**Homework 5**

1. **A**

   (a) The digits are as shown below:

   

   

(b) The examples in 2D are:

2. **B**

    (a)    i. Dataset 1

Dataset 1

ii. Dataset 2

Dataset 2

Dataset 2

iii. Dataset 3

Dataset 3

(b) To decide the value of k for each dataset, we can look at the elbow point in the plots of WC-SSD versus k, because beyond this value o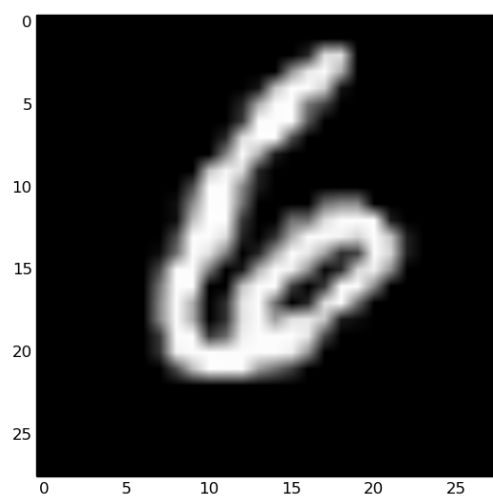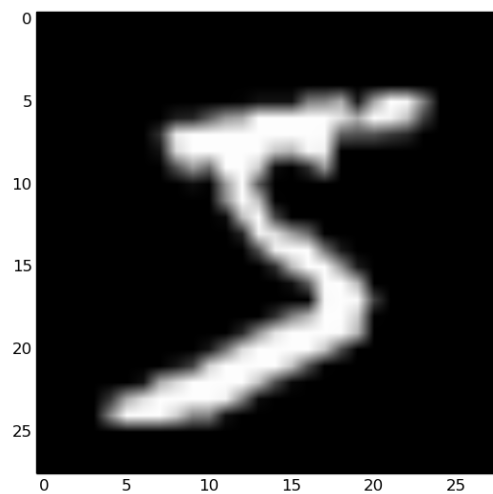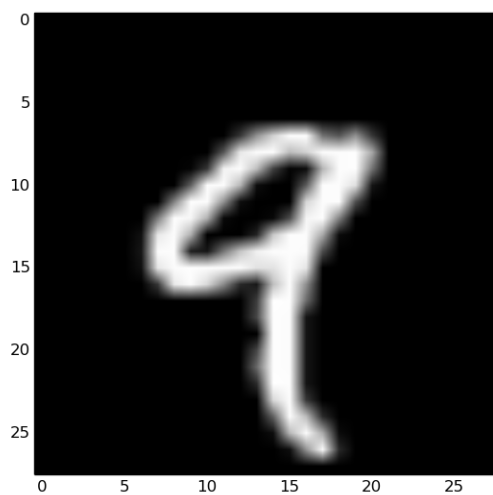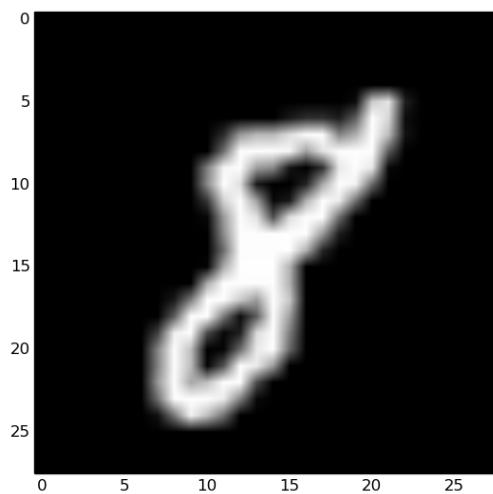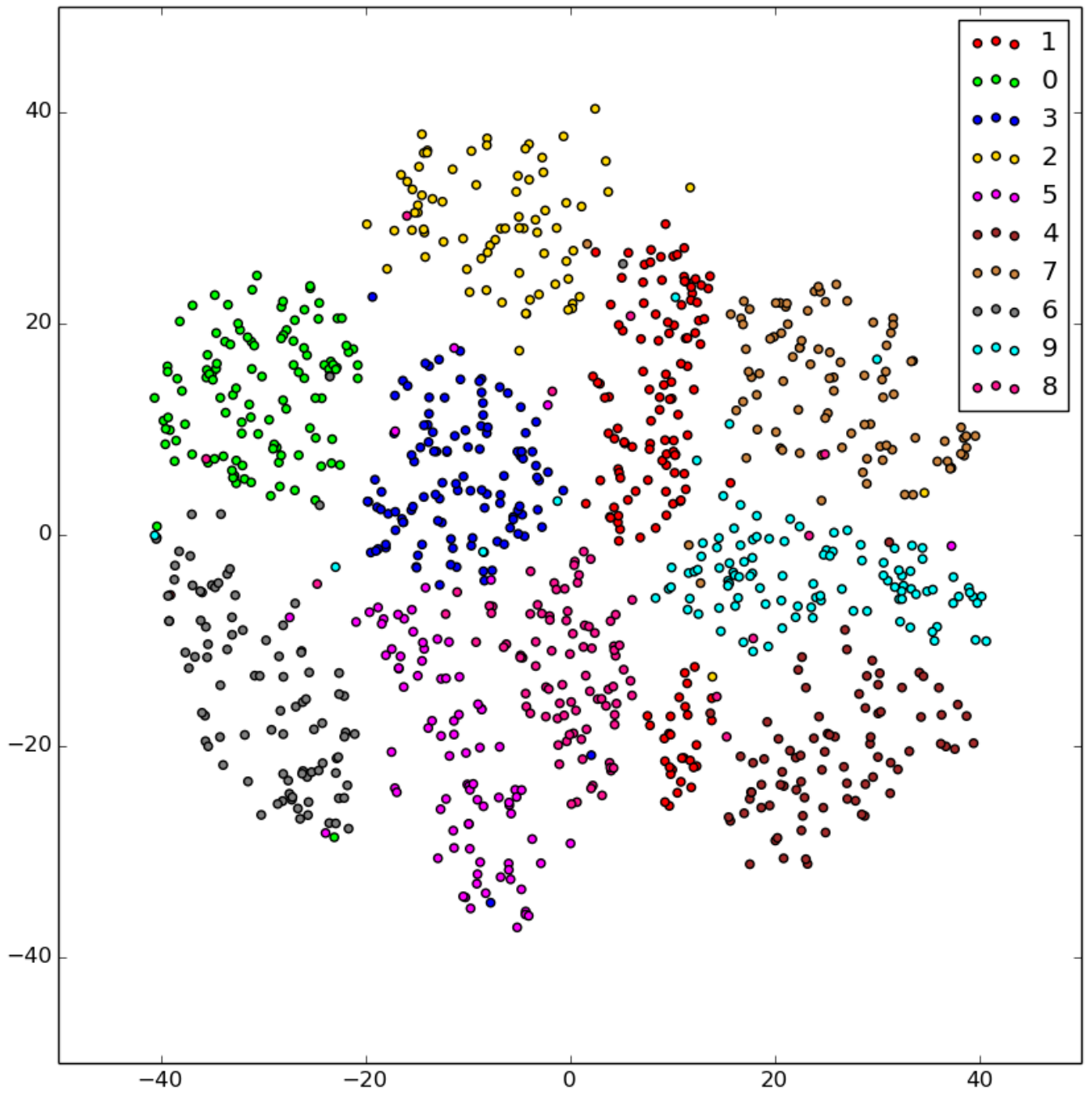