

Course > Week 8... > Proble... > Proble...

## **Problem Set 8**

Problems 1-5 correspond to "Clustering with the k-means algorithm I"

## Problem 1

4/4 points (graded)

Consider the following data set consisting of five points in  $\mathbb{R}^1$ :

$$-10, -8, 0, 8, 10.$$

We would like to cluster these points into k=3 groups. Determine the optimal k-means solution.

a) What is the location of the leftmost center?



b) What is the location of the middle center?



c) What is the location of the rightmost center?



d) What is the k-means cost of this optimal solution?



**1** Answers are displayed within the problem

#### Problem 2

4/4 points (graded)

Recall that the previous problem dealt with a data set of five points in  $\mathbb{R}^1$ ,

$$-10, -8, 0, 8, 10,$$

and for k=3 we figured out the optimal k-means solution for this data. Unfortunately, Lloyd's algorithm does not always return the optimal solution; in fact, its output depends, in part, on the initialization.

In this problem, suppose we run Lloyd's algorithm, starting with the initialization  $\mu_1=-10, \mu_2=-8, \mu_3=0$ . Determine the final set of cluster centers obtained by the algorithm.

a) What is the location of the leftmost center?



b) What is the location of the middle center?

c) What is the location of the rightmost center?	
6 <b>✓ Answer:</b> 6	
6	
d) What is the $k$ -means cost of this particular solution?	
<b>✓ Answer:</b> 56	
56	
Explanation In the first iteration of $k$ -means, only $\mu_3$ changes. It moves to the mean $6$ . Thereafter, nothing moves.	of $0,8,10$ , which is
Submit	
Answers are displayed within the problem	
Problem 3	
1/1 point (graded) Suppose we have $n$ data points in $\mathbb{R}^d$ , and want to find the $k$ -means clawhat is the running time of a <i>single iteration</i> of Lloyd's algorithm?	ustering of this data.
$\bigcirc nd$	
$\bigcirc kd$	
$\bigcirc nk$	

$\bigcirc nkd$
<b>✓</b>
Explanation Step 1: For each data point, we need to compute its distance to each of the $k$ centers. This is a total of $nk$ distance computations. Each distance computation takes $O\left(d\right)$ time. So, this step takes $O\left(nkd\right)$ time. Step 2: Assign each point to its closest center. For each point, look at the (just-computed) distances to the $k$ centers and pick the smallest one. This takes time $O\left(k\right)$ for each data point, and thus $O\left(nk\right)$ overall. Step 3: Update the centers. Each center is reset to the mean of the points assigned to it. Computing the mean of $m$ points in $d$ -dimensions takes time $O\left(md\right)$ , so the total time for this step is $O\left(nd\right)$ .
Answers are displayed within the problem
Problem 4
1/1 point (graded) Which of the following properties are true of the $k\text{-means}++$ algorithm, when run on its own? Select all that apply.
$oxedsymbol{oxed}$ If called multiple times, it always returns the same $k$ centers.
✓ The centers it returns are always actual data points.
$oxedsymbol{oxed}$ It is guaranteed to return the optimal $k$ -means solution.
$oxedsymbol{oxed}$ It selects a set of $k$ centers, then spends some time iteratively improving them.

Submit

Answers are displayed within the problem
Problem 5
1/1 point (graded) Suppose we have $n$ data points in $\mathbb{R}^d$ . What is the running time of the $k$ -means $++$ algorithm?
$igcirc n^2kd$
$igcirc n^2 k$
$\boxed{ \bigcirc  nkd }$
$igcirc nk^2d$
Answer Correct: At any given time, the algorithm maintains a table $T$ of size $n$ with the squared distance from each data point to the closest of the centers chosen so far. On each iteration: (1) A new center is chosen at random, picking point $x_i$ with probability proportional to $T\left[i\right]$ ; this take time $O\left(n\right)$ . (2) The table $T$ is updated: this requires computing the distance from each data point to the newly chosen center, and since there are $n$ data points, and each distance computation takes time $O\left(d\right)$ , the total time for this step is $O\left(nd\right)$ . Since there are $k$ iterations, the overall running time is $O\left(nkd\right)$ .
<b>?</b> Hint (1 of 1): Think of it in this way: At any given time, the algorithm maintains a table $T$ of size $n$ with the squared distance from each data point to the closest of the centers chosen so far.
Submit

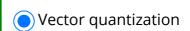
**1** Answers are displayed within the problem

Problems 6-8 correspond to "Clustering with the k-means algorithm II"

### Problem 6

1/1 point (graded)

In lecture, we described two general uses of clustering: for vector quantization and for finding natural groups in data. Which of these two uses is most directly captured by the k-means cost function?



Finding natural groups in data



Submit

**1** Answers are displayed within the problem

### Problem 7

1/1 point (graded)

In the streaming model of computation, we need to cluster data but are unable to hold the entire data set in memory. Suppose we only have room for M data points, while our full data set has size  $n\gg M$ . Here's one suggestion for a streaming algorithm for k-means clustering (assuming  $k\ll M$ ):

- ullet Pick M of the data points at random and store them in memory.
- ullet Run k-means clustering on these M points, to get centers  $\mu_1,\ldots,\mu_k$  .
- ullet For each point in the full training set: assign it to the cluster with the closest  $\mu_j$ .

