

## Homework 4: Clustering and EM

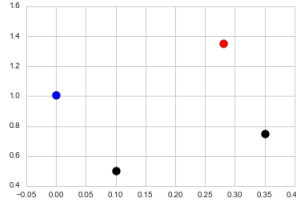
This homework assignment focuses on different unsupervised learning methods from a theoretical and practical standpoint. In Problem 1, you will explore Hierarchical Clustering and experiment with how the choice of distance metrics can alter the behavior of the algorithm. In Problem 2, you will derive from scratch the full expectation-maximization algorithm for fitting a Poisson mixture model. In Problem 3, you will implement PCA on a dataset of handwritten images and analyze the latent structure learned by this algorithm.

There is a mathematical component and a programming component to this homework. Please submit your PDF, tex, and Python files to Canvas, and push all of your work to your GitHub repository. If a question requires you to make any plots, please include those in the writeup.

**Problem 1** (Hierarchical Clustering, 7 pts)

At each step of hierarchical clustering, the two most similar clusters are merged together. This step is repeated until there is one single group. We saw in class that hierarchical clustering will return a different result based on the pointwise-distance and cluster-distance that is used. In this problem you will examine different choices of pointwise distance (specified through choice of norm) and cluster distance, and explore how these choices change how the HAC algorithm runs on a toy data set.

Consider the following four data points in  $\mathbb{R}^2$ , belonging to three clusters: the black cluster consisting of  $\mathbf{x}_1 = (0.1, 0.5)$  and  $\mathbf{x}_2 = (0.35, 0.75)$ , the red cluster consisting of  $\mathbf{x}_3 = (0.28, 1.35)$ , and the blue cluster consisting of  $\mathbf{x}_4 = (0, 1.01)$ .



Different pointwise distances  $d(\mathbf{x}, \mathbf{x}') = \|\mathbf{x} - \mathbf{x}'\|_p$  can be used. Recall the definition of the  $\ell_1$ ,  $\ell_2$ , and  $\ell_\infty$  norm:

$$\|\mathbf{x}\|_1 = \sum_{j=1}^m |x_j| \quad \|\mathbf{x}\|_2 = \sqrt{\sum_{j=1}^m x_j^2} \quad \|\mathbf{x}\|_\infty = \max_{j \in \{1, \dots, m\}} |x_j|$$

Also recall the definition of min-distance, max-distance, centroid-distance, and average-distance between two clusters (where  $\boldsymbol{\mu}_G$  is the center of a cluster  $G$ ):

$$\begin{aligned} d_{\min}(G, G') &= \min_{\mathbf{x} \in G, \mathbf{x}' \in G'} d(\mathbf{x}, \mathbf{x}') \\ d_{\max}(G, G') &= \max_{\mathbf{x} \in G, \mathbf{x}' \in G'} d(\mathbf{x}, \mathbf{x}') \\ d_{\text{centroid}}(G, G') &= d(\boldsymbol{\mu}_G, \boldsymbol{\mu}_{G'}) \\ d_{\text{avg}}(G, G') &= \frac{1}{|G||G'|} \sum_{\mathbf{x} \in G} \sum_{\mathbf{x}' \in G'} d(\mathbf{x}, \mathbf{x}') \end{aligned}$$

1. Draw the 2D unit sphere for each norm, defined as  $\mathcal{S} = \{\mathbf{x} \in \mathbb{R}^2 : \|\mathbf{x}\| = 1\}$ . Feel free to do it by hand, take a picture and include it in your pdf.
2. For each norm  $(\ell_1, \ell_2, \ell_\infty)$  and each clustering distance, specify which two clusters would be the first to merge.
3. Draw the complete dendrograms showing the order of agglomerations for the  $\ell_2$  norm and each of the clustering distances. We have provided some code to make this easier for you. You are not required to use it.

**Solution**

**Problem 2** (Expectation-Maximization for Poisson Mixture Models, 15pts)

In this problem we will explore expectation-maximization for the Poisson Mixture model. Each observation  $\mathbf{x}_n$  is a non-negative integer  $\mathbb{Z}^*$ . We posit that each observation comes from *one* mixture component. For this problem, we will assume there are  $K$  components. Each component  $k \in \{1, \dots, K\}$  will be associated with a mean  $\lambda_k \in \mathbb{R}^+$ . Finally let the (unknown) overall mixing proportion of the components be  $\boldsymbol{\theta} \in [0, 1]^K$ , where  $\sum_{k=1}^K \theta_k = 1$ .

Our generative model is that each of the  $N$  observations comes from a single component. We encode observation  $n$ 's component-assignment as a one-hot vector  $\mathbf{z}_n \in \{0, 1\}^K$  over components. This one-hot vector is drawn from  $\boldsymbol{\theta}$ ; then,  $\mathbf{x}_n$  is drawn from  $\text{Poisson}(\lambda_{z_n})$ , which simply means that if  $\mathbf{z}_{nj} = 1$  for some  $j \in \{1, \dots, K\}$  (i.e. the  $j$ th element of  $\mathbf{z}_n$  equals 1), then  $\mathbf{x}_n \sim \text{Poisson}(\lambda_j)$ .