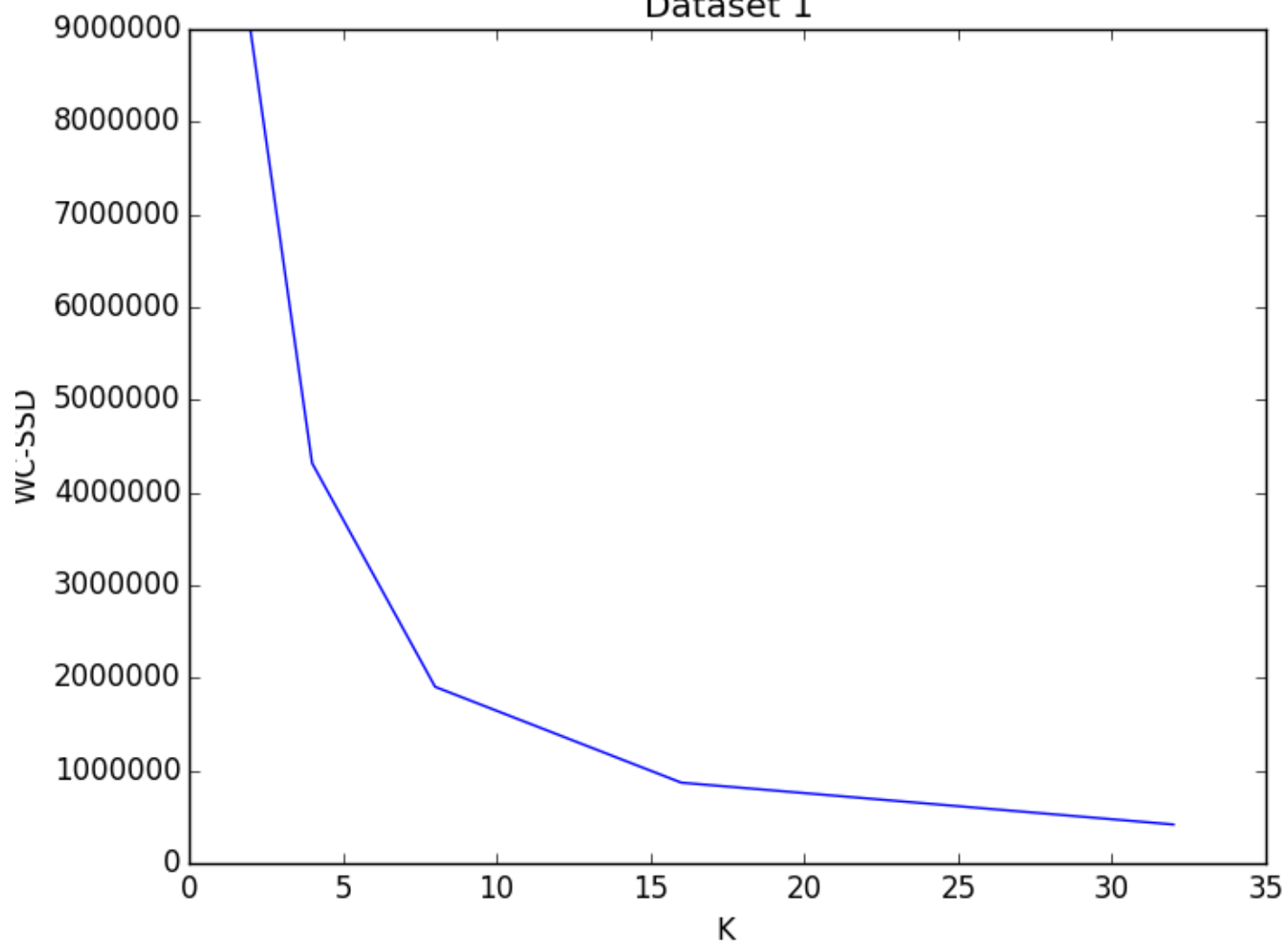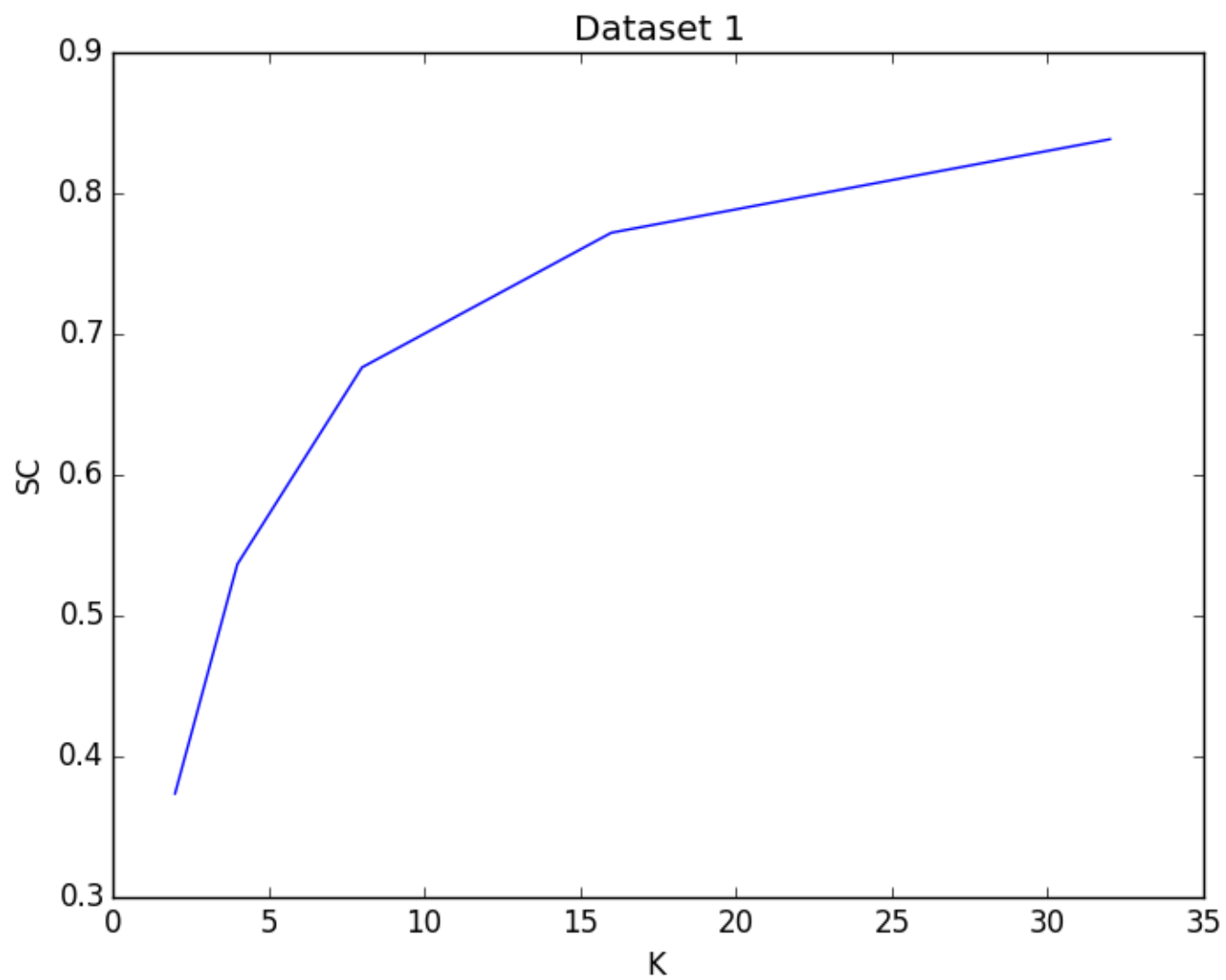f k, the reduction in WC-SSD is minimal for an increased computational complexity of running k-means. We also look at the knee-point in the plots of SC versus k (based on the definition of SC used).