Algorithms like this, based on random sampling, often work reasonably well. In which of the following situations might it fail badly, though? Select all that apply.

When there are important clusters in the full data set that contain very few points.
When there are duplicate copies of some of the data points.
When the optimal cluster centers are far away from each other.
<b>✓</b>
<b>Explanation</b> Methods that sample a data set uniformly at random are in danger of overlooking structure that is limited to small groups of points. In this case, the random subset might not include any representatives from clusters with few points. As a result, these clusters will not be correctly identified.
Submit
Answers are displayed within the problem
Problem 8
1/1 point (graded) The $sequential\ k$ -means $algorithm$ presented in lecture is able to handle an infinite stream of data, arriving one point at a time. It always maintains a set of $k$ centers that are based on the data seen so far. If the data is $d$ -dimensional, how long does this algorithm take to process the $n$ th data point that arrives?
$\bigcirc k$
$\boxed{ \bigcirc kd }$
$\bigcirc nd$
$\bigcirc nkd$
✓

When a new data point arrives, the algorithm computes its distance to the k current centers. Since each distance computation takes time  $O\left(d\right)$ , the overall time for the distance computations is  $O\left(kd\right)$ . The subsequent operation of updating the closest center takes time only  $O\left(d\right)$ .

		٠,
<u> </u>	ını	mit
ור	,,,,,	

**1** Answers are displayed within the problem

Problems 9-12 correspond to "Clustering with mixtures of Gaussians"

#### Problem 9

1/1 point (graded)

A *soft clustering* of a data set into k clusters is a clustering in which each data point isn't necessarily assigned to a single cluster, but rather is allowed to be fractionally assigned to all k clusters. We can represent a soft clustering of n data points by an  $n \times k$  matrix P, where

 $P_{ij}$  = fraction of the ith data point that is assigned to cluster j

and  $P_{i1} + \cdots + P_{ik} = 1$  for all  $i = 1, \ldots, n$ . A more traditional clustering, in which each point chooses a single cluster, is sometimes called a *hard clustering*. If matrix P corresponds to a hard clustering, which of the following conditions must hold? Select all that apply.

Each en	trv of I	) is ei	ther (	) or	1

 $\checkmark$  Each row of P has a single nonzero entry.

 $oxedsymbol{\square}$  Each column of P has a single nonzero entry.

 $oxedsymbol{oxed}$  No two rows of P are identical.



#### **Explanation**

In a hard clustering, each data point is assigned to exactly one cluster: thus, the corresponding row of P contains exactly 1 one and k-1 zeros.

Submit
Answers are displayed within the problem
Problem 10  1/1 point (graded) Which of the following is true of the EM algorithm for fitting a mixture of Gaussians to data? Select all that apply.
It finds the optimal maximum-likelihood solution.
✓ Different initializations could lead to different running times.
✓ Each iteration consists of a soft clustering of the data, followed by an update of the mixture parameters.
The algorithm eventually converges to a hard clustering of the data.
Submit
Answers are displayed within the problem
Problem 11 1/1 point (graded) Which of the following describe ways in which the mixture-of-Gaussians clustering method tries to improve upon $k$ -means?
✓ It can accommodate clusters of more general shapes and sizes.
It is faster.

✓ The final model doesn't just cluster the data, but also gives probabilities for cluster labels.
It is easier to implement for massive data sets.
Submit
Answers are displayed within the problem
Problem 12
0 points possible (ungraded) For you to think about: Although we talked specifically about mixtures of Gaussians, it is also common to look at mixtures of other kinds of distributions. And indeed, the EM algorithm can quite easily be adapted to these other cases as well. Can you flesh out how this might be done? (No need to enter any answer.)
<b>Explanation</b> We can think of a similar scheme that <i>alternates</i> between assigning points to clusters (soft partitioning) and then updating the parameters of the cluster distributions.
Submit
Answers are displayed within the problem
Problems 13-15 correspond to "Hierarchical clustering"

# Problem 13

1/1 point (graded)

Which of the following are reasons for which hierarchical clustering might be preferred to flat clustering? Select all that apply.

It does not require the number of clusters to be specified.
It captures the structure of the data at multiple scales.
It is computationally simpler.
It has an intuitive cost function for which an optimal solution can efficiently be obtained.
Submit
Answers are displayed within the problem
Problem 14
1/1 point (graded) Bottom-up linkage" methods for hierarchical clustering work by (i) initially putting each data point in its own cluster and then (ii) successively merging two existing clusters. How many merge steps are needed when there are $n$ data points? Your answer should be a function of $n$ .
[Your answer will not be correctly parsed if you include spaces in it.]
n-1 ✓ Answer: n-1
<b>Explanation</b> We start with $n$ clusters: each data point is in a cluster by itself. Each merge step reduces the overall number of clusters by 1. At the end of the process, there is just one big cluster containing all the data points. Thus the number of merge steps must be $n-1$ .

Submit

Answers are displayed within the problem
Problem 15
1/1 point (graded) Some hierarchical clustering algorithms can immediately be applied to data in any metric space. Which of the following algorithms have this property? Select all that apply.
✓ Single linkage
✓ Complete linkage
✓ Average linkage based on average pairwise distance between clusters
Ward's method
Explanation All of these methods can be used with any distance function, except for Ward's method, which explicitly uses Euclidean distance.  Submit
Answers are displayed within the problem
© All Rights Reso