
CSE 151A - Discussion 08

K-Means Clustering

- Goal : find k cluster centers $\{\vec{\mu}^{(1)} \dots \vec{\mu}^{(k)}\}$ that minimizes the K-Means cost
- Cost ($\{\vec{\mu}^{(1)} \dots \vec{\mu}^{(k)}\}$) = $\frac{1}{n} \sum_{i=1}^n \min_{j \in \{1 \dots k\}} \|\vec{x}^{(i)} - \vec{\mu}^{(j)}\|^2$
 - This is the average squared distance from each data point to its nearest cluster center

Lloyd's Algorithm for K-Means

- Initialize $\{\vec{\mu}^{(1)} \dots \vec{\mu}^{(k)}\}$ in some way
- Repeat until no change in cost :
 - Assign each point $\vec{x}^{(i)}$ to closest center
 - Update $\vec{\mu}^{(i)}$ to be mean of points assigned to it
- Key Facts
 - Converges to local optimum of K-Means cost
 - Cost monotonically decreases as the algorithm progresses
 - Number of iterations unknown
 - Quality of solution depends heavily on initialization

K-Means++ Initialization

- Pick $\vec{\mu}^{(1)}$ uniformly at random from data
- Let $C = \{\vec{\mu}^{(1)}\}$ be the centers chosen so far
- Repeat $k - 1$ times :
 - Pick \vec{x} at random, with probability $P(\vec{x}) \propto \min_{\vec{\mu} \in C} \|\vec{x} - \vec{\mu}\|^2$
 - Add \vec{x} to C

Gaussian Mixture Models + EM Algorithm

- Mixture of k Gaussians : $\mathbb{P}(\vec{x}) = \sum_{j=1}^k \pi_j P_j(\vec{x}^{(i)})$
- Single Gaussian : $P_j = \mathcal{N}(\vec{\mu}^{(j)}, C_j)$
 - Mean : $\vec{\mu}^{(j)} = \frac{1}{\sum_{i=1}^n w_{ij}} \sum_{i=1}^n w_{ij} \vec{x}^{(i)}$
 - Covariance matrix : $C_j = \frac{1}{\sum_{i=1}^n w_{ij}} \sum_{i=1}^n w_{ij} (\vec{x}^{(i)} - \vec{\mu}^{(j)}) (\vec{x}^{(i)} - \vec{\mu}^{(j)})^T$
- Mixing weight : $\pi_j = \frac{1}{n} \sum_{i=1}^n w_{ij}$

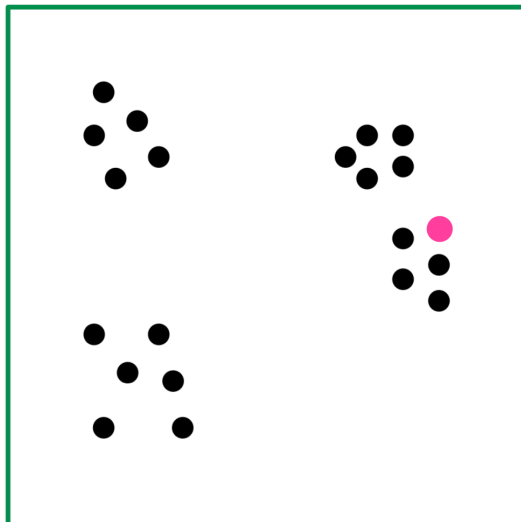
- Responsibility of cluster j for point i : $w_{ij} = \frac{\pi_j P_j(\vec{x}^{(i)})}{\sum_l \pi_l P_l(\vec{x}^{(i)})}$
- Algorithm:
 - Initialize $\pi_1 \cdots \pi_k, \vec{\mu}^{(1)} \cdots \vec{\mu}^{(k)}, C_1 \cdots C_k$
 - Make soft assignment (update responsibilities w_{ij})
 - Update mixing weights (π_j), means ($\vec{\mu}^{(j)}$), covariances (C_j)

Hierarchical Clustering Basics

- Single Linkage : $\mathcal{L}(C, C') = \min_{x, x' \in C, C'} d(x, x') \leftarrow$ smallest distance between any pair of points
- Complete Linkage : $\mathcal{L}(C, C') = \max_{x, x' \in C, C'} d(x, x') \leftarrow$ largest distance between any pair of points
- Average Linkage : $\mathcal{L}(C, C') = \frac{1}{|C||C'|} \sum_{x, x' \in C, C'} d(x, x') \leftarrow$ average distance between all point pairs
- Density Cluster Tree : For a probability density function f , assign clusters $C_f(\lambda)$ to connected components of $\{f \geq \lambda\}$ for any $\lambda > 0$

Problem 1.