Formally, documents are generated in two steps:

$$\begin{aligned}\mathbf{z}_n &\sim \text{Categorical}(\boldsymbol{\theta}) \\ \mathbf{x}_n &\sim \text{Poisson}(\lambda_{z_n})\end{aligned}$$

1. **Intractability of the Data Likelihood** We are generally interested in finding a set of parameters  $\lambda_k$  that maximize the data likelihood  $\log p(\{\mathbf{x}_n\}_{n=1}^N | \{\lambda_k\}_{k=1}^K)$ . Expand the data likelihood to include the necessary sums over observations  $\mathbf{x}_n$  and latents  $\mathbf{z}_n$ . Why is optimizing this loss directly intractable?
2. **Complete-Data Log Likelihood** Define the complete data for this problem to be  $D = \{(\mathbf{x}_n, \mathbf{z}_n)\}_{n=1}^N$ . Write out the complete-data negative log likelihood. Note that optimizing this loss is now computationally tractable if we know  $\mathbf{z}_n$ .

$$\mathcal{L}(\boldsymbol{\theta}, \{\lambda_k\}_{k=1}^K) = -\ln p(D | \boldsymbol{\theta}, \{\lambda_k\}_{k=1}^K).$$

3. **Expectation Step** Our next step is to introduce a mathematical expression for  $\mathbf{q}_n$ , the posterior over the hidden topic variables  $\mathbf{z}_n$  conditioned on the observed data  $\mathbf{x}_n$  with fixed parameters, i.e  $p(\mathbf{z}_n | \mathbf{x}_n; \boldsymbol{\theta}, \{\lambda_k\}_{k=1}^K)$ .
  - **Part 3.A** Write down and simplify the expression for  $\mathbf{q}_n$ .
  - **Part 3.B** Give an algorithm for calculating  $\mathbf{q}_n$  for all  $n$ , given the observed data  $\{\mathbf{x}_n\}_{n=1}^N$  and settings of the parameters  $\boldsymbol{\theta}$  and  $\{\lambda_k\}_{k=1}^K$ .
4. **Maximization Step** Using the  $\mathbf{q}_n$  estimates from the Expectation Step, derive an update for maximizing the expected complete data log likelihood in terms of  $\boldsymbol{\theta}$  and  $\{\lambda_k\}_{k=1}^K$ .
  - **Part 4.A** Derive an expression for the expected complete-data log likelihood in terms of  $\mathbf{q}_n$ .
  - **Part 4.B** Find an expression for  $\boldsymbol{\theta}$  that maximizes this expected complete-data log likelihood. You may find it helpful to use Lagrange multipliers in order to enforce the constraint  $\sum \theta_k = 1$ . Why does this optimized  $\boldsymbol{\theta}$  make intuitive sense?
  - **Part 4.C** Apply a similar argument to find the values of  $\{\lambda_k\}_{k=1}^K$  that maximizes the expected complete-data log likelihood.
5. Suppose that this had been a classification problem, that is, you were provided the “true” categories  $\mathbf{z}_n$  of each document, and you were going to perform the classification by inverting the provided generative model. Could you reuse any of your inference derivations above?

## Solution

**Problem 3** (PCA, 15 pts)

For this problem you will implement PCA from scratch. Using `numpy` to call SVDs is fine, but don't use a third-party machine learning implementation like `scikit-learn`.

We return to the MNIST data set from T3. You have been given representations of 6000 MNIST images, each of which are  $28 \times 28$  greyscale handwritten digits. Your job is to apply PCA on MNIST, and discuss what kinds of structure is found.

As before, the given code loads the images into your environment as a 6000x28x28 array.

- Compute the PCA. Plot the eigenvalues corresponding to the most significant 500 components in order from most significant to least. Make another plot that describes the cumulative proportion of variance explained by the first  $k$  most significant components for values of  $k$ , 1 through 500. How much variance is explained by the first 500 components? Describe how the cumulative proportion of variance explained changes with  $k$ .
- Plot the mean image as well as the images corresponding the first 10 principle components. How does images compare to the cluster centers from K-means? Discuss any similarities and differences.
- Compute the reconstruction error on the data set using the mean image as well as the first 10 principle components. How does this error compare to running K-means and using the cluster centers as the reconstructions for each image? Discuss any similarities and differences.

As in past problem sets, please include your plots in this document. (There may be several plots for this problem, so feel free to take up multiple pages.)

**Solution**

- Name:
- Email:
- Collaborators:
- Approximately how long did this homework take you to complete (in hours):