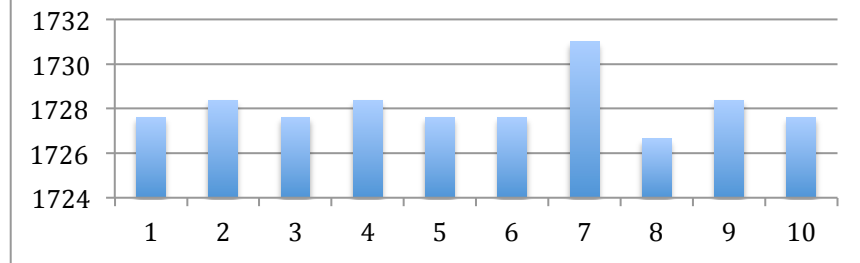
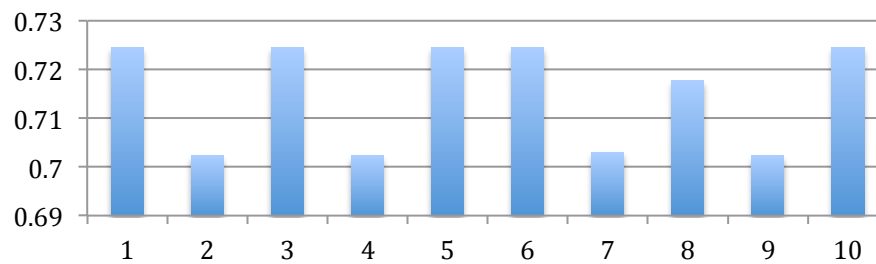
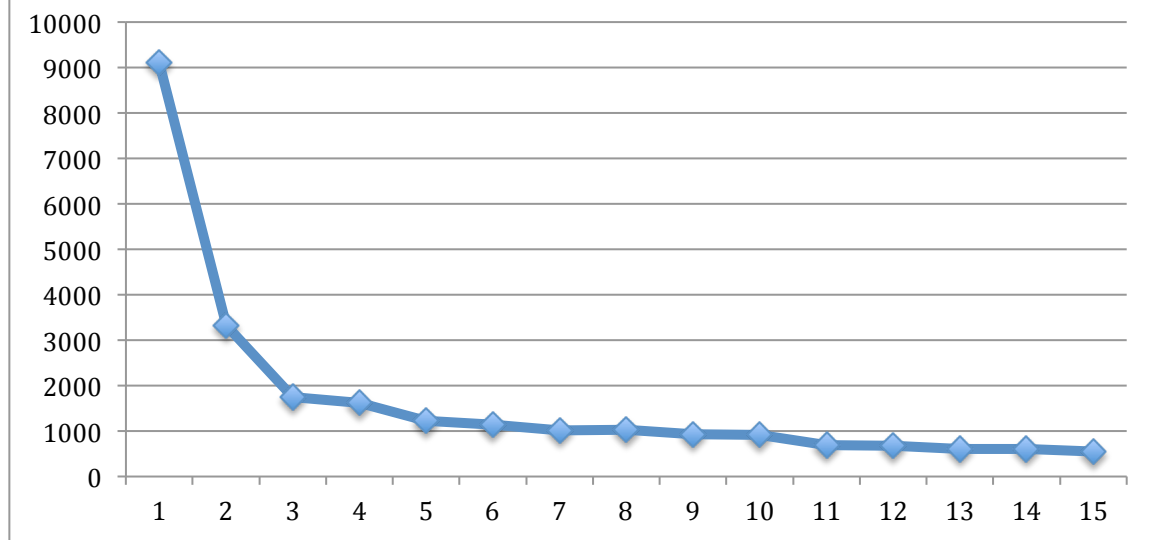
