

## Unsupervised Learning

LATEST SUBMISSION GRADE

100%

1. For which of the following tasks might K-means clustering be a suitable algorithm? Select all that apply.

1 / 1 point

- ☐ Given many emails, you want to determine if they are Spam or Non-Spam emails.
- ☐ Given historical weather records, predict if tomorrow's weather will be sunny or rainy.
- ☒ Given a set of news articles from many different news websites, find out what are the main topics covered.

✓ Correct

K-means can cluster the articles and then we can inspect them or use other methods to infer what topic each cluster represents

- ☒ From the user usage patterns on a website, figure out what different groups of users exist.

✓ Correct

We can cluster the users with K-means to find different, distinct groups.

2. Suppose we have three cluster centroids  $\mu_1 = \begin{bmatrix} 1 \\ 2 \end{bmatrix}$ ,  $\mu_2 = \begin{bmatrix} -3 \\ 0 \end{bmatrix}$  and  $\mu_3 = \begin{bmatrix} 4 \\ 2 \end{bmatrix}$ . Furthermore, we have a

1 / 1 point

- ☐ Using the elbow method to choose K.
- ☒ Move the cluster centroids, where the centroids  $\mu_k$  are updated.

✓ Correct

The cluster update is the second step of the K-means loop.

- ☐ Feature scaling, to ensure each feature is on a comparable scale to the others.

- ☒ The cluster assignment step, where the parameters  $c^{(i)}$  are updated.

✓ Correct

This is the correct first step of the K-means loop.

✓ Correct

This function is the distortion function. Since a lower value for the distortion function implies a better clustering, you should choose the clustering with the smallest value for the distortion function.

5. Which of the following statements are true? Select all that apply.

1 / 1 point

- ☒ If we are worried about K-means getting stuck in bad local optima, one way to ameliorate (reduce) this problem is if we try using multiple random initializations.

✓ Correct

Since each run of K-means is independent, multiple runs can find different optima, and some should avoid bad local optima.

- ☐ Since K-Means is an unsupervised learning algorithm, it cannot overfit the data, and thus it is always better to have as large a number of clusters as is computationally feasible.

- ☒ For some datasets, the "right" or "correct" value of K (the number of clusters) can be ambiguous, and hard even for a human expert looking carefully at the data to decide.

✓ Correct

In many datasets, different choices of K will give different clusterings which appear quite reasonable. With no labels on the data, we cannot say one is better than the other.

- ☐ The standard way of initializing K-means is setting  $\mu_1 = \dots = \mu_k$  to be equal to a vector of zeros.