# Lab 3 - Parallelizing k-means Stat 215A, Fall 2017

## Melody Liu

## October 23, 2018

- File list: Lab3.R: R file with similarity function that returns the Jaccard coefficient of two binary matrices and the file also contains calculation of similarity score with 100 iterations and maximum number of k equals to 10.

  Similarity_output_R.csv: I stored the result after running lab3.R file. The result is a 100 times 9 data frame, where row indicates all iterations and column indicates the number of clusters.

  Lab3_C.cpp: CPP file with similarity function.

  lab3_C.R: I sourced CPP similarity function and the file also contains calculation of simialrity score with 100 iterations and maximum number of k equals to 10.

  Similarity_output_C.csv: I stored the result after running lab3_C.R file. The result is a 100 times 9 data frame, where row indicates all iterations and column indicates the number of clusters.
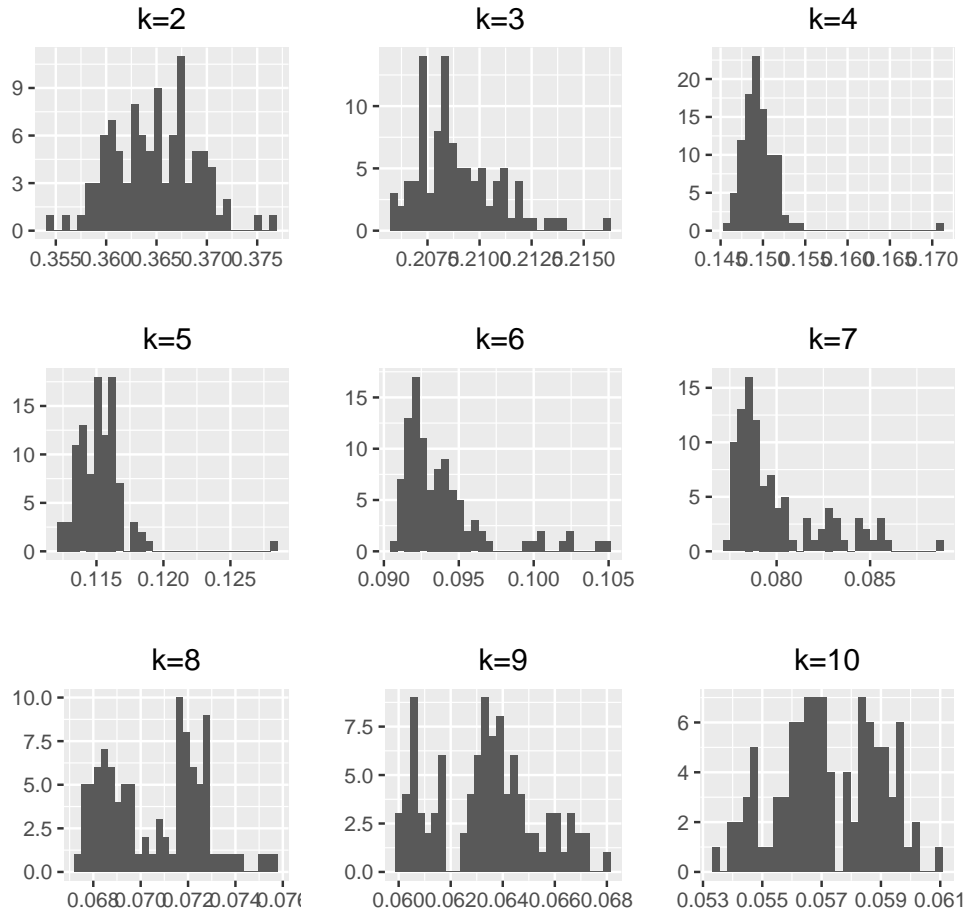
- Compare CPP. and R versions of the similarity matrix Due to the limit of nodes available on server, I subsetted first 10,000 rows of data from lingBianry data set and ran the code on my local machine. It took me around 23 mins to finish running the R code. The similarity score output is saved into a csv file. It took me 17 mins to finish running the code with similarity function in CPP. In CPP, I used a nested for loop to calculate N01, N11 and N10 instead of storing matrices. CPP is a compiled language while R is an interpreted language thus it is much faster for CPP. to run loops than R. The two methods yielded similar results: the similarity matrix generated from two methods are in line with each other.

