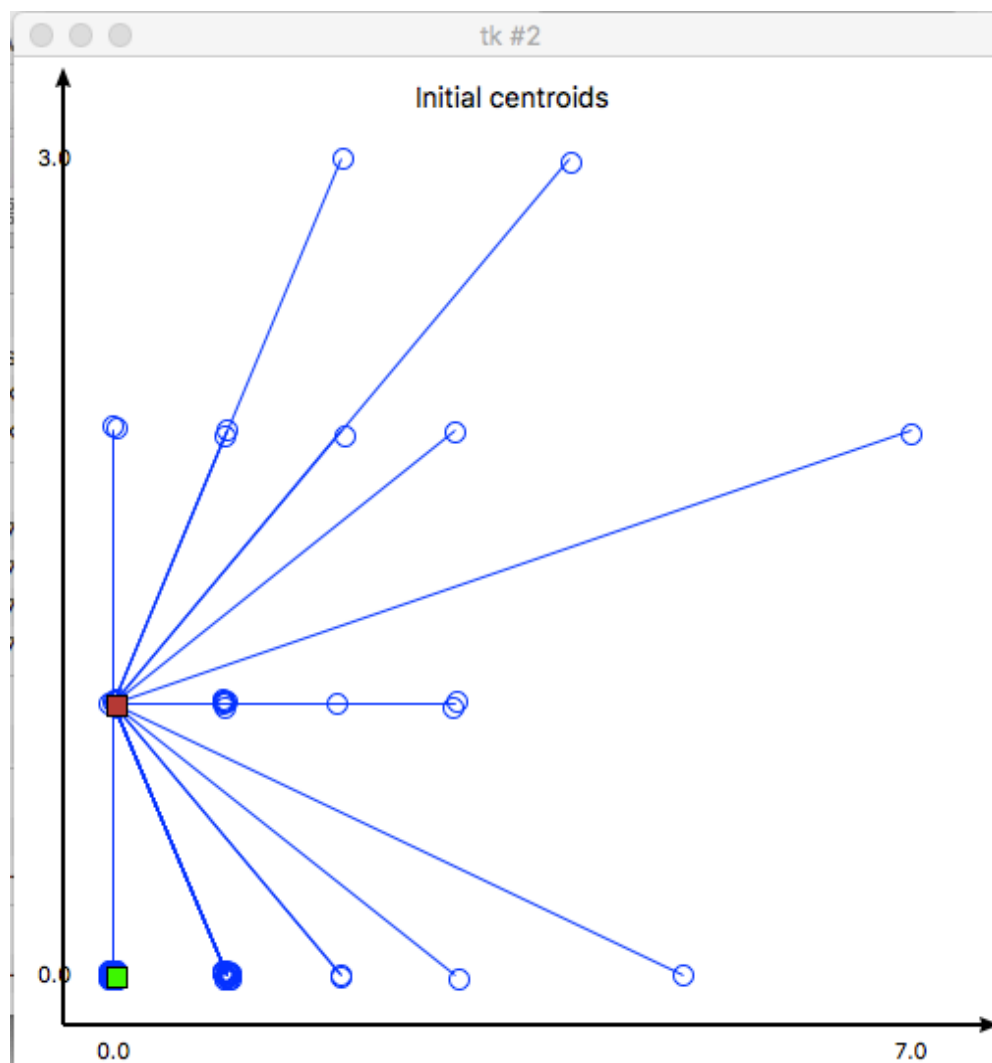
Printing centroid
0      0.60 0.00 15.20 0.00 1.60 1.60 2.80 1.00 1.60 2.20 1.60 1.80 0.00 0.00 0.60 1.00 1.20 1.00 1.00 3.40 1.00 0.00 0.00 0.00 0.60 0.80 2.40 0.40 0.60 0.40
0.00 0.00 0.20 1.00 0.40 0.40 0.40 0.00 0.20 0.80 0.00 0.00 0.00 0.00 0.60 0.00 4.20 2.00 0.20 1.20 1.40 0.00 0.20 0.20 2.20 2.20 1.20 0.00 0.40 0.20 0.40 0.
60 2.80 1.00 4.40 1.00 0.20 0.40 0.00 0.00 0.00 0.00 1.60 0.60 1.40 0.80 2.00 2.40 0.20 1.20 0.00 0.00 0.20 1.20 0.40 0.00 0.80 0.00 0.00 0.00 0.00 0.00
0.20 0.00 0.80 0.20 0.20 0.00 0.00 0.20 0.40 0.20 0.00 0.00 0.20 0.00 0.00 0.40 0.80 0.40 0.80 0.40 0.40 0.40 0.40 0.80 1.00 1.00 0.40 0.60 2.20 1.0
0 1.80 1.00 1.00 0.00 0.40 5.80 26.20 0.80 0.00 0.20 1.20 0.00 0.40 1.80 0.20 0.60 0.20 0.20 0.20 0.40 0.40 1.60 0.00 0.20 0.00 4.40 0.60 17.20 0.40 4.20 0.00
0.00 5.40 0.00 0.00 0.00 0.00 0.40 0.00 0.00 0.00 0.00 0.20 0.00 0.20 0.00 0.80 0.20 0.00 2.20 0.60 0.20 0.40 0.00 0.40 0.00 0.20 1.00 2.00 0.40 1.
40 2.00 1.00 0.20 0.40 0.40 0.00 0.20 0.20 0.00 1.00 0.40 0.00 2.60 1.40 1.60 0.40 1.60 4.20 0.40 2.40 1.00 0.00 0.60 0.80 5.60 1.40 5.40 0.40 0.20 0.80 0.20
0.20 0.20 3.80 0.80 0.80 1.00 0.00 0.80
1      0.78 1.00 10.22 0.00 0.00 0.00 0.00 0.00 0.67 0.11 0.00 0.00 0.00 0.00 0.44 1.00 1.00 1.00 1.33 2.89 2.11 2.56 2.33 2.44 2.22 0.11 1.00 0.33 0.67 0.56
1.11 0.11 0.22 0.00 1.00 0.89 2.56 0.67 0.56 1.56 0.78 0.44 1.11 2.33 3.11 0.67 0.00 2.78 1.78 0.22 0.00 4.44 1.56 0.56 2.67 1.11 2.11 0.33 2.44 0.44 2.89 0.
67 0.11 1.00 1.56 1.11 0.33 0.33 0.00 0.78 1.11 0.00 0.33 1.00 2.89 5.00 0.22 1.44 0.89 0.33 0.00 1.22 0.67 0.67 1.33 0.11 0.00 0.56 0.00 0.00 0.00 0.00 0.00
0.00 0.78 1.22 0.33 0.33 0.11 0.00 0.33 0.44 0.33 0.67 0.00 1.78 0.44 0.33 0.44 0.89 4.33 0.89 1.78 1.00 1.33 0.89 0.89 1.11 1.78 1.22 1.67 0.89 2.67 2.00 1.3
3 1.33 1.00 1.33 3.44 2.89 30.78 3.11 11.00 0.89 0.89 12.44 1.22 0.56 9.11 9.44 1.00 2.22 7.00 2.22 0.78 1.56 1.11 0.11 1.00 0.67 1.67 1.11 1.56 1.00 0.33 0.7
8 0.33 1.44 0.44 0.11 3.22 0.78 0.67 0.56 1.33 0.44 0.78 0.22 1.11 2.78 0.44 0.33 0.33 1.89 1.44 0.89 1.78 2.22 0.56 1.11 1.44 2.22 1.22 0.78 2.33 1.22 0.44 0
.56 1.33 0.89 0.89 1.56 0.44 1.00 0.22 1.33 0.44 0.89 0.22 1.44 1.33 0.56 1.11 0.56 0.22 2.33 0.33 0.11 1.11 0.67 2.33 0.89 3.44 3.56 0.33 0.89 0.89 0.56 0.22
0.44 0.33 0.89 0.67 0.67 0.22 0.11 1.00
2      0.00 0.00 9.50 0.00 2.00 2.00 3.50 1.00 2.00 2.50 2.00 2.00 0.00 0.00 3.00 1.00 1.00 1.00 1.00 4.00 0.00 0.00 0.50 1.50 0.50 0.50 1.50 0.00 1.00 1.00
0.50 0.50 0.50 0.00 1.00 0.00 1.50 0.00 0.00 1.50 1.00 0.00 0.00 0.00 0.00 4.00 0.00 0.50 2.50 0.00 2.00 0.00 0.00 6.00 1.50 2.00 0.00 4.00 0.00 2.00 1.0
0 0.50 1.00 4.50 1.00 0.00 0.50 0.00 0.00 0.00 0.50 0.00 1.50 0.50 1.00 0.50 5.50 0.00 1.50 0.00 0.00 0.50 0.00 1.00 1.00 0.00 1.00 0.00 0.00 0.00 0.50 0.00 0
.00 0.00 0.50 0.00 1.50 0.00 0.00 0.00 0.50 0.00 1.00 0.00 0.00 0.50 0.00 0.00 1.00 1.00 1.00 1.00 1.00 2.50 1.50 1.00 1.00 1.00 2.00 1.00
2.50 1.00 1.00 0.00 0.00 0.00 8.50 1.50 0.00 0.50 0.00 0.50 1.50 0.00 0.00 3.50 0.50 1.50 0.50 1.00 1.00 0.50 0.00 0.00 0.00 0.00 0.50 0.00 0.00 0.00 0.00 0.