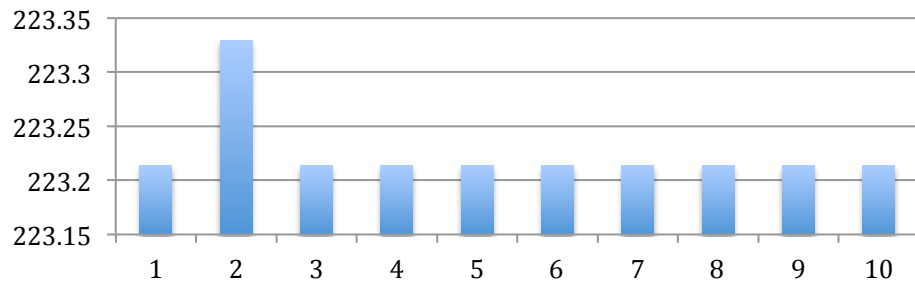
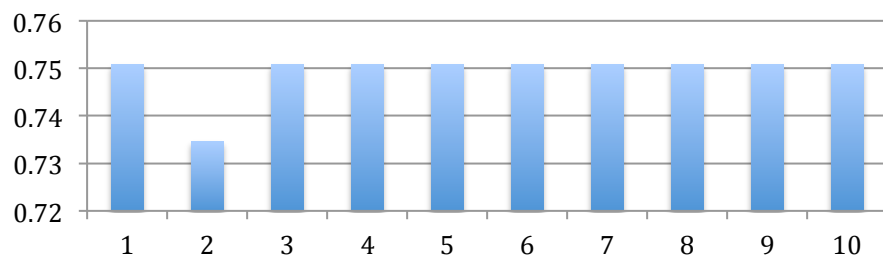
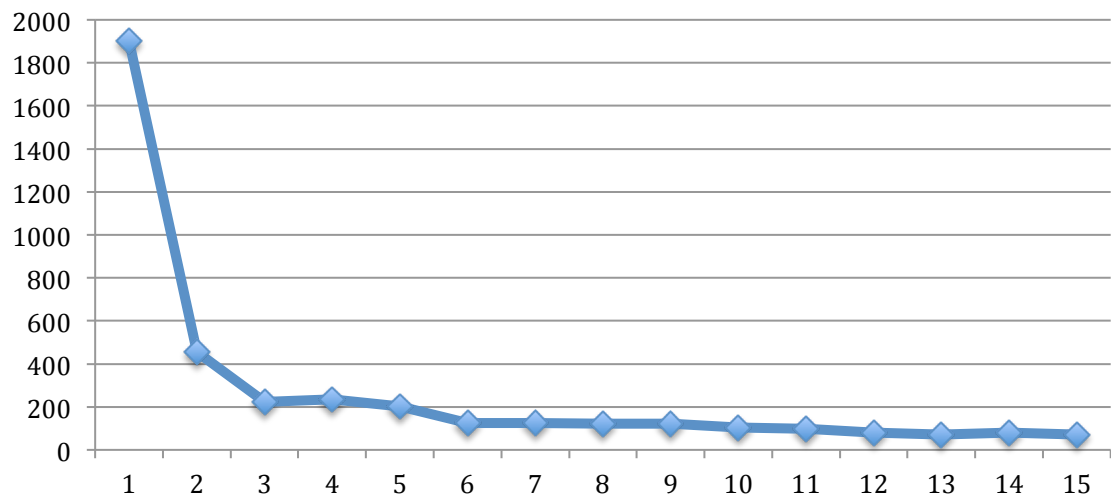