  I also ran the full data set with similarity function in CPP. It took me around 7 hours in my local machine. Results and plots will be shared in the report.
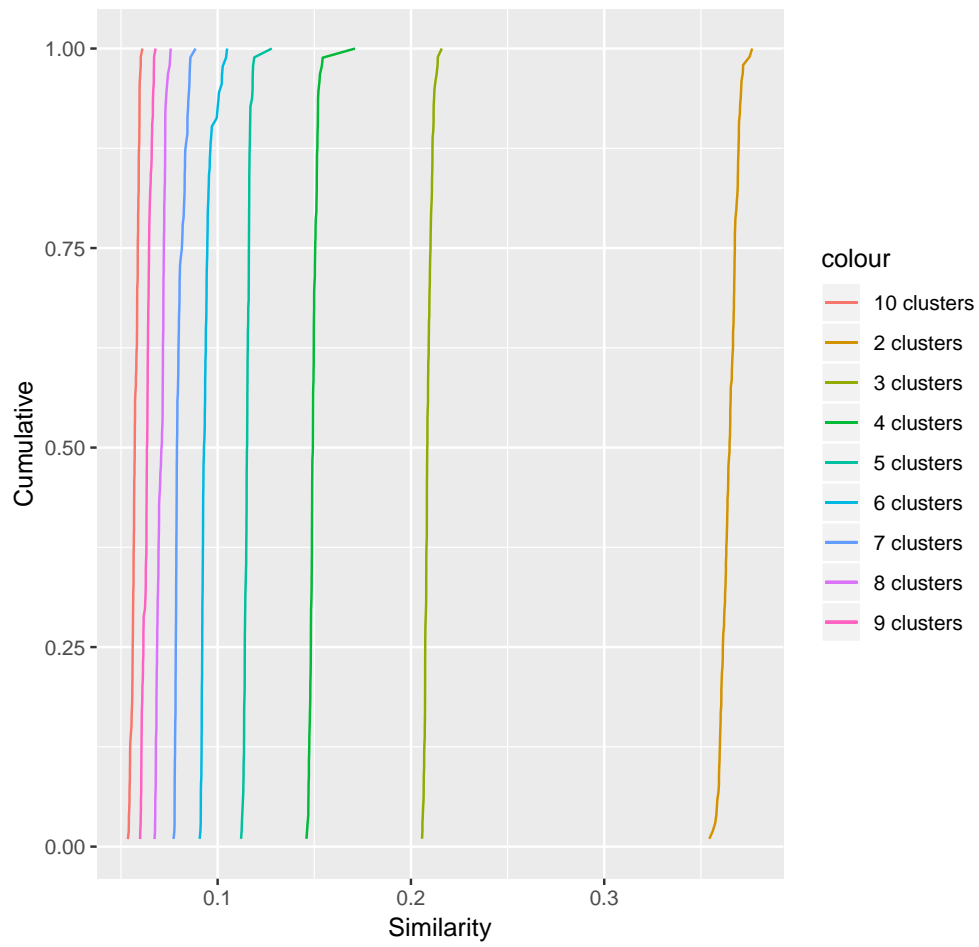
- make the plot for part 3