00 0.00 0.00 0.00 1.50 0.00 0.00 0.00 0.00 0.00 0.00 0.00 0.00 0.00 0.00 0.00 0.00 0.00 0.00 0.00 3.50 6.50 0.00 2.00 0.00 0.00 0.50 0.50 0.00 5.00 1.00 0.00
3.50 0.50 1.00 0.00 0.00 0.00 0.50 1.00 0.00 1.00 1.00 0.00 4.00 1.50 1.50 0.50 7.00 7.00 0.00 0.50 2.00 0.00 1.00 1.50 19.00 16.50 4.00 1.00 1.50 2.50 0.00 0
.50 0.00 0.00 2.00 0.50 0.50 0.00 0.50
3      0.64 0.32 1.71 0.45 1.42 1.42 1.89 0.73 1.15 1.19 1.30 1.49 0.12 0.07 0.35 0.79 1.33 0.71 1.07 2.81 1.08 0.73 0.82 0.85 0.76 0.24 0.98 0.20 0.49 0.65
0.70 0.17 1.17 0.46 1.07 1.02 1.67 0.04 0.27 0.38 0.19 0.13 0.51 1.30 0.60 0.31 0.43 0.23 0.37 0.43 0.65 0.38 1.13 0.27 1.24 0.88 1.68 0.42 0.40 0.70 1.45 0.3
2 0.62 1.58 0.73 0.24 0.49 0.81 0.26 0.36 0.61 0.13 0.06 0.43 0.76 0.60 0.33 1.08 0.08 0.37 0.38 0.52 0.29 0.52 1.39 0.40 0.13 0.27 0.14 0.18 0.31 1.33 0.31 0
.20 0.18 0.81 0.74 0.36 0.35 0.42 0.43 1.40 0.08 0.10 0.06 0.67 0.51 0.10 0.08 1.07 0.88 1.31 1.55 0.89 0.70 1.10 1.14 1.61 2.25 1.25 0.94 0.75 1.00 1.89 1.14
1.98 1.86 1.95 0.07 0.19 3.45 1.44 0.20 0.39 0.04 0.10 0.10 0.18 0.44 0.05 0.48 0.14 0.32 0.07 0.08 0.24 0.15 0.00 0.29 0.00 0.27 0.33 0.40 0.30 0.20 0.48 0.
14 0.23 0.27 0.19 0.13 0.00 0.23 0.05 0.13 0.46 0.01 0.04 0.13 0.10 0.19 0.19 0.06 0.06 0.07 0.04 0.27 0.36 0.19 0.05 0.35 0.55 0.42 0.32 1.00 1.40 0.23 0.95
1.67 0.51 0.54 0.20 0.12 0.33 0.20 0.20 0.23 0.30 0.10 0.21 0.49 0.21 0.52 0.18 0.69 0.64 0.39 0.40 0.38 0.04 0.73 1.01 0.64 0.79 0.13 0.38 0.14 0.10 0.24 0.5
2 0.58 0.15 0.05 0.37 0.40 0.07 0.58

```

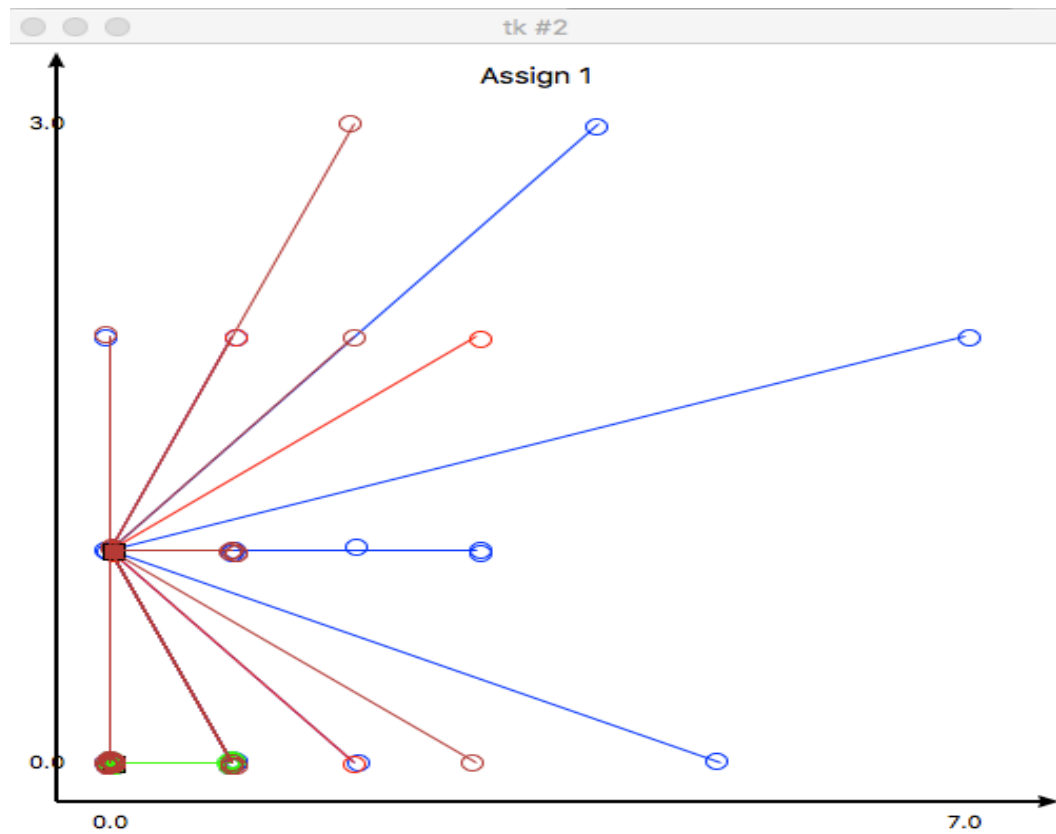
