

Machine Learning

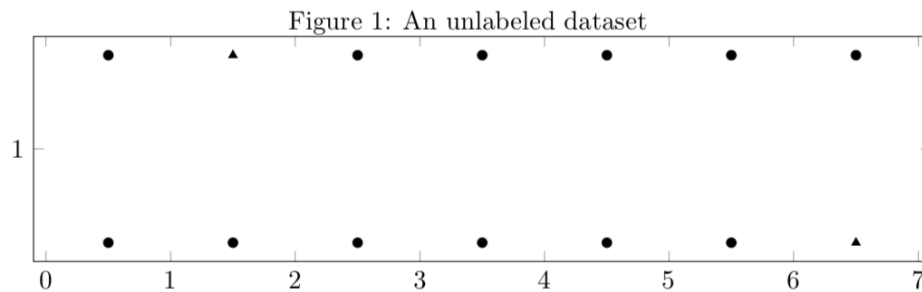
CSCI 567 Spring 2021

Discussion: K-means, Naive Bayes

Date: Mar 26th, 2021

1. Consider the following dataset. All points are unlabeled and part of the same set. The triangles are used to distinguish two points later.

Suppose we run the K-means algorithm on this dataset with $K = 2$ and the two points indicated by triangles as the initial centroids. When the algorithm converges, there will be two clearly separated clusters. Directly on Figure 1, draw a straight line that separates these two clusters, as well as the centroids of these two clusters.



Next, find two other different sets of initialize centroids that will converge to the exact same result if we apply K-means. Please follow the instructions below

- the initialize centroids have to be points of the dataset;
- these two sets of points can overlap with each other, but of course cannot be the same;
- similarly these sets can overlap with the initialization of Figure 1 but cannot be the same;
- do not pick those that lead to ambiguous results due to different ways of breaking ties.

2. Suppose we have the following training data. Each data point has three features (Weather, Emotion, Homework), where $\text{Weather} \in \{\text{Sunny}, \text{Cloudy}\}$, $\text{Emotion} \in \{\text{Happy}, \text{Normal}, \text{Unhappy}\}$, $\text{Homework} \in \{\text{Much}, \text{Little}\}$. The label PlayBasketball indicates whether it is suitable to play basketball. You are asked to build a naive Bayes classifier. Recall the naive Bayes assumption is

$$P(\text{Weather}, \text{Emotion}, \text{Homework} | \text{PlayBasketball}) = P(\text{Weather} | \text{PlayBasketball}) \times P(\text{Emotion} | \text{PlayBasketball}) \times P(\text{Homework} | \text{PlayBasketball}).$$

(a) Write down the MLE for the following parameters (you only need to provide the final number):

| Weather | Emotion | Homework | PlayBasketball |
|---------|---------|----------|----------------|
| Sunny | Happy | Little | Yes |
| Sunny | Normal | Little | Yes |
| Cloudy | Happy | Much | Yes |
| Cloudy | Unhappy | Little | Yes |
| Sunny | Unhappy | Little | No |
| Cloudy | Normal | Much | No |

- $P(\text{PlayBasketball} = \text{No})$
- $P(\text{Weather} = \text{Cloudy} | \text{PlayBasketball} = \text{Yes})$
- $P(\text{Emotion} = \text{Unhappy} | \text{PlayBasketball} = \text{No})$
- $P(\text{Homework} = \text{Little} | \text{PlayBasketball} = \text{Yes})$

(b) Given a new data instance $\mathbf{x} = (\text{Weather} = \text{Cloudy}, \text{Emotion} = \text{Unhappy}, \text{Homework} = \text{Little})$, compute $P(\text{PlayBasketball} = \text{No} | \mathbf{x})$ (show your work).