**CS204P : Data Structures & Algorithms Lab**

For a given set of data points in $\mathbf{R}^2$ space find K-Means Clustering of the points using a Red-Black Tree.

**Input** : $[set\,of\,data\,points : \{(x_i, y_i)\,|\,i = 0, 1, \ldots, n\},\ K : number\,of\,clusters]$

**Output** : $[set\,of\,data\,points\,with\,their\,associated\,cluster : \{(x_i, y_i, k)\,|\,i = 0, 1, \ldots, n\,;\ k \in K\}]$
**K-means Clustering Algorithm**

---

**Algorithm 1** K-Means Clustering (Lloyd's Algorithm)　　　*Note: written for clarity, not efficiency.*

1: **Input:** Data vectors $\{\boldsymbol{x}_n\}_{n=1}^N$, number of clusters $K$
2: **for** $n \leftarrow 1 \ldots N$ **do**　　　　　　　　　　　　　▷ Initialize all of the responsibilities.
3:　　$\boldsymbol{r}_n \leftarrow [0, 0, \cdots, 0]$　　　　　　　　　　　▷ Zero out the responsibilities.
4:　　$k' \leftarrow \text{RandomInteger}(1, K)$　　　　　▷ Make one of them randomly one to initialize.
5:　　$r_{nk'} = 1$
6: **end for**
7: **repeat**
8:　　**for** $k \leftarrow 1 \ldots K$ **do**　　　　　　　　　　　　▷ Loop over the clusters.
9:　　　　$N_k \leftarrow \sum_{n=1}^N r_{nk}$　　　　　　　▷ Compute the number assigned to cluster $k$.
10:　　　　$\boldsymbol{\mu}_k \leftarrow \frac{1}{N_k} \sum_{n=1}^N r_{nk}\boldsymbol{x}_n$　　　　　　▷ Compute the mean of the $k$th cluster.
11:　　**end for**
12:　　**for** $n \leftarrow 1 \ldots N$ **do**　　　　　　　　　　　　▷ Loop over the data.
13:　　　　$\boldsymbol{r}_n \leftarrow [0, 0, \cdots, 0]$　　　　　　　　▷ Zero out the responsibilities.
14:　　　　$k' \leftarrow \arg\min_k ||\boldsymbol{x}_n - \boldsymbol{\mu}_k||^2$　　　　　　▷ Find the closest mean.
15:　　　　$r_{nk'} = 1$
16:　　**end for**
17: **until** none of the $\boldsymbol{r}_n$ change
18: **Return** assignments $\{\boldsymbol{r}_n\}_{n=1}^N$ for each datum, and cluster means $\{\boldsymbol{\mu}_k\}_{k=1}^K$.

---

Figure 1: K-Means Clustering



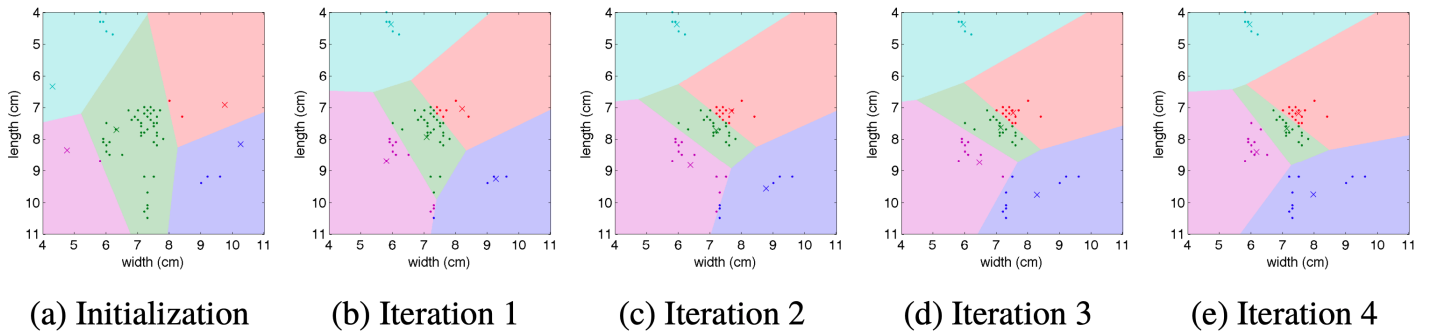(a) Initialization　　(b) Iteration 1　　(c) Iteration 2　　(d) Iteration 3　　(e) Iteration 4

Figure 2: Clustering Convergence