c. KMeans algorithm with Cosine distance: The Cosine distance is used as distance metric with KMeans algorithm. Using Cosine as the distance metric KMeans algorithm converges in 6 iterations.

- Below is the graphical representation of cluster formations:

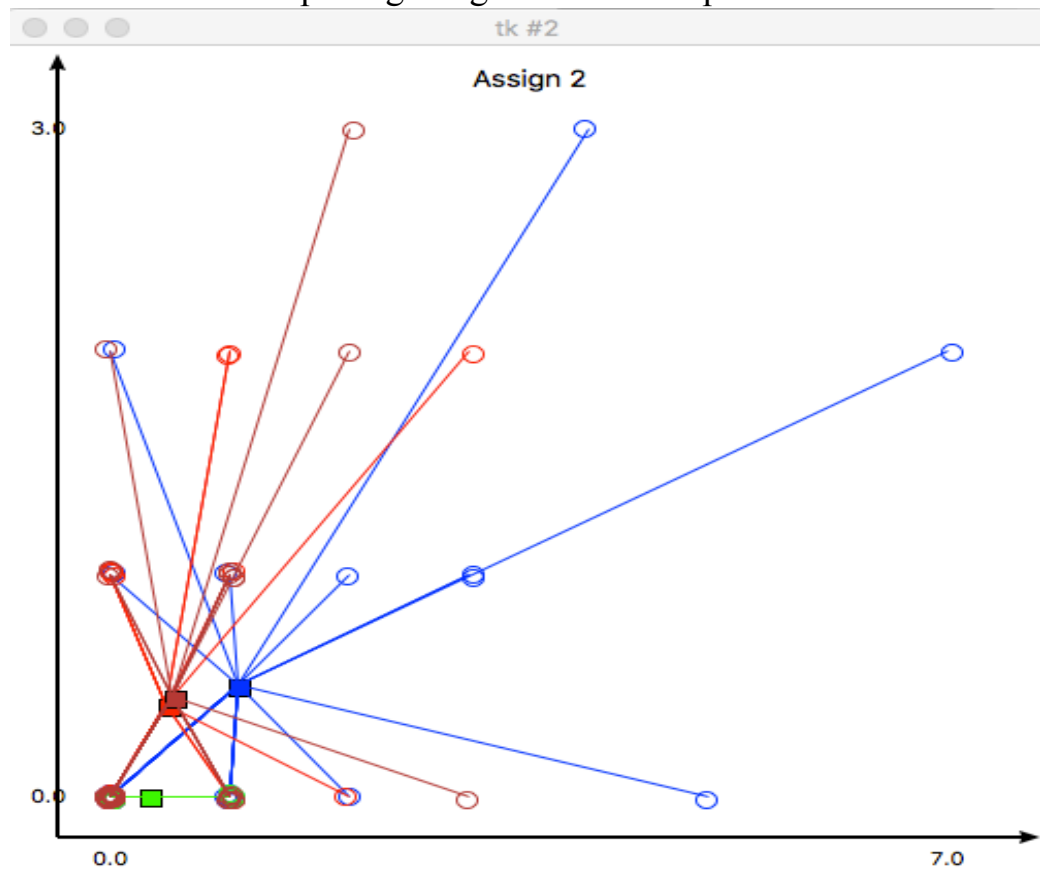
- Initial Centroids:



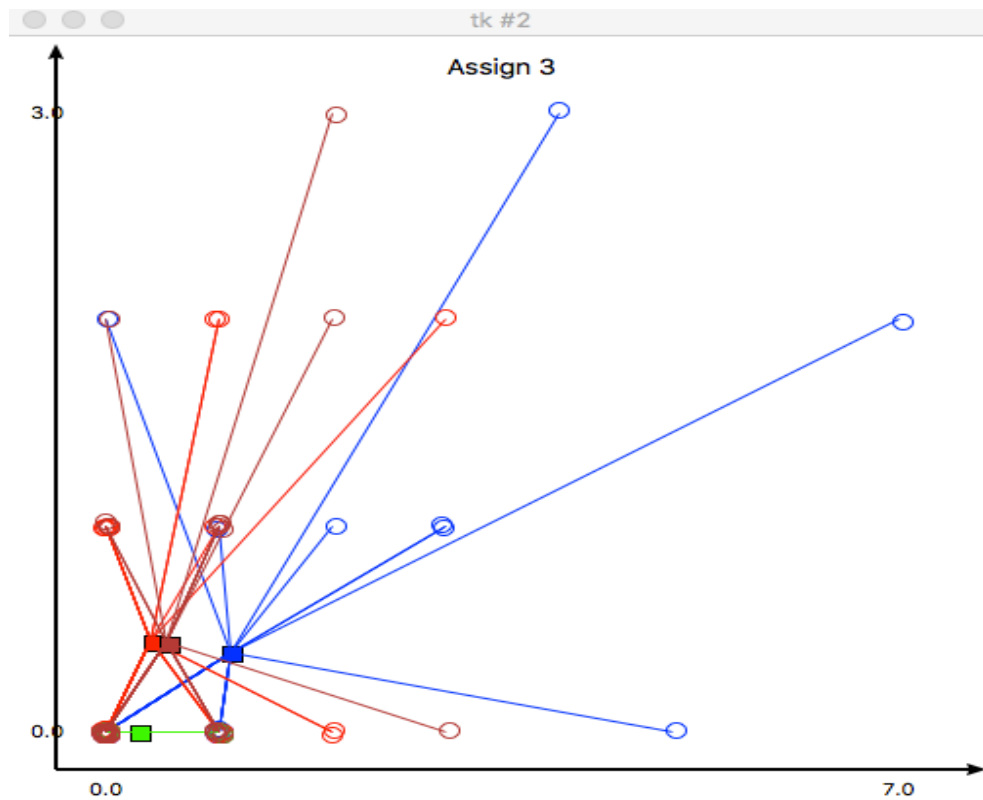
- First iteration depicting assignment of data points to centroids:



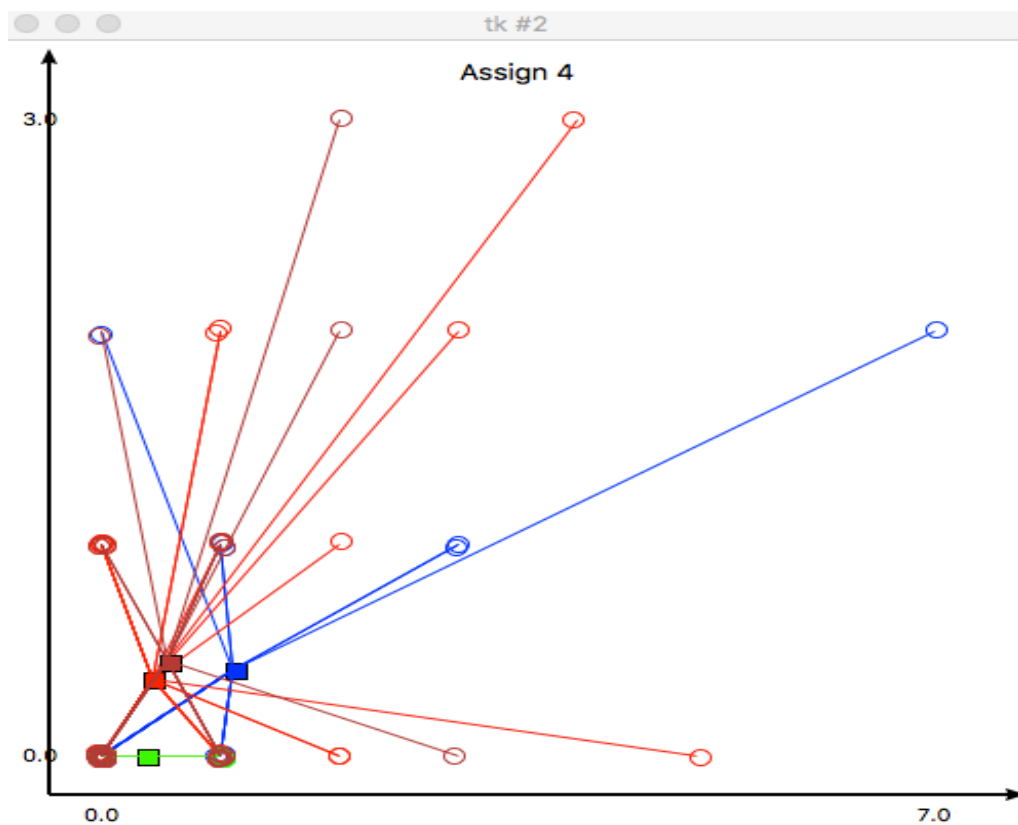
- Second iteration depicting assignment of data points to centroids:



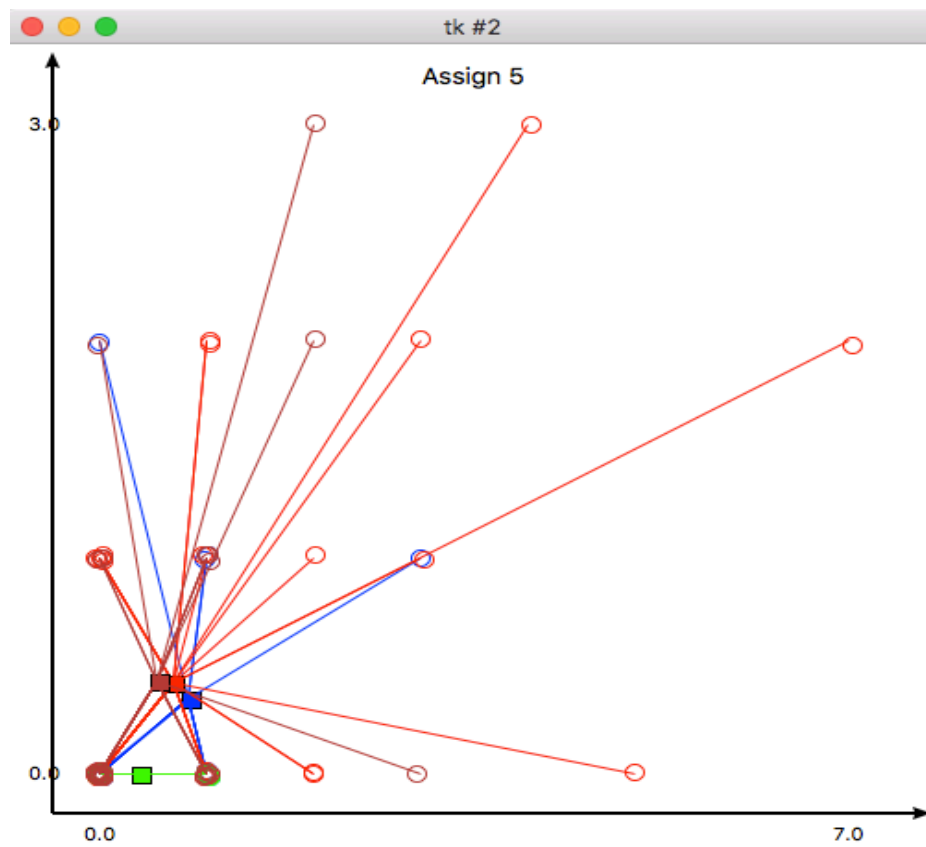
- Third iteration depicting assignment of data points to centroids:



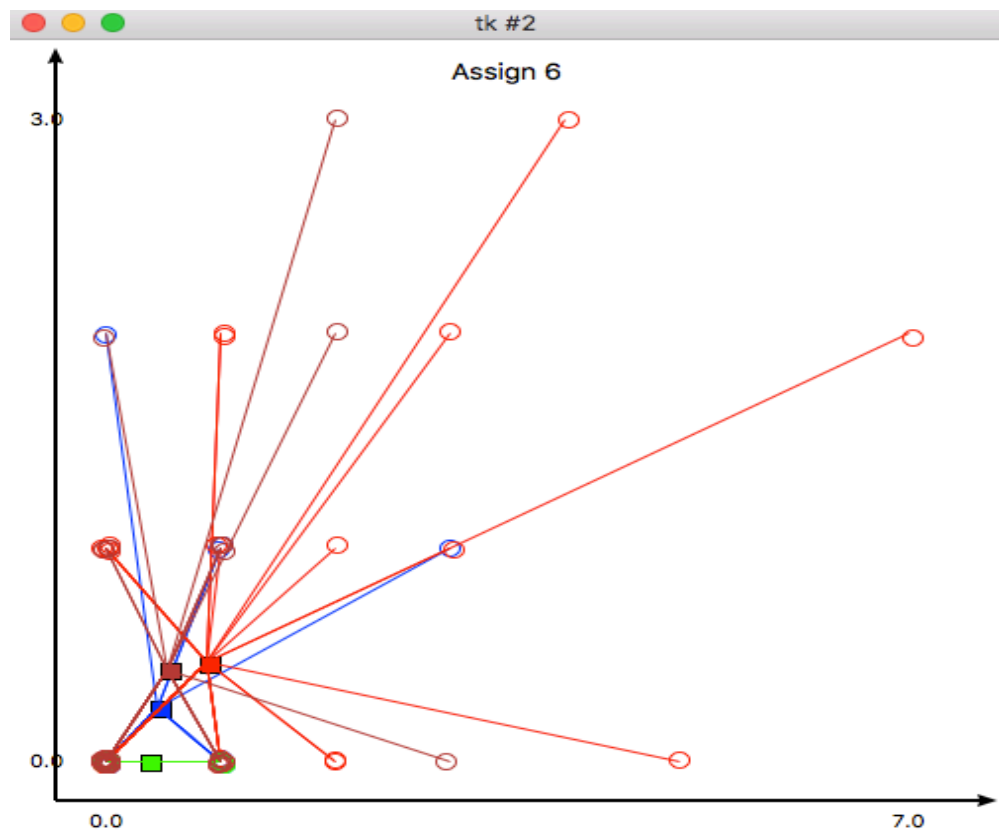
- Fourth iteration depicting assignment of data points to centroids:



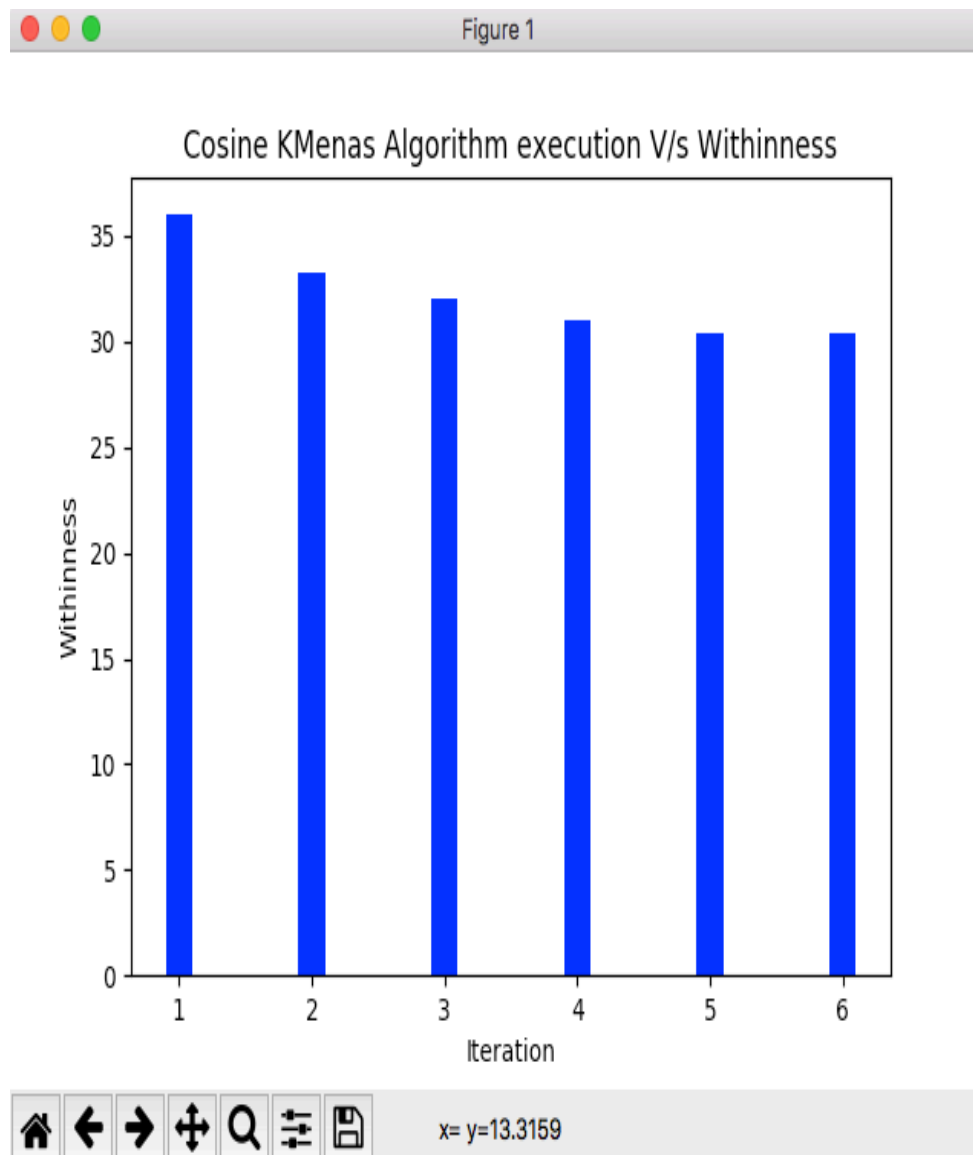
- Fifth iteration depicting assignment of data points to centroids:



- Sixth iteration depicting assignment of data points to centroids:



- The **withinness calculation** for Cosine KMeans algorithm is displayed below:



- Clearly visible from the above graph the withinness of cluster formation decreases along the successive assignments for Cosine distance metric as well.

- Centroid allocation: Below are the centroids allocated once the algorithm converges (The text file Cosine_Points.txt contains the centroid allocation information):

```

Printing centroid
0      0.45 0.25 0.65 0.00 0.00 0.00 0.05 0.55 0.00 0.00 0.60 0.30 0.05 0.00 0.10 0.35 0.70 0.05 0.30 0.75 0.55 0.45 0.35 0.10 0.15 0.10 0.65 0.00 0.35 0.05
0.05 0.00 2.45 0.10 1.10 0.40 1.40 0.00 0.20 0.15 0.05 0.20 0.40 0.90 0.25 0.40 0.05 0.00 0.05 0.05 1.00 0.45 0.55 0.40 0.95 1.25 2.65 0.30 0.65 1.15 0.70 0.0