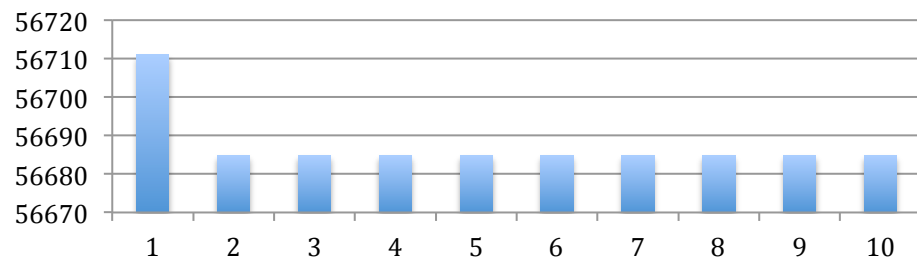
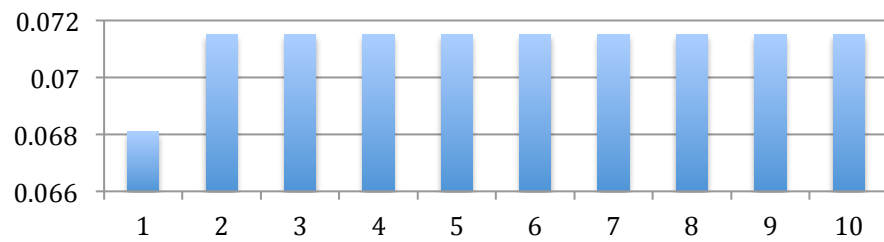
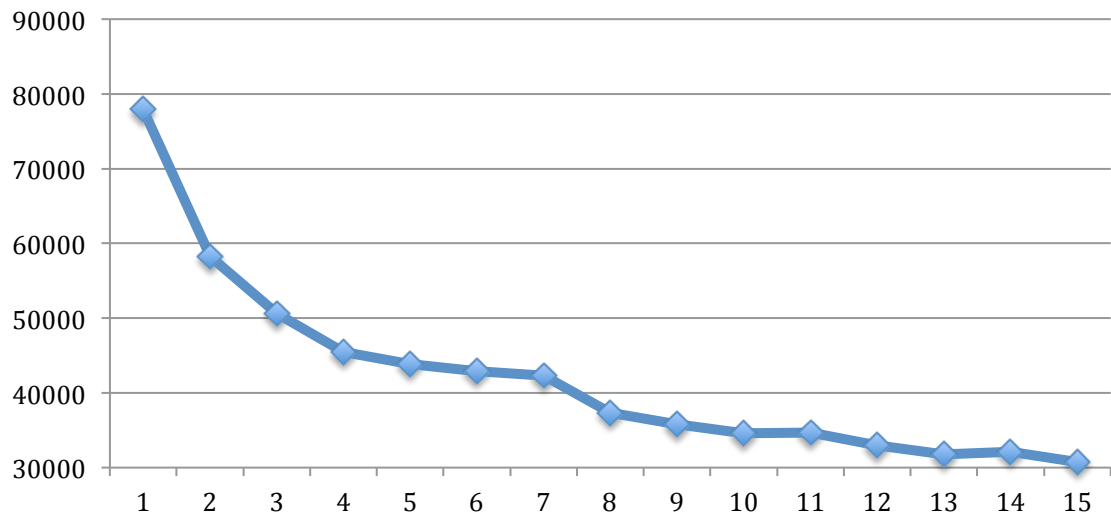
## SEEDS.ARFF