  The first two plots are for similarity function in R and last two plots are for similarity function in CPP.
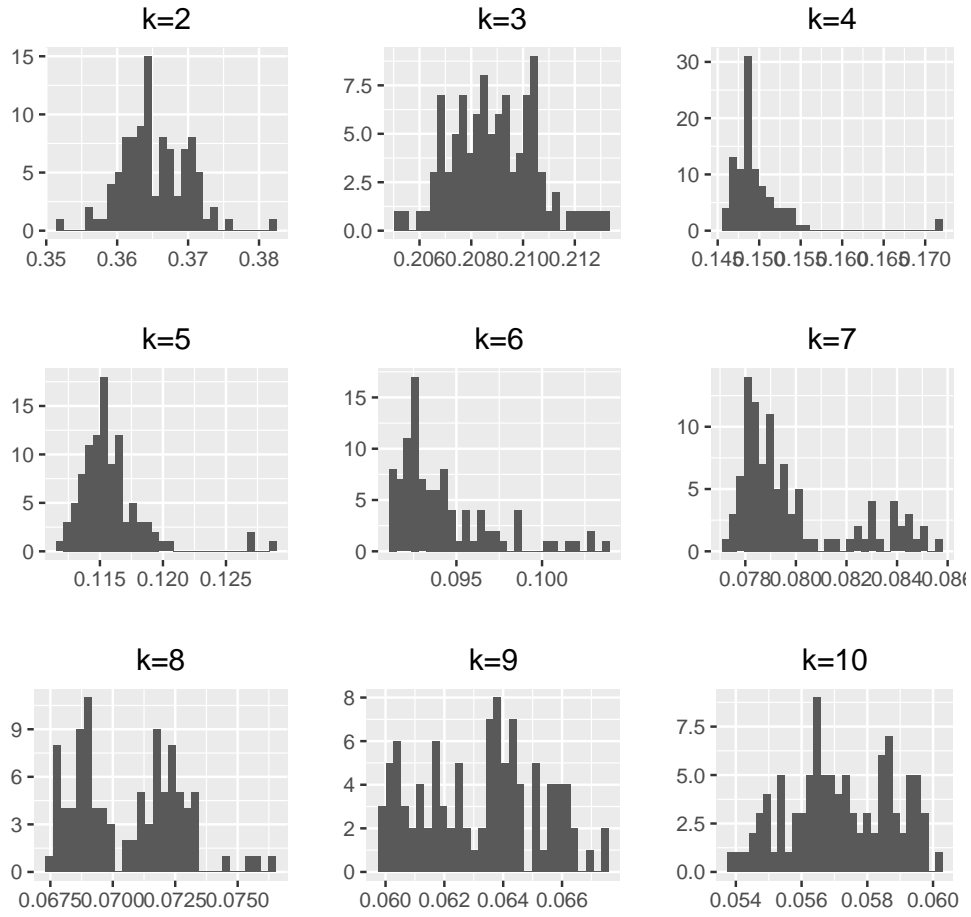
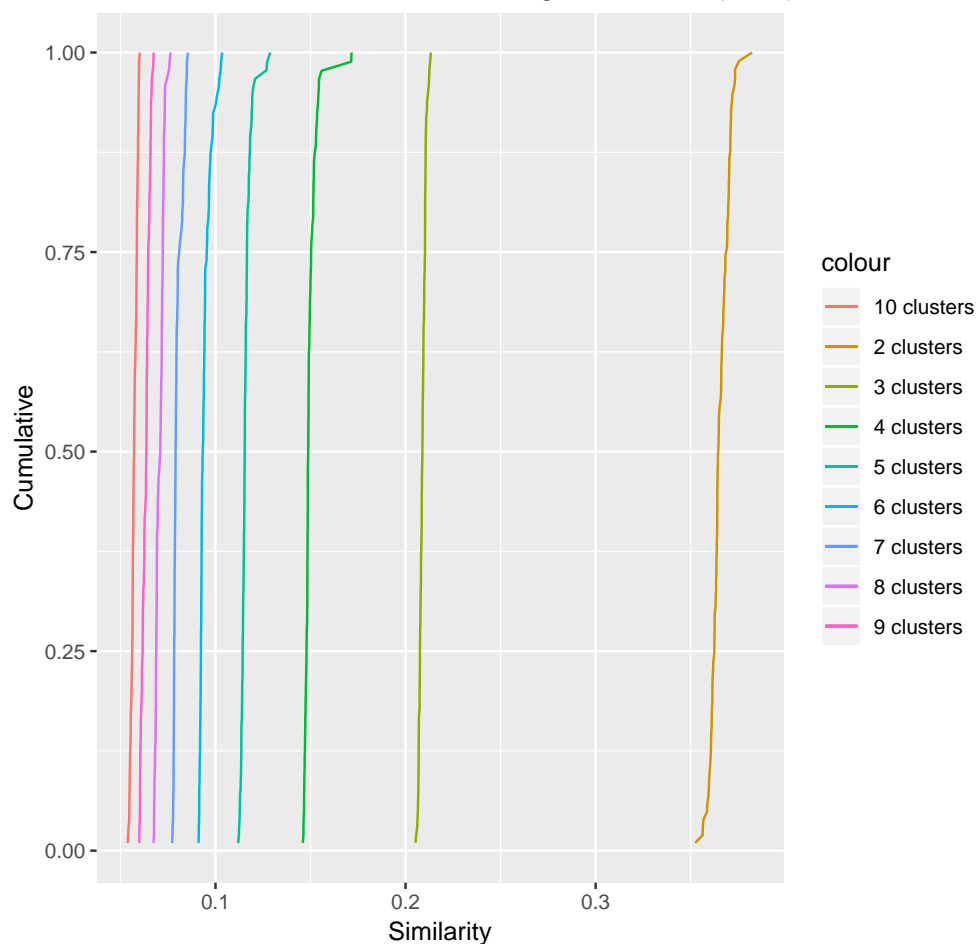## Histogram of the correlation similarity measure (R)

Cumulative distributions for increasing values of k (R)....

# Histogram of the correlation similarity measure (CPP)

## Cumulative distributions for increasing values of k (CPP)....



- discuss how many clusters to choose

  From the cumulative plot, we can conclude that similarity score is higher and more robust when k is smaller. The highest similarity score is achieved when k=2. When k=2, similarity scores mostly fall into the range of 0.35 to 0.375, which is relatively the highest among all k values. Thus k=2 is the optimal value for this data set.

- discuss whether you trust the method or not

  I trust this clustering method, as my last lab report showed 3 noticable cluster centers. Thus k=2 falls into the reasonable range. I also adopted both R and CPP functions to generate similarity coefficients and both gave me similar results.