Hence, from the above plots, we can see the appropriate value of K will be 16 for Dataset 1, 8 for Dataset 2, and 2 for Dataset 3.

We can see that the selected value of K is close to the number of digits (classes) in the each dataset. Also, we can observe that as the number of digits (classes) decrease, the WC-SSD values also decrease.

(c) The plots are

Dataset 2

Dataset 2

Dataset 3

From the above plots, we can conclude that k-means is sensitive to the initial starting conditions, since the standard deviation (or variance) is high in both the plots: WC-SSD versus K, and SC versus K.

(d) Solution

    i. Dataset 1

      $K = 16$

      $NMI = 0.368292602769$

ii. Dataset 2

K = 4

NMI = 0.45465341281

iii. Dataset 3

K = 2

NMI = 0.490710990204



We can see that the NMI value increases as the number of digits (classes) in the dataset decreases. The NMI value is highest for dataset 3 (2 classes), followed by dataset 2 ( 4 classes) followed by dataset 1 (10 classes).
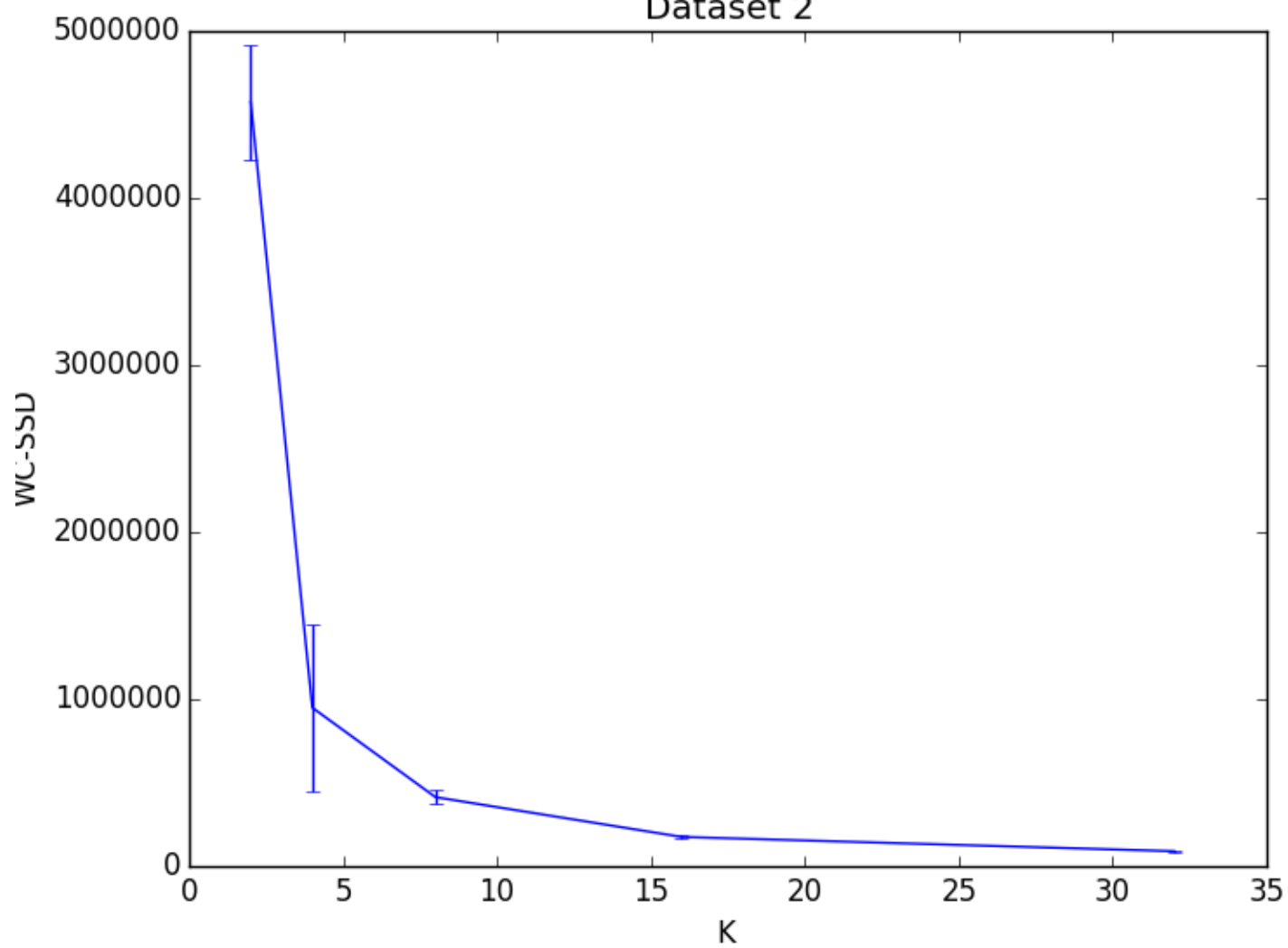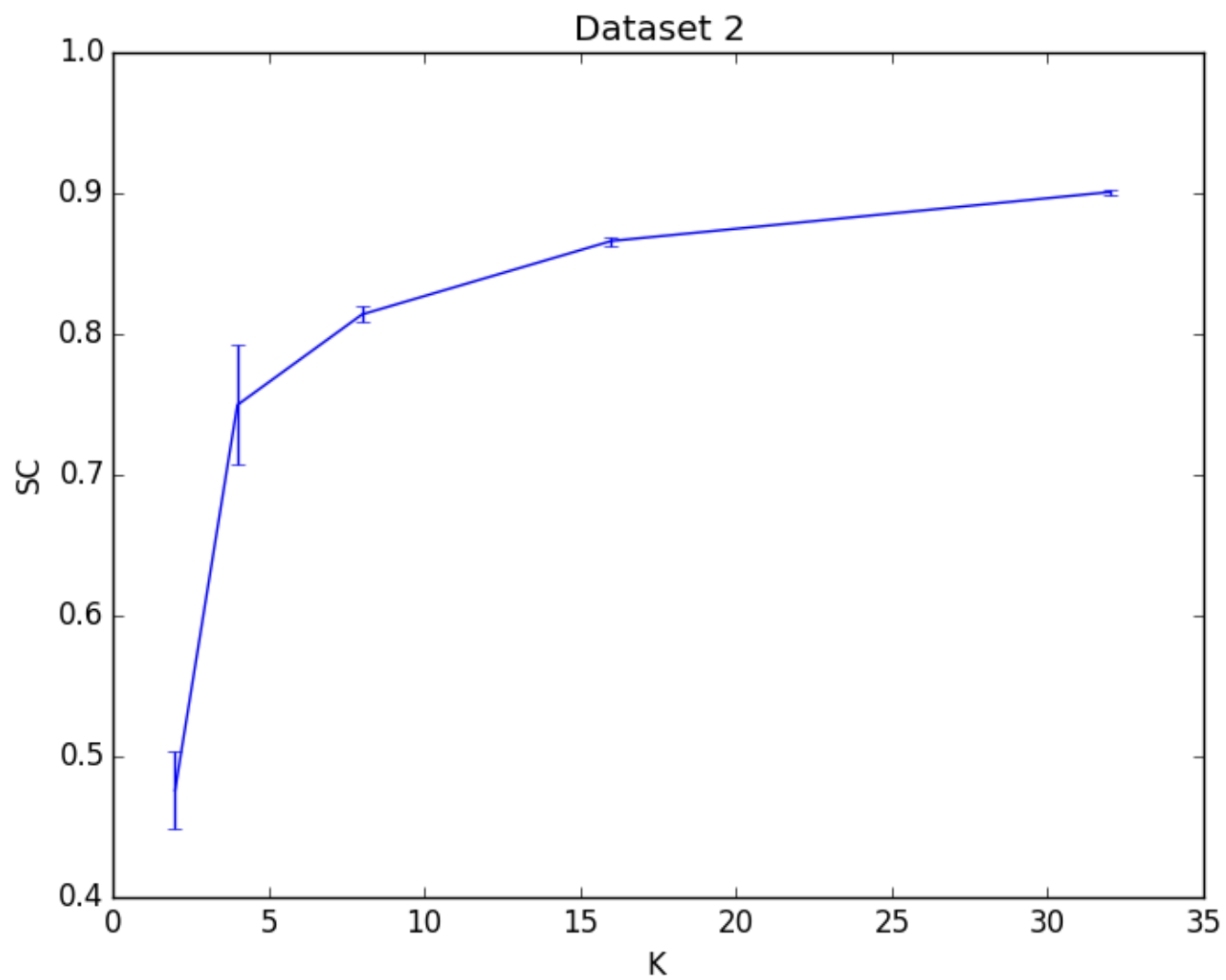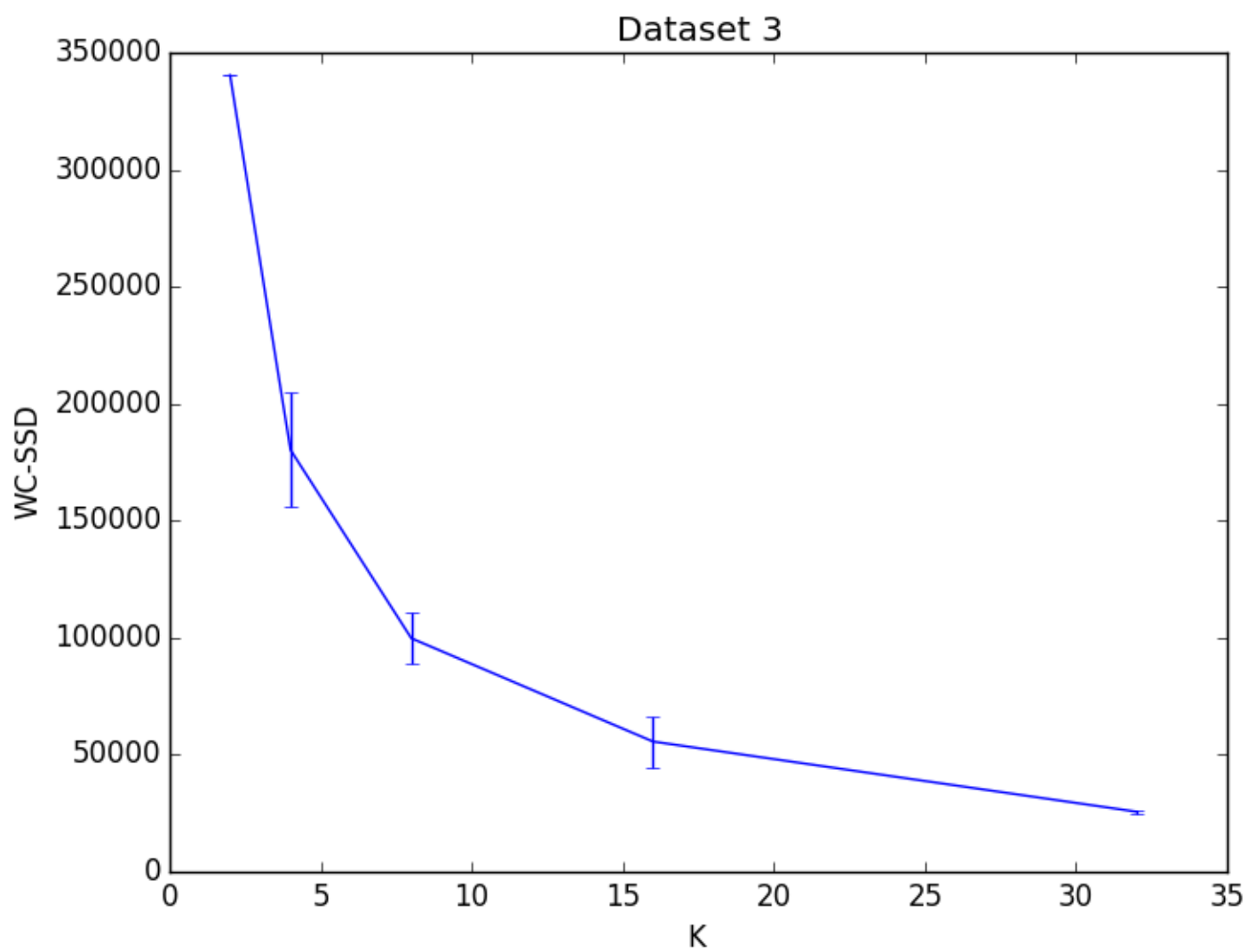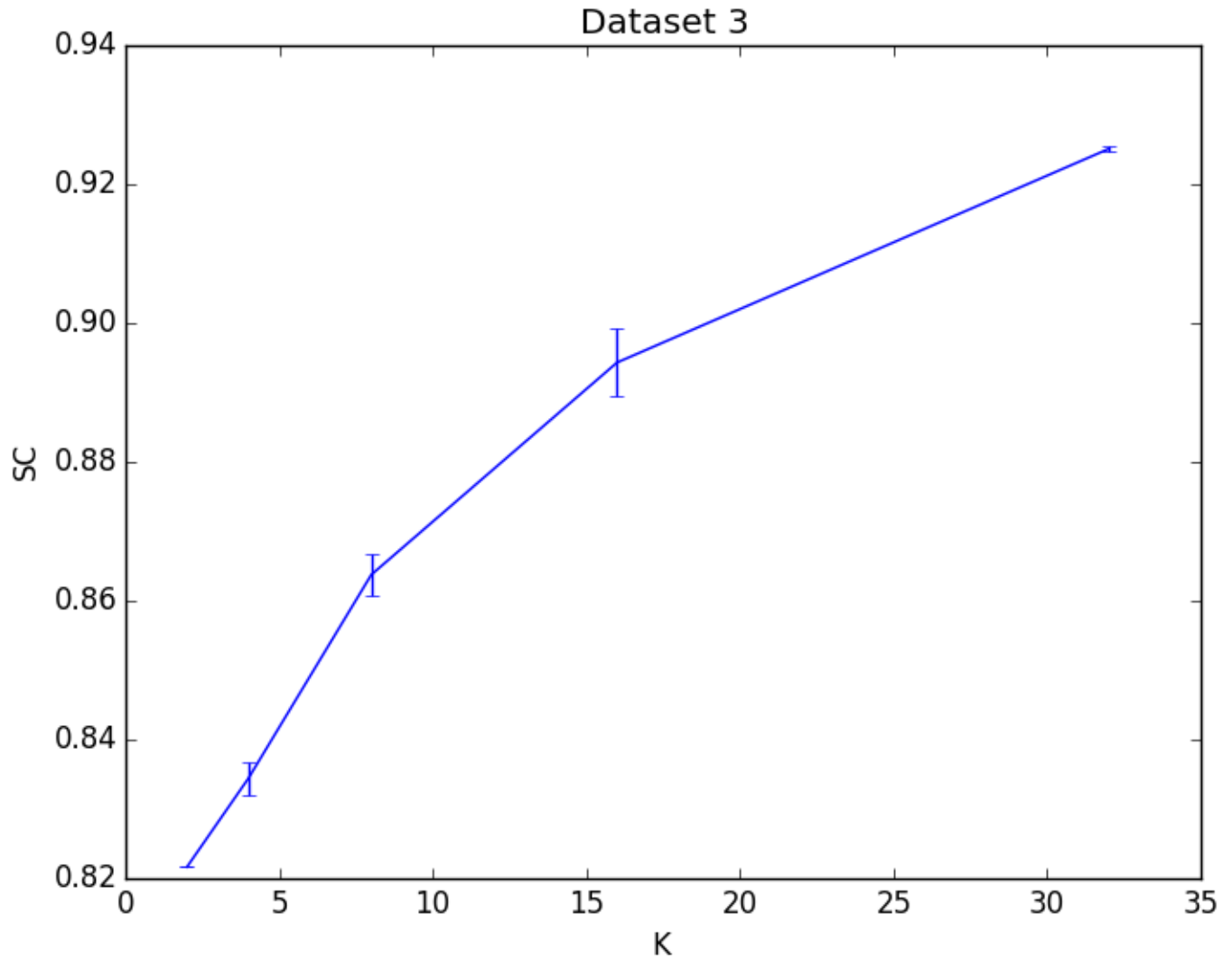
Also, in terms of visualization, the clusters become more separable when the number of classes in the dataset decreases. The clusters in dataset 3 (2 classes) are more clearly separated than those in dataset 2 (4 classes) which are more clearly separated than those in dataset 1 (10 classes).