5 0.10 2.90 0.30 0.05 0.25 0.45 0.15 0.15 0.90 0.00 0.00 0.15 0.35 0.40 0.10 0.65 0.05 0.40 0.25 0.25 0.20 0.40 3.50 0.10 0.05 0.00 0.00 0.50 0.30 0.00 0.20 0
.05 0.05 0.80 1.10 0.10 0.05 0.10 0.05 0.40 0.05 0.05 0.00 0.40 0.95 0.05 0.00 1.15 0.30 2.00 0.20 0.20 0.05 0.15 1.45 3.00 3.05 0.45 0.20 0.05 0.20 1.50 0.20
4.05 3.85 3.85 0.00 0.10 1.85 0.25 0.00 0.15 0.00 0.00 0.05 0.05 0.00 0.05 0.40 0.00 0.30 0.00 0.15 0.05 0.00 0.00 0.15 0.00 0.10 0.20 0.40 0.10 0.40 0.45 0.
20 0.10 0.15 0.15 0.00 0.00 0.05 0.05 0.25 0.85 0.00 0.00 0.00 0.05 0.55 0.00 0.00 0.00 0.00 0.05 0.30 0.10 0.15 0.00 0.65 0.60 0.15 0.10 1.25 1.90 0.05 1.45
1.90 0.50 0.15 0.15 0.00 0.20 0.45 0.00 0.15 0.35 0.00 0.10 0.15 0.00 0.40 0.00 0.00 0.45 0.25 1.35 0.05 0.15 0.90 0.50 0.10 0.35 0.00 0.30 0.00 0.10 0.10 0.5
0 1.25 0.00 0.00 0.30 0.20 0.10 0.40
1      0.88 0.46 1.20 0.90 2.59 2.59 3.39 1.02 2.00 2.02 2.07 2.51 0.22 0.15 0.20 1.00 1.78 1.00 1.51 4.24 1.63 0.88 1.07 1.24 1.12 0.41 1.29 0.34 0.78 1.29
1.05 0.34 0.63 0.59 1.12 1.80 2.24 0.05 0.46 0.24 0.34 0.17 0.73 1.95 0.56 0.37 0.46 0.27 0.37 0.78 0.76 0.15 1.59 0.22 1.29 0.83 1.39 0.61 0.32 0.66 2.15 0.2
4 1.05 1.02 0.98 0.12 0.59 1.27 0.39 0.49 0.46 0.02 0.10 0.61 0.73 0.22 0.34 1.51 0.07 0.32 0.63 0.73 0.22 0.78 0.22 0.71 0.24 0.34 0.27 0.12 0.44 2.66 0.54 0
.37 0.12 1.12 0.68 0.59 0.61 0.71 0.61 2.61 0.05 0.10 0.02 1.05 0.41 0.10 0.05 1.05 1.20 1.05 2.17 1.27 1.00 1.68 1.10 1.00 1.93 1.85 1.24 1.10 1.46 2.10 1.71
1.12 1.00 1.17 0.12 0.07 0.71 1.39 0.29 0.63 0.02 0.10 0.00 0.32 0.24 0.02 0.56 0.17 0.02 0.07 0.07 0.41 0.27 0.00 0.37 0.00 0.20 0.49 0.34 0.49 0.07 0.56 0.
