

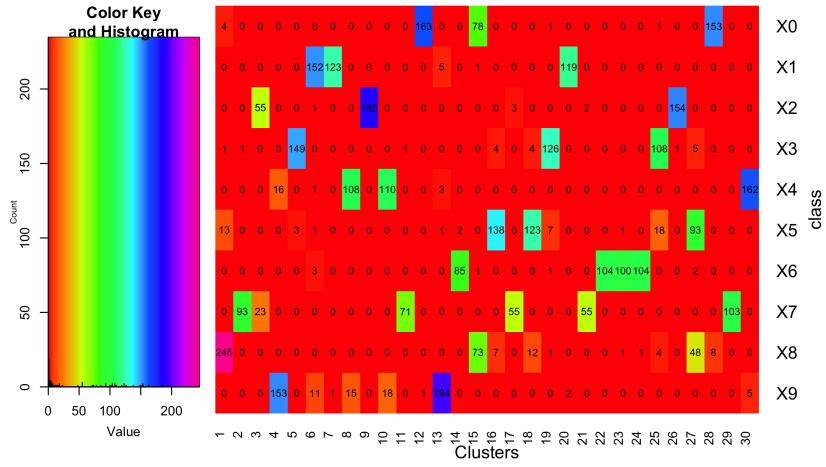
1 K - means clustering

Clustering is the task of grouping a set of data points in such a way that data points in the same group (called a cluster) are more similar (in some sense or another) to each other than to those in other groups (clusters). The clustering is done by minimizing the sum of square distance between each data point and their corresponding cluster centroid.

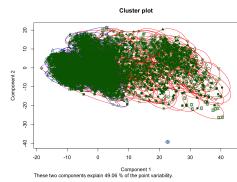
The number of clusters needed is determined by the user. It usually depends on how the data points is distributed. Choosing a large number of cluster will result in a reduced squared distance, and ideally zero as each data point gets its own cluster. The optimal choice of k will consist of having a balance between having the data compressed as possible, and still retain a decent amount of accuracy. More formally is the goal to partition the data into k clusters in order to minimize the total within cluster sum of squares. For which the elbow method might become handy.



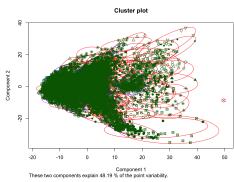
Class distribution within each cluster



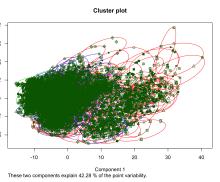
(a) Class distribution within each cluster



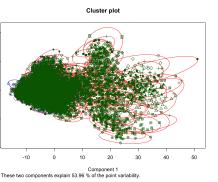
(b) Clustering data consisting of digit 0



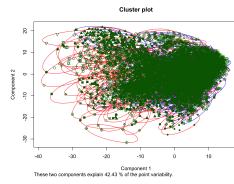
(c) Clustering data con-



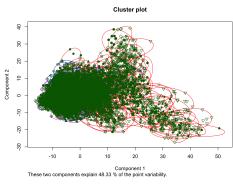
(d) Clustering data con-



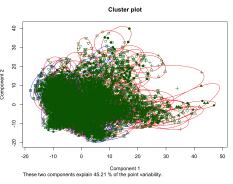
(e) Clustering data con-



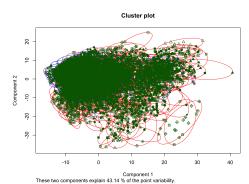
(f) Clustering data con-



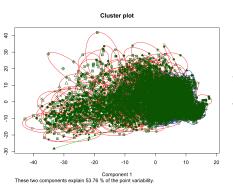
(g) Clustering data con-



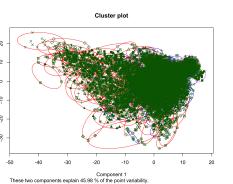
(h) Clustering data con-



(i) Clustering data con-



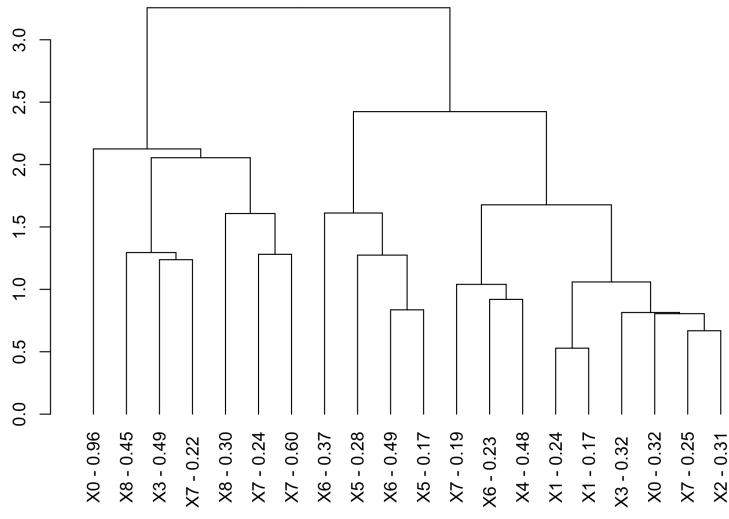
(j) Clustering data con-



(k) Clustering data con-

Figur 1: This and that

some text
some text



(a) Dendrogram



Insert some text