Given the data points below, assign $k = 4$ cluster centers using the K-Means++ initialization algorithm, with the initial choice of $\vec{\mu}^{(1)}$ shown in pink. (Note that there are many possible correct solutions.) Once you have chosen the cluster centers, draw the boundary lines to define the $k = 4$ convex regions.



Solution:

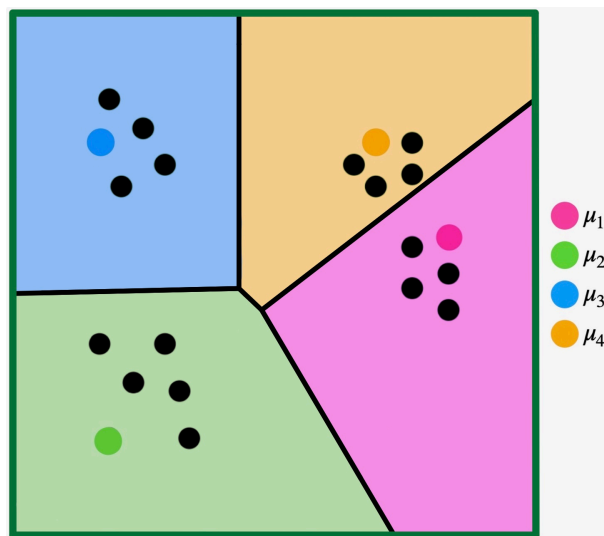
The K-Means++ algorithm chooses subsequent points $\vec{x}^{(i)}$ to act as cluster centers $\vec{\mu}^{(j)}$ with probability proportional to their distance from current cluster centers.

We are given $\vec{\mu}^{(1)}$, represented here as a pink dot.

We can reasonably select the green dot as $\vec{\mu}^{(2)}$, as it is far away from $\vec{\mu}^{(1)}$.

The same logic follows for the selection of $\vec{\mu}^{(3)}$ (it is far from $\vec{\mu}^{(1)}$ and $\vec{\mu}^{(2)}$), and finally for $\vec{\mu}^{(4)}$.

To define the cluster regions, simply draw perpendicular bisectors between each of the $\vec{\mu}^{(j)}$ and shade the area with the color corresponding to the center. This is often called a "voronoi" pattern.



Problem 2.

For a cluster S consisting of n points, and any arbitrary cluster center $\vec{\mu}$, we define the K-Means cost as :

$$\text{Cost}(\{\vec{\mu}\}) = \frac{1}{n} \sum_{i=1}^n \|\vec{x}^{(i)} - \vec{\mu}\|^2$$

We can also define the following lemma over any cluster S and any arbitrary cluster center $\vec{\mu}$:

$$\text{Cost}(\{\vec{\mu}\}) = \text{Cost}(\{\bar{x}\}) + \|\vec{\mu} - \bar{x}\|^2, \text{ where } \bar{x} = \text{mean}(S)$$

Using the above lemma, prove that for any data point $\vec{x}^{(i)}$ chosen randomly from S , it is true that :

$$\mathbb{E}_{\vec{x}^{(i)}}[\text{Cost}(\vec{x}^{(i)})] = 2 \cdot \text{Cost}(\bar{x})$$

This is to say that the expected value of the cost of clustering set S with any randomly chosen center $\vec{x}^{(i)}$ is twice the cost of a clustering with center $\bar{x} = \text{mean}(S)$.