10 0.12 0.41 0.02 0.07 0.00 0.37 0.00 0.02 0.15 0.00 0.05 0.00 0.05 0.07 0.32 0.12 0.00 0.02 0.02 0.15 0.44 0.27 0.05 0.07 0.61 0.54 0.44 1.07 1.32 0.05 0.78
2.17 0.51 0.76 0.17 0.24 0.56 0.05 0.22 0.39 0.24 0.05 0.29 0.59 0.15 0.59 0.02 0.29 0.61 0.59 0.10 0.32 0.00 0.83 1.39 1.05 1.20 0.07 0.51 0.12 0.02 0.41 0.7
3 0.29 0.17 0.10 0.61 0.44 0.02 0.49
2      0.36 0.00 10.55 0.00 1.82 1.82 3.09 1.18 1.82 2.36 1.82 2.00 0.00 0.00 2.18 1.00 1.09 1.00 1.18 3.27 0.45 0.00 0.18 0.36 0.45 0.45 2.09 0.27 0.64 0.36
0.27 0.09 0.18 0.55 0.64 0.18 0.82 0.00 0.09 1.27 0.27 0.00 0.09 0.09 0.45 0.00 3.45 1.18 0.18 1.18 0.73 0.55 0.09 0.09 2.64 1.27 1.36 0.09 0.91 0.27 0.91 0.
73 1.73 1.00 3.55 0.91 0.18 0.36 0.00 0.00 0.00 0.09 0.00 1.00 0.36 0.91 1.09 2.09 1.18 0.55 0.55 0.00 0.27 0.18 0.91 0.55 0.00 0.73 0.09 0.00 0.00 0.27 0.00
0.09 0.09 0.45 0.09 0.36 0.00 0.00 0.09 0.27 0.09 0.18 0.09 0.09 0.18 0.00 0.00 0.64 0.82 0.64 1.27 0.64 0.64 0.73 0.64 0.64 1.36 1.00 0.91 0.64 0.73 2.18 1.1
8 1.64 1.00 1.09 0.00 0.27 2.64 16.00 0.64 0.00 0.18 0.73 0.09 0.45 0.82 0.09 1.09 0.27 0.36 0.18 0.45 0.55 1.00 0.00 0.09 0.00 2.27 0.36 7.91 0.27 2.18 0.00
0.00 2.73 0.00 0.00 0.27 0.00 0.36 0.00 0.00 0.00 0.00 0.09 0.00 0.09 0.00 0.18 0.00 0.45 0.09 0.09 2.27 1.73 0.09 0.55 0.00 0.27 0.09 0.27 0.64 2.82 0.64 0.9
1 1.64 0.82 0.73 0.18 0.18 0.00 0.27 0.27 0.00 0.64 0.36 0.00 2.00 1.00 1.18 0.45 5.55 4.55 0.18 1.45 0.91 0.00 0.64 0.91 6.36 3.91 3.91 0.55 0.36 0.91 0.09 0
.18 0.18 1.91 0.73 0.45 0.64 0.00 0.55
3      0.54 0.43 5.46 0.04 0.18 0.18 0.21 0.07 0.46 0.29 0.14 0.25 0.00 0.00 0.29 0.82 1.07 0.82 1.00 2.21 1.14 1.39 1.36 1.39 1.21 0.07 0.61 0.18 0.21 0.21
0.82 0.04 0.89 0.43 1.00 0.43 1.39 0.25 0.18 0.93 0.25 0.14 0.50 1.11 1.71 0.32 0.25 1.07 1.11 0.11 0.11 1.96 1.18 0.39 1.79 0.89 1.61 0.21 1.07 0.39 1.50 0.6
8 0.14 1.36 0.75 0.75 0.50 0.32 0.11 0.50 0.86 0.36 0.14 0.61 1.89 2.75 0.25 0.96 0.36 0.39 0.04 0.71 0.54 0.32 1.71 0.07 0.00 0.43 0.00 0.00 0.07 0.04 0.00 0
.04 0.54 0.61 0.54 0.25 0.14 0.14 0.46 0.25 0.25 0.32 0.11 0.75 0.36 0.21 0.32 1.04 1.96 1.14 1.68 0.89 0.93 0.86 0.96 1.46 2.11 1.00 1.29 0.79 1.46 1.86 1.00
1.68 1.54 1.64 1.14 1.29 17.89 2.11 3.71 0.43 0.36 4.07 0.64 0.21 3.89 3.11 0.57 0.86 2.96 0.82 0.25 0.50 0.36 0.04 0.54 0.21 0.89 0.50 0.89 0.39 0.21 0.54 0
.25 0.79 0.25 0.46 1.32 0.25 0.25 0.29 0.61 0.71 0.29 0.11 0.75 1.07 0.21 0.18 0.11 0.75 0.64 0.29 0.71 0.96 0.25 0.43 0.93 1.00 0.75 0.46 1.21 0.93 0.61 0.75
0.86 0.61 0.50 0.75 0.14 0.36 0.25 0.71 0.14 0.57 0.29 0.61 0.89 0.57 0.71 0.61 0.32 0.93 0.25 0.04 0.96 0.21 1.00 0.82 1.29 1.39 0.11 0.39 0.54 0.32 0.11 0.
29 0.50 0.43 0.21 0.21 0.46 0.14 1.04

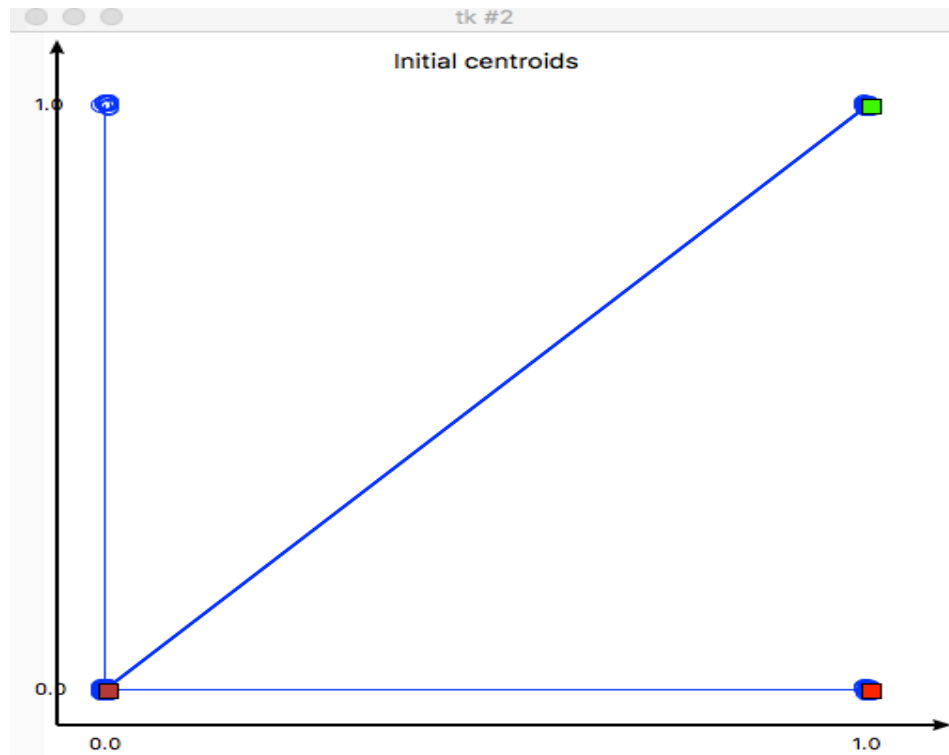
```

- d. KMeans algorithm with Jaccard distance: The Jaccard distance is used as distance metric with KMeans algorithm. Using Jaccard as the distance metric KMeans algorithm converges in 3 iterations, however the results are not at par with other algorithms.