**TotalCS****NMI****Cluster Scatter as function of K**

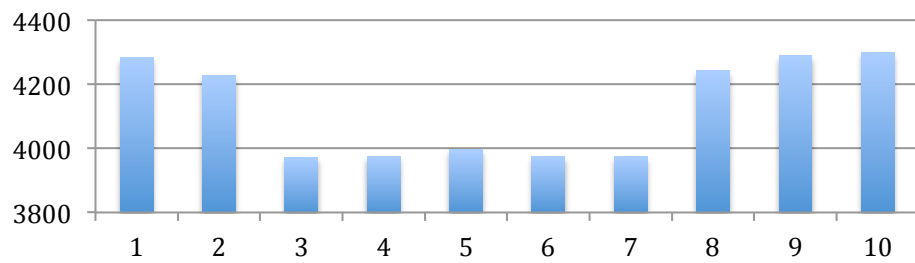
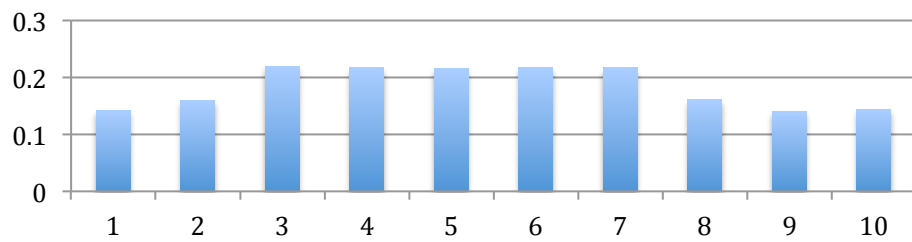
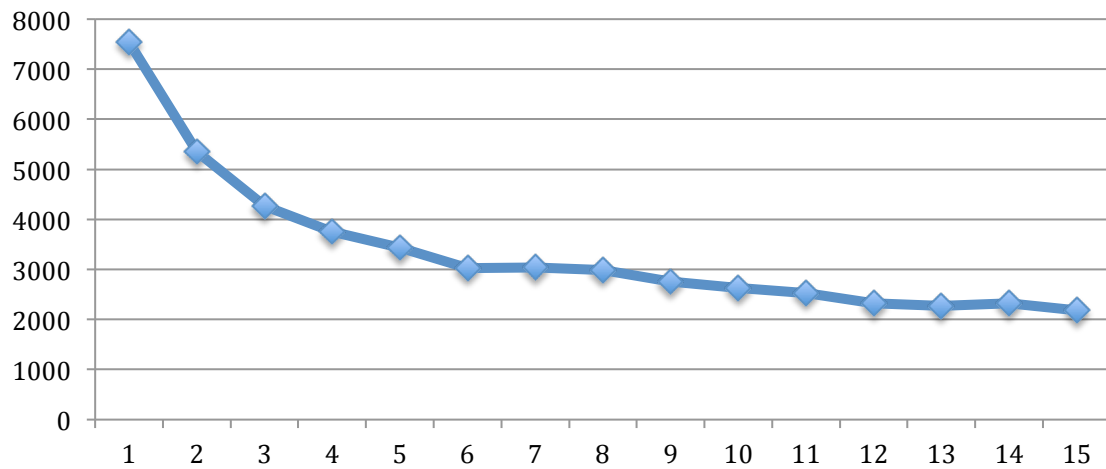
## IRIS.ARFF

**Cluster Scatter over 10 runs****NMI over 10 runs****Cluster Scatter as function of K**

## IONOSPHERE.ARFF

**TotalCS over 10 runs****NMI over 10 runs****Cluster scatter as function of K**

## ARTDATA.ARFF

**TotalCS over 10 runs****NMI over 10 runs****Cluster scatter as function of K**

## Experiment Results:

### Selecting K:

In this experiment, we wanted to see how the value of  $k$  affects the CS of the clusters. For all datasets, we see a really nice curve that shows the decrease of CS as  $k$  improves, as we expect it to. Then they all plateau eventually, and the values of CS start leveling out. We could use these lines to find the value of  $k$  that gives us a good CS without sacrificing performance (aka keeping  $k$  small). Iris.arff is a good example of this. We see a sharp decline in CS up until 3 and then it generally levels out. Therefore, picking  $k = 3$  or possibly  $k = 6$  would be a good choice.

### Sensitivity to initialization:

I was more surprised by the values I obtained from this experiment. We were testing the affects of initializations on CS and NMI, but what I saw from most of the datasets is that regardless of the initialization, the values eventually settled down and remained almost constant. I expected a bigger variation, but the values are very similar. The only one that shows any variation is the artData.arff. This is the dataset that was artificially generated and includes a high level of noise. Therefore, we see that the variation was much higher for NMI and CS if the data is noisy. Otherwise, the NMI and CS should remain relatively stable for multiple initializations.