Solution: We can write the expected value of $\text{Cost}(\vec{x}^{(i)})$ as the average over all possible choices of $\vec{x}^{(i)} \in S$ as follows:

$$\mathbb{E}_{\vec{x}^{(i)}}[\text{Cost}(\{\vec{x}^{(i)}\})] = \frac{1}{n} \sum_{i=1}^n [\text{Cost}(\{\vec{x}^{(i)}\})]$$

We then substitute for $\text{Cost}(\vec{x}^{(i)})$ using the equation provided in the lemma above.

$$= \frac{1}{n} \sum_{i=1}^n [\text{Cost}(\{\bar{x}\}) + \|\vec{x}^{(i)} - \bar{x}\|^2]$$

Distributing the summation through to both terms yields:

$$= \frac{1}{n} \sum_{i=1}^n \text{Cost}(\{\bar{x}\}) + \frac{1}{n} \sum_{i=1}^n \|\vec{x}^{(i)} - \bar{x}\|^2$$

We can then pull $\text{Cost}(\{\bar{x}\})$ out of the sum as it is independent of $\vec{x}^{(i)}$.

$$= \frac{1}{n} \text{Cost}(\{\bar{x}\}) \sum_{i=1}^n (1) + \frac{1}{n} \sum_{i=1}^n \|\vec{x}^{(i)} - \bar{x}\|^2$$

Simplifying $\sum_{i=1}^n (1) = n$ gives us the following :

$$= \frac{1}{n} \text{Cost}(\{\bar{x}\})n + \frac{1}{n} \sum_{i=1}^n \|\vec{x}^{(i)} - \bar{x}\|^2$$

We can cancel terms on the left, leaving us with :

$$= \text{Cost}(\{\bar{x}\}) + \frac{1}{n} \sum_{i=1}^n \|\vec{x}^{(i)} - \bar{x}\|^2$$

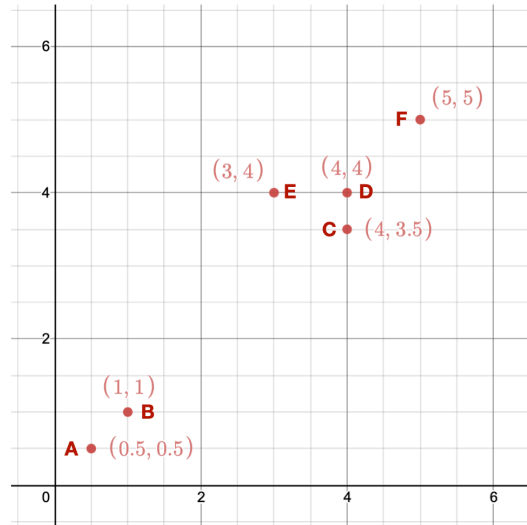
The right hand term is exactly the equation for $\text{Cost}(\{\bar{x}\})$, and making that substitution yields :

$$= \text{Cost}(\{\bar{x}\}) + \text{Cost}(\{\bar{x}\}) \\ = 2 \cdot \text{Cost}(\{\bar{x}\})$$

Thus we have shown that $\mathbb{E}_{\vec{x}^{(i)}}[\text{Cost}(\{\vec{x}^{(i)}\})] = 2 \cdot \text{Cost}(\{\bar{x}\})$.

Problem 3.

Run the single-linkage clustering algorithm on the data points below, drawing the cluster links as you go. Once you are left with a single cluster, draw the corresponding dendrogram.



Solution:

Recall the distance between two clusters C and C' for single linkage : $\mathcal{L}(C, C') = \min_{x, x' \in C, C'} d(x, x')$.

Using the single linkage function above, we will group clusters in the following order:

$\{C\}$ and $\{D\}$ with $\mathcal{L}(\{C\}, \{D\}) = 0.5$

$\{A\}$ and $\{B\}$ with $\mathcal{L}(\{A\}, \{B\}) = \frac{\sqrt{2}}{2}$

$\{C, D\}$ and $\{E\}$ with $\mathcal{L}(\{C, D\}, \{E\}) = 1$

$\{C, D, E\}$ and $\{F\}$ with $\mathcal{L}(\{C, D, E\}, \{F\}) = \sqrt{2}$

$\{C, D, E, F\}$ and $\{A, B\}$ with $\mathcal{L}(\{C, D, E, F\}, \{A, B\}) = \sqrt{13}$

This results in the dendrogram shown below.