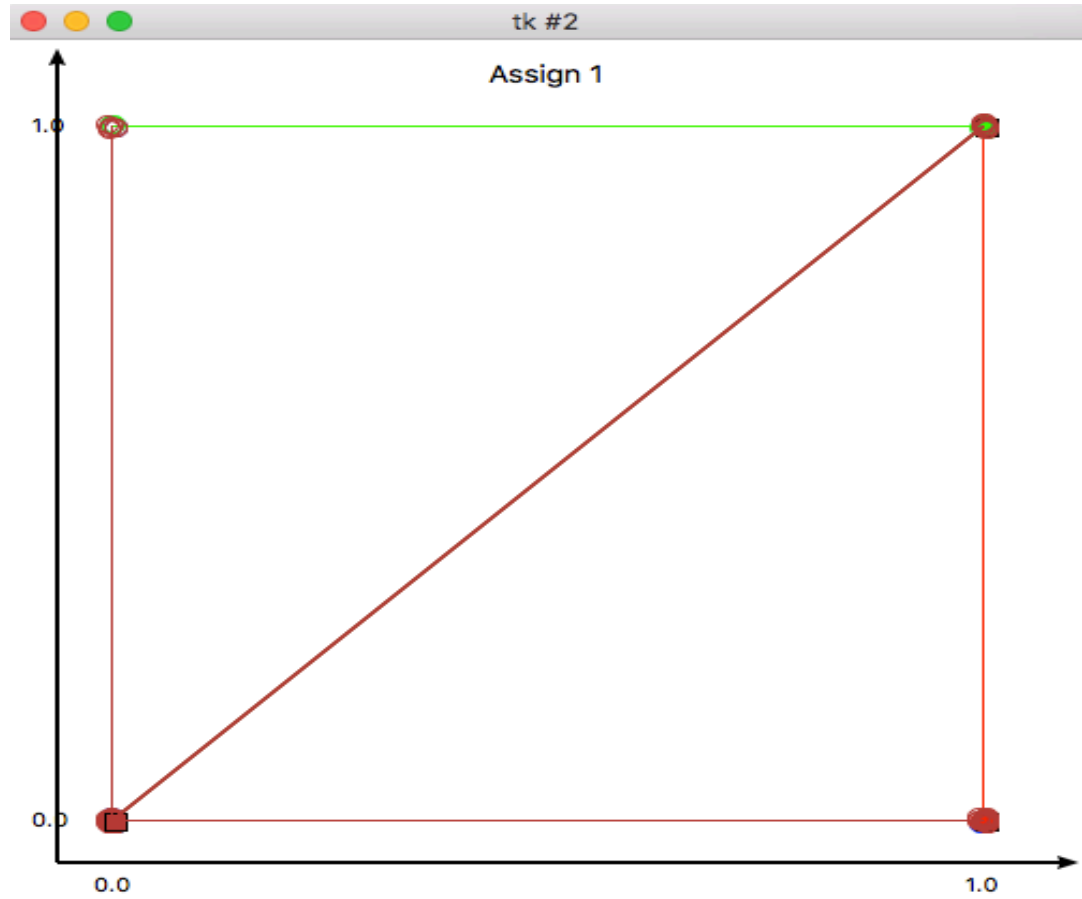
NOTE – For jaccard distance calculation we used similarity matrix rather than the distance matrix as the jaccard similarity works on presence or absence of data.

- Below is the graphical representation of cluster formations:

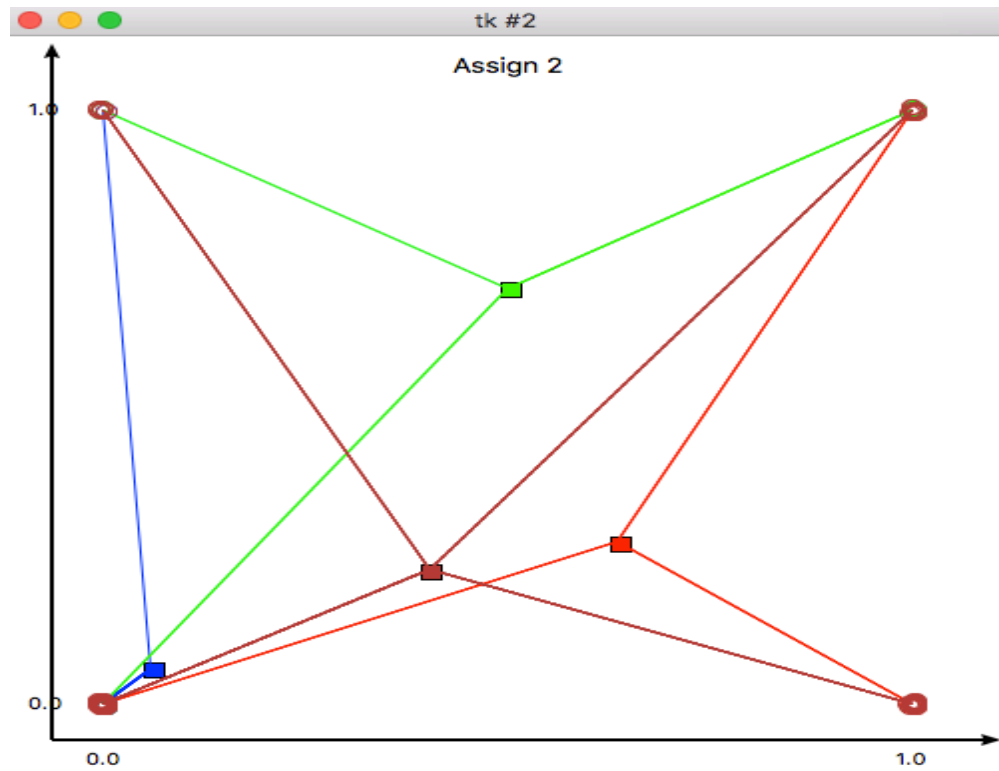
- Initial Centroids:



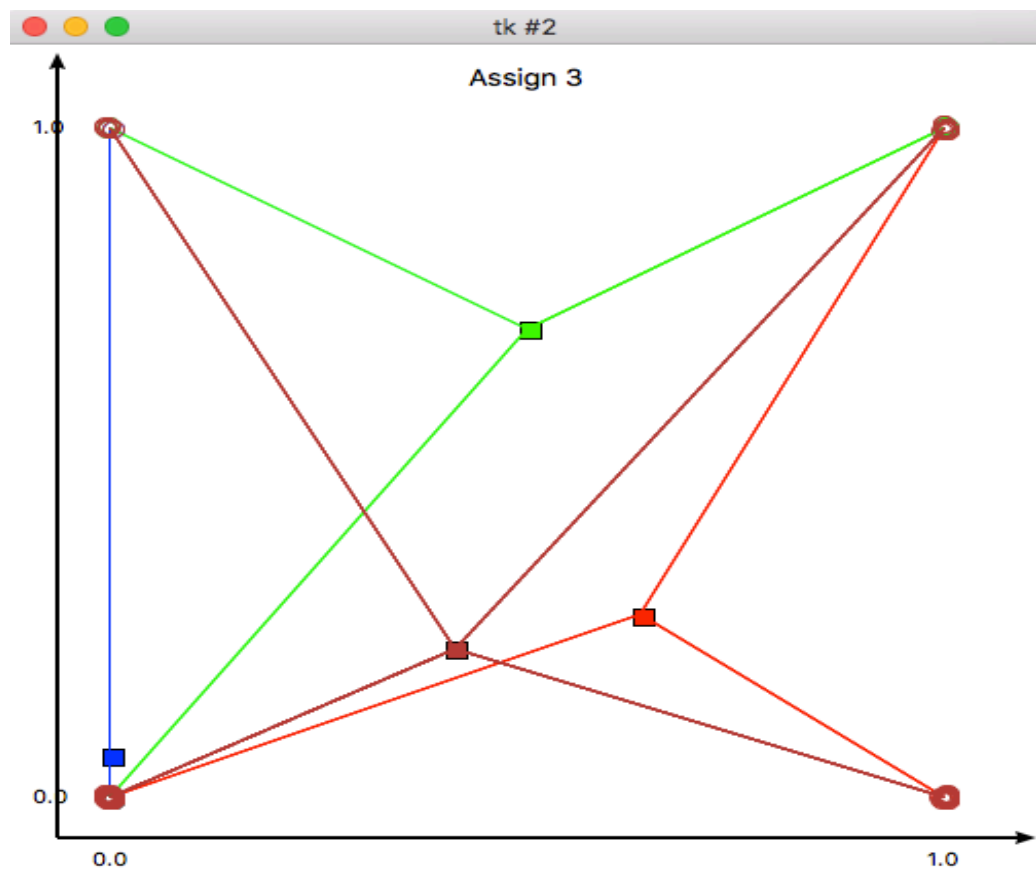
- First iteration depicting assignment of data points to centroids:



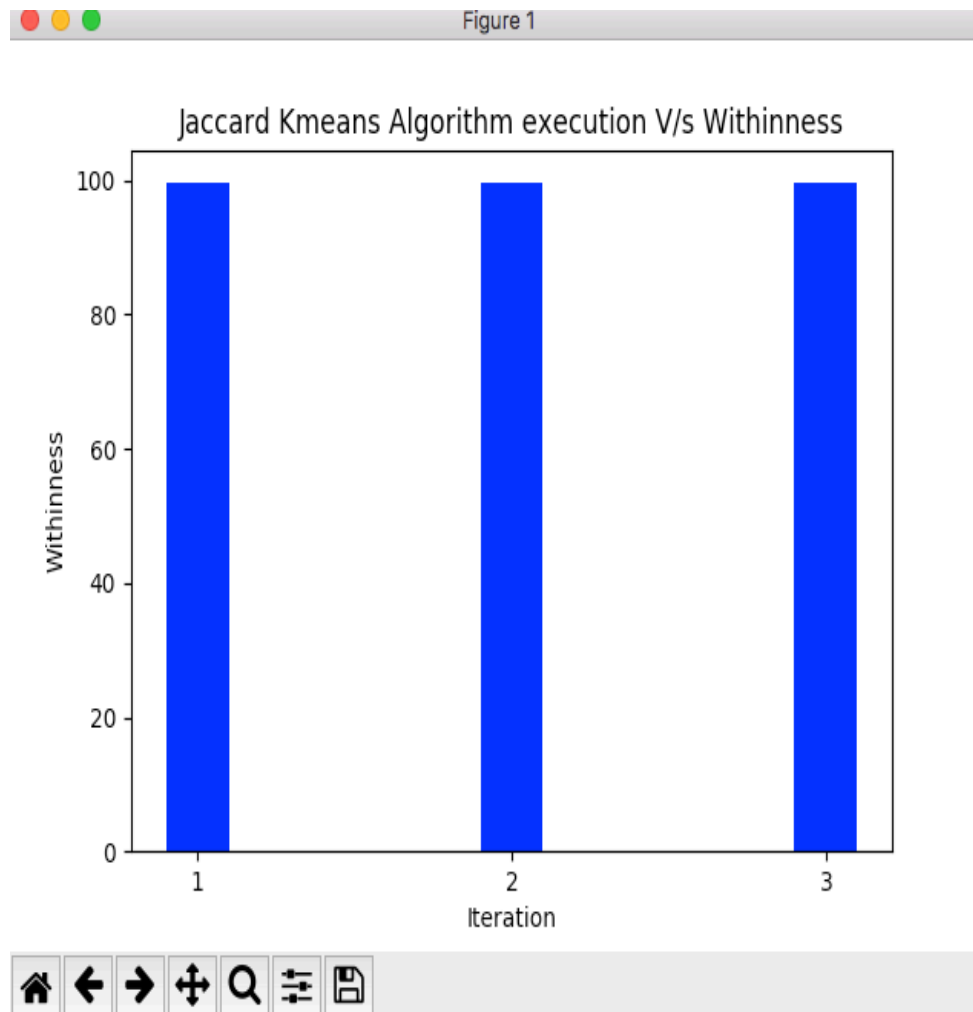
- Second iteration depicting assignment of data points to centroids:



- Third iteration depicting assignment of data points to centroids:



- The **withinness calculation** for Jaccard KMeans algorithm is displayed below:



- Clearly visible from the above graph the withinness of cluster formation does not change along the successive assignments for Jaccard distance metric, indicating that jaccard distance is not a very good measure of distance when implemented with KMeans algorithm.