3. **C**

(a) Dendrogram using Single Linkage as the distance measure.

Hierarchical Clustering Dendrogram (Single Linkage)

(b)  i. Dendrogram using Complete Linkage as the distance measure.



Hierarchical Clustering Dendrogram (Complete Linkage)

ii. Dendrogram using Average Linkage as the distance measure.

Hierarchical Clustering Dendrogram (Average Linkage)

(c)    i. Single Linkage

Hierarchical Clustering (Single Linkage)

ii. Complete Linkage

Hierarchical Clustering (Complete Linkage)

iii. Average Linkage

Hierarchical Clustering (Average Linkage)

Hierarchical Clustering (Average Linkage)

(d) To decide the value of k for each dataset, we can look at the elbow point in the plots of WC-SSD versus k, because beyond this value of k, the reduction in WC-SSD is minimal for an increased computational c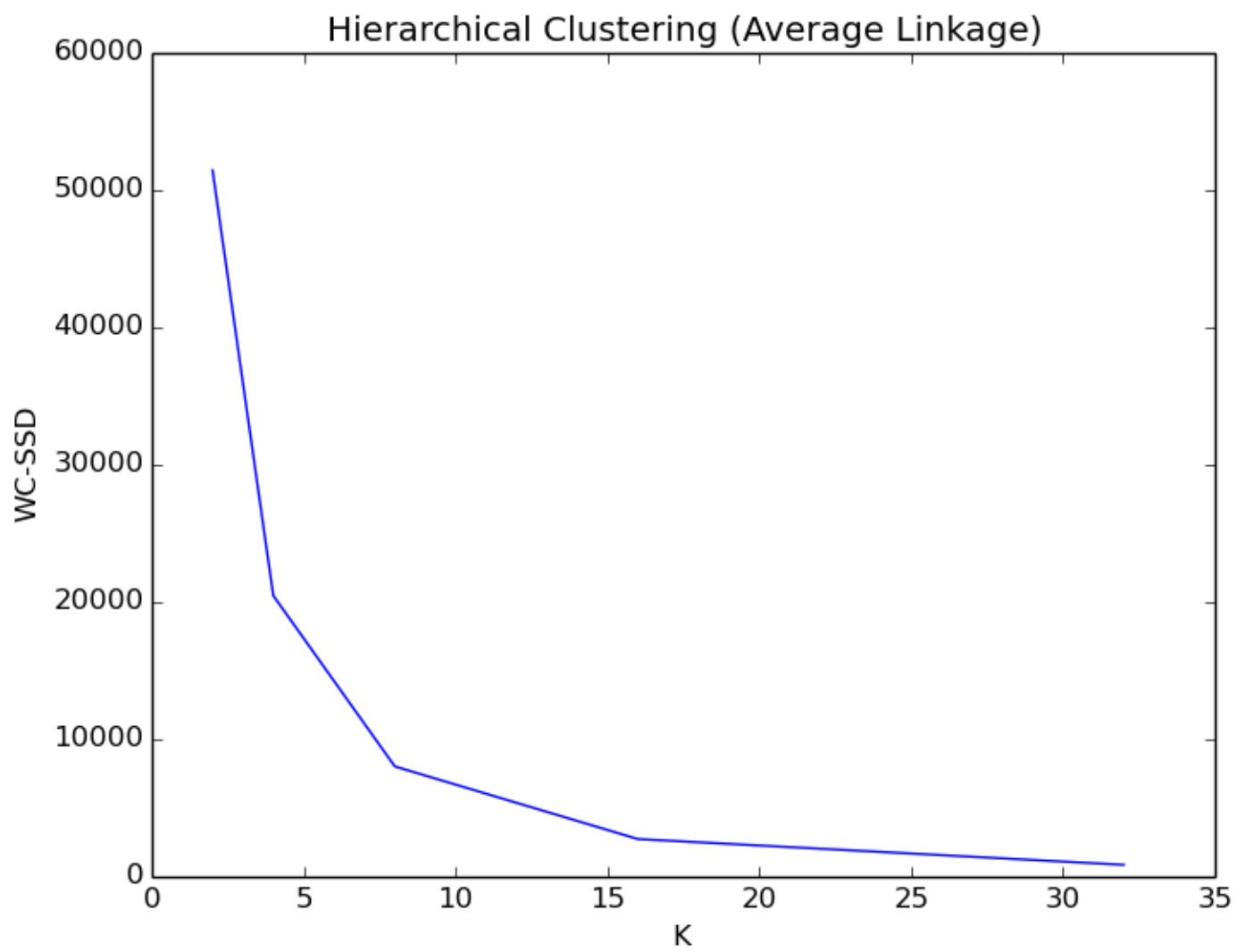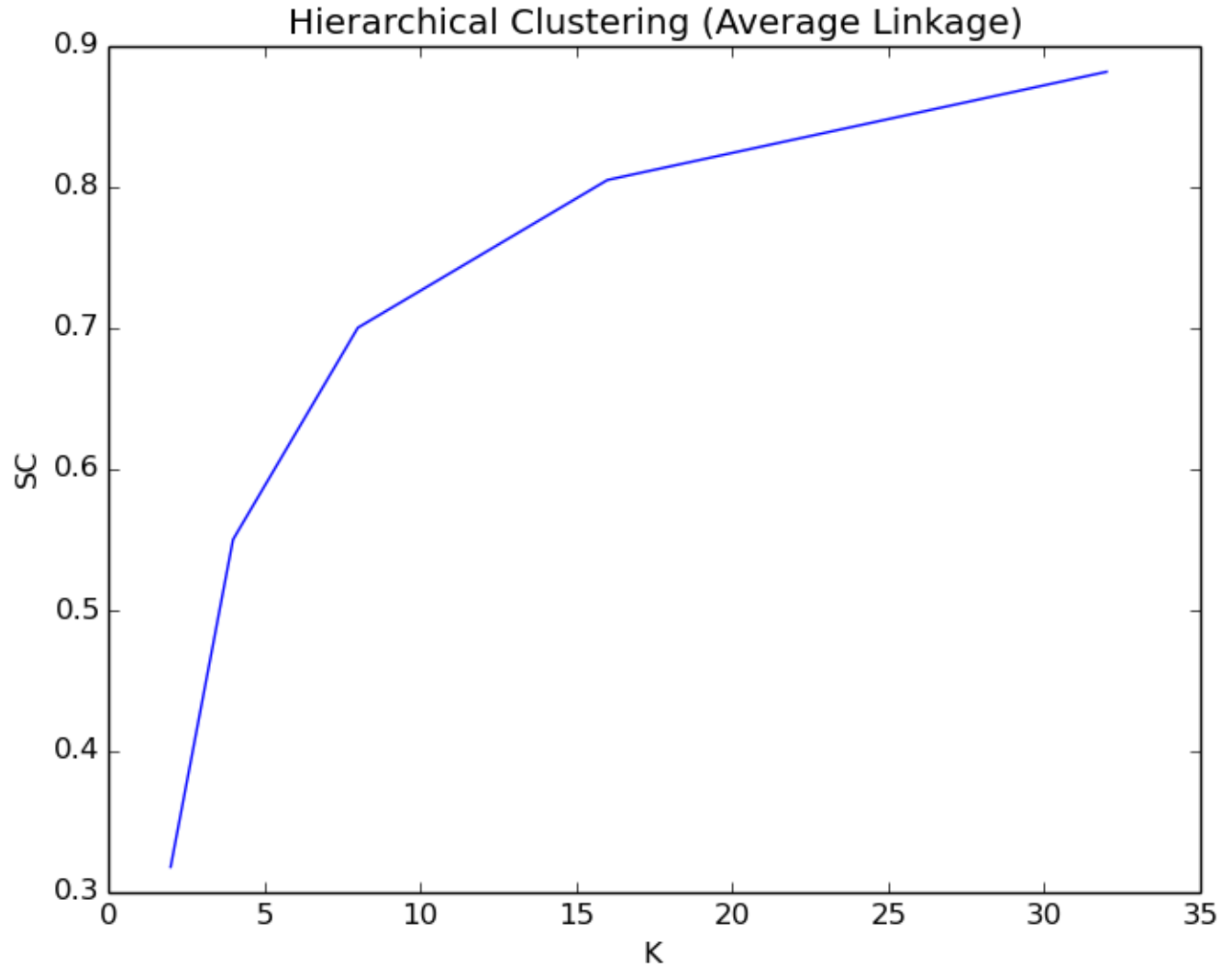omplexity of running k-means. We also look at the knee-point in the plots of SC versus k (based on the definition of SC used).

Hence, from the above plots, we can see the appropriate value of K will be 16 for all types of linkages : single, complete and average.

The selected value of K, which is 16, is same as that for the K chosen for kmeans which was also 16, for dataset 1 (complete dataset with 10 digits).

(e) The value of K chosen is 16 for all types of linkages: single, average and complete. The NMI values are:

Single Linkage: 0.3782359727

Complete Linkage: 0.4088758492

Average Linkage: 0.4373220577

The NMI score is highest for average linkage, followed by complete linkage, followed by single linkage. But, the differences in scores are not significant, and the dataset is too small in size (100 examples), and there is a lot of variation in the results across different runs (since the points are chosen randomly). Hence, these scores are rough estimates.

The NMI scores for hierarchical clustering using all (single, complete and hierarchical) linkages are higher than the NMI scores for k-means in part B.