- Intuitively this seems to be right as well. Jaccard is a similarity matrix which calculates similarity on the basis of presence or absence of features which does not fit well within the concept of KMeans algorithm.

- Centroid allocation: Below are the centroids allocated once the algorithm converges (The text file jacard.txt contains the centroid allocation information):

```
centroid0      0.00 0.06 0.62 0.00 0.12 0.12 0.12 0.19 0.12 0.19 0.12 0.25 0.06 0.00 0.06 0.19 0.50 0.19 0.31 0.50 0.31 0.31 0.12 0.06 0.06 0.06 0.31 0.06 0.
25 0.00 0.00 0.00 0.25 0.12 0.81 0.06 0.62 0.00 0.19 0.25 0.06 0.00 0.19 0.12 0.25 0.25 0.00 0.06 0.06 0.06 0.31 0.25 0.44 0.12 0.75 0.75 0.88 0.06 0.31 0.81
0.06 0.06 0.12 1.00 0.19 0.00 0.12 0.50 0.00 0.06 0.69 0.00 0.00 0.31 0.31 0.44 0.00 0.19 0.12 0.44 0.12 0.19 0.00 0.12 0.88 0.12 0.06 0.00 0.00 0.00 0.00 0.1
2 0.12 0.12 0.06 0.38 0.81 0.06 0.12 0.00 0.06 0.19 0.06 0.06 0.00 0.06 0.69 0.00 0.00 0.81 0.19 0.81 0.00 0.12 0.00 0.12 0.81 0.81 0.81 0.19 0.12 0.00 0.06 1
.00 0.31 1.00 1.00 1.00 0.00 0.19 0.62 0.25 0.00 0.00 0.00 0.00 0.12 0.06 0.12 0.00 0.25 0.00 0.50 0.00 0.06 0.12 0.00 0.00 0.31 0.00 0.25 0.12 0.31 0.12 0.31
0.00 0.19 0.12 0.06 0.12 0.00 0.00 0.12 0.00 0.25 0.31 0.00 0.00 0.00 0.12 0.06 0.06 0.00 0.06 0.06 0.00 0.12 0.19 0.12 0.06 0.50 0.38 0.19 0.19 0.88 0.88 0.
12 0.75 0.88 0.25 0.19 0.12 0.00 0.06 0.25 0.06 0.00 0.12 0.00 0.12 0.31 0.06 0.25 0.00 0.00 0.12 0.06 0.00 0.06 0.06 0.19 0.12 0.12 0.12 0.00 0.12 0.00 0.06
0.00 0.38 0.69 0.12 0.00 0.06 0.25 0.06 0.19

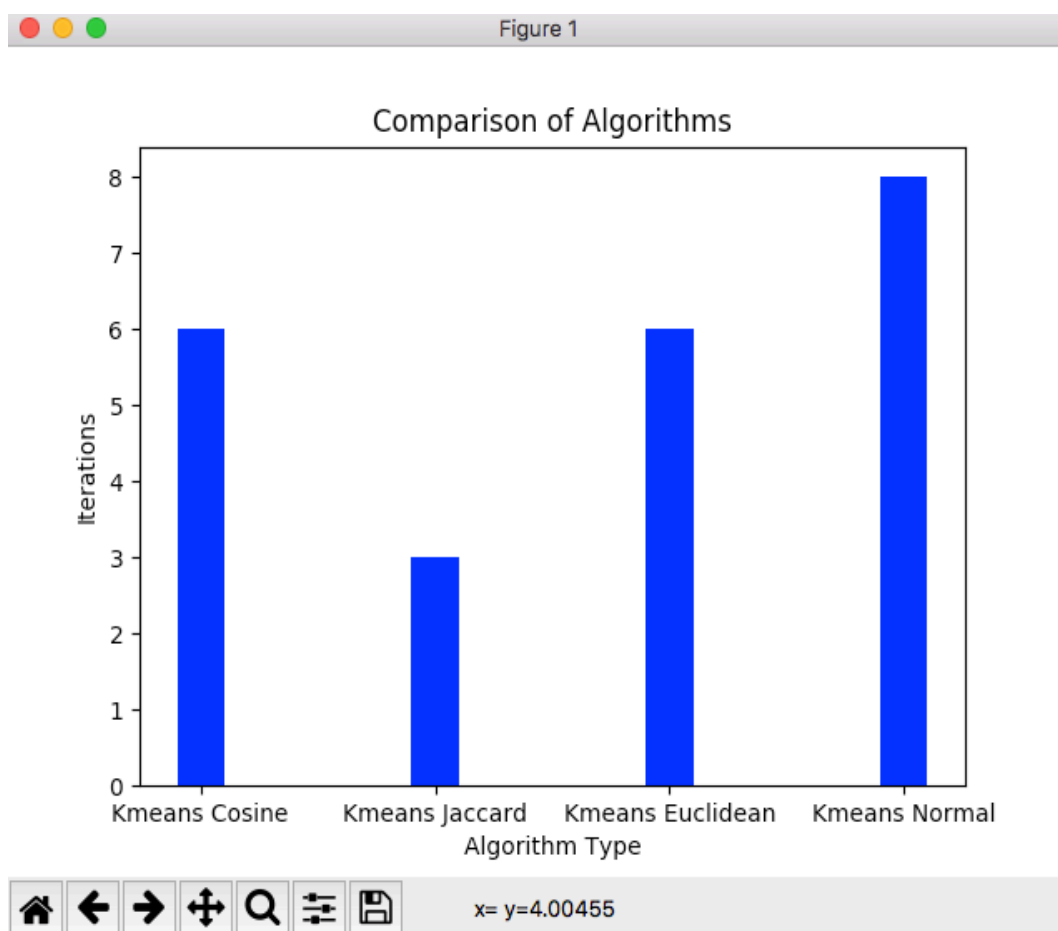
centroid1      0.64 0.27 0.36 0.00 0.00 0.00 0.09 0.18 0.00 0.00 0.27 0.27 0.00 0.00 0.09 0.09 0.36 0.00 0.36 0.45 0.27 0.09 0.09 0.18 0.18 0.18 0.55 0.00 0.
00 0.18 0.00 0.00 0.64 0.00 1.00 0.55 0.91 0.00 0.09 0.18 0.00 0.27 0.36 0.64 0.36 0.00 0.18 0.00 0.09 0.09 0.64 0.18 0.36 0.36 1.00 1.00 1.00 0.36 0.27 1.00
0.64 0.18 0.00 1.00 0.18 0.18 0.27 0.09 0.18 0.18 1.00 0.09 0.00 0.27 0.09 0.18 0.18 0.45 0.00 0.09 0.18 0.36 0.18 0.45 1.00 0.09 0.00 0.09 0.00 0.18 0.18 0.0
0 0.18 0.09 0.18 0.55 1.00 0.09 0.00 0.27 0.09 0.09 0.00 0.00 0.00 0.55 1.00 0.09 0.00 1.00 0.27 1.00 0.27 0.18 0.00 0.00 1.00 1.00 1.00 0.55 0.18 0.00 0.27 1
.00 0.09 1.00 1.00 1.00 0.00 0.27 0.45 0.27 0.00 0.09 0.00 0.09 0.09 0.00 0.18 0.18 0.36 0.00 0.00 0.00 0.09 0.00 0.00 0.00 0.00 0.09 0.00 0.36 0.18 0.27
0.64 0.09 0.09 0.09 0.36 0.00 0.00 0.00 0.18 0.00 0.45 0.00 0.00 0.09 0.00 0.00 0.09 0.00 0.09 0.00 0.09 0.00 0.09 0.18 0.00 0.18 0.64 0.09 0.00 1.00 1.00 0.
00 1.00 1.00 0.36 0.18 0.18 0.00 0.27 0.27 0.00 0.09 0.00 0.09 0.09 0.27 0.00 0.36 0.09 0.00 0.27 0.27 0.18 0.09 0.09 0.55 0.45 0.09 0.27 0.00 0.36 0.00 0.09
0.27 0.18 0.82 0.09 0.00 0.27 0.18 0.09 0.45

centroid2      0.50 0.70 1.00 0.10 0.10 0.10 0.10 0.10 0.20 0.10 0.10 0.10 0.10 0.00 0.30 1.00 1.00 1.00 1.00 1.00 0.90 0.80 0.90 0.90 1.00 0.20 0.40 0.10 0.
40 0.20 0.90 0.20 0.20 0.10 0.80 0.40 0.90 0.00 0.30 0.90 0.30 0.30 0.70 0.80 0.90 0.40 0.00 0.40 0.80 0.30 0.00 0.60 0.80 0.50 0.80 0.70 0.90 0.20 0.80 0.40
1.00 0.50 0.20 1.00 0.60 0.40 0.20 0.40 0.00 0.30 0.70 0.00 0.10 0.70 1.00 1.00 0.30 0.60 0.30 0.40 0.00 0.40 0.50 0.30 1.00 0.20 0.00 0.30 0.00 0.00 0.10 0.1
0 0.00 0.10 0.50 0.50 0.20 0.40 0.10 0.00 0.20 0.30 0.30 0.50 0.00 0.60 0.40 0.30 0.10 0.80 1.00 0.90 0.80 1.00 0.80 0.80 0.90 0.80 0.90 0.80 0.90 0.80 1.00 1
.00 1.00 1.00 1.00 1.00 0.50 0.60 1.00 0.90 0.90 0.60 0.40 0.70 0.40 0.60 0.90 0.70 0.20 0.80 0.70 0.60 0.40 0.70 0.40 0.10 0.50 0.30 0.50 0.60 0.70 0.60 0.20
0.50 0.20 0.50 0.50 0.10 0.80 0.40 0.50 0.20 0.50 0.40 0.30 0.30 0.50 0.70 0.50 0.20 0.40 0.60 0.50 0.50 0.70 0.60 0.40 0.50 0.30 0.70 0.60 0.50 1.00 0.80 0.
40 0.50 0.70 0.40 0.40 0.60 0.30 0.30 0.20 0.60 0.20 0.50 0.20 0.70 0.50 0.40 0.60 0.40 0.10 0.70 0.50 0.20 0.50 0.40 0.80 0.60 0.80 0.80 0.20 0.60 0.50 0.40
0.20 0.40 0.20 0.40 0.20 0.40 0.30 0.10 0.80

centroid3      0.41 0.22 0.57 0.41 0.78 0.78 0.78 0.78 0.81 0.83 0.78 0.78 0.06 0.03 0.22 1.00 1.00 1.00 1.00 1.00 0.56 0.25 0.68 0.75 0.71 0.17 0.57 0.19 0.
32 0.29 0.81 0.10 0.44 0.29 0.97 0.37 0.67 0.06 0.16 0.22 0.11 0.08 0.33 0.48 0.33 0.11 0.27 0.17 0.21 0.32 0.29 0.19 0.46 0.08 0.63 0.35 0.60 0.25 0.27 0.27
0.97 0.22 0.25 1.00 0.40 0.24 0.22 0.38 0.14 0.25 0.21 0.08 0.08 0.33 0.37 0.29 0.29 0.57 0.13 0.25 0.17 0.33 0.14 0.35 0.38 0.29 0.11 0.22 0.05 0.05 0.19 0.1
6 0.13 0.08 0.16 0.44 0.27 0.22 0.37 0.38 0.46 0.41 0.06 0.06 0.05 0.35 0.14 0.10 0.05 0.95 0.97 0.97 0.97 0.97 0.97 0.97 0.97 0.97 0.97 0.97 0.98 0.98 0.97 0.98 1
.00 1.00 1.00 1.00 1.00 0.02 0.11 0.40 0.90 0.16 0.21 0.06 0.08 0.03 0.19 0.35 0.05 0.21 0.17 0.13 0.10 0.10 0.16 0.16 0.00 0.25 0.00 0.27 0.27 0.29 0.19 0.14
0.22 0.11 0.25 0.21 0.03 0.06 0.00 0.17 0.03 0.03 0.06 0.03 0.05 0.06 0.11 0.06 0.10 0.03 0.05 0.05 0.03 0.16 0.32 0.13 0.11 0.06 0.33 0.27 0.25 0.35 0.51 0.
16 0.41 0.49 0.27 0.41 0.17 0.08 0.21 0.11 0.14 0.11 0.17 0.06 0.16 0.44 0.17 0.46 0.11 0.24 0.30 0.21 0.13 0.22 0.00 0.35 0.43 0.46 0.46 0.14 0.27 0.11 0.08
0.17 0.29 0.16 0.16 0.06 0.25 0.24 0.03 0.43
```

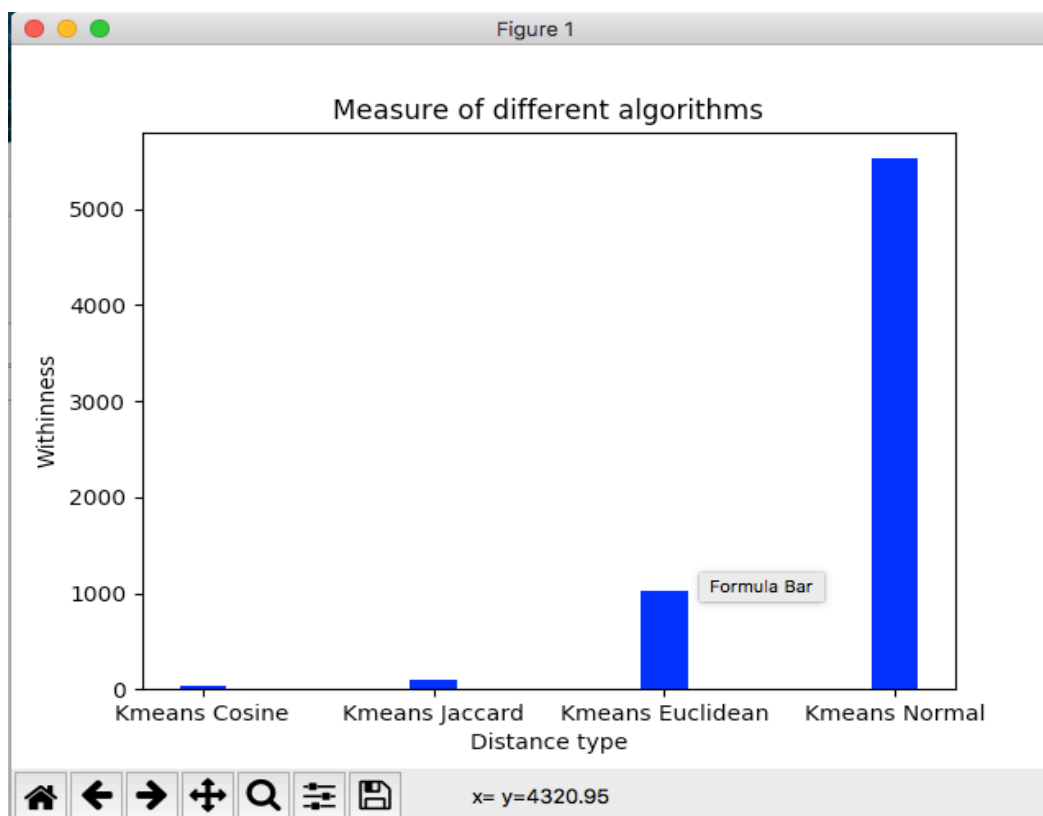
8.Conclusion:

- On the basis of number of iterations taken to converge KMeans: The number of iterations taken for converging KMeans algorithm are similar for Cosine distance metric with Kmeans and normal Kmeans algorithm, Jaccard KMeans algorithm converges a lot earlier then all the other metric distances and Euclidean KMeans taking the largest number of iterations.
- The graphical presentation of the same is as below:



- We see the lowest withinness for Cosine metric, then for jaccard, then Euclidean and highest for kmeans_original.

- We use Sum of Squared Error formula to calculate the withinness and report it in the form of graph below. Higher SSE means lower accuracy, lowest SSE means maximum accuracy.
- Below is the graphical representation for withinness for different algorithms, the accuracy will be opposite of the below, meaning Cosine KMeans will have highest accuracy whereas Normal KMeans will have lowest accuracy.



- As per the analysis Cosine is the best distance metric to run KMeans algorithm, then comes Jaccard, then Euclidean and final is the kmeans normal algorithm