

Assignment #5

Due: 23:55pm Friday 26 April 2019

Upload at: https://dropbox.cs.princeton.edu/COS324_S2019/HW5

Problem 1 (1pt)

Use the instructions above to put your name on the assignment when you compile this document.

Problem 2 (24pts)

Imagine that we have N data $\{\mathbf{x}_n\}_{n=1}^N$, where $\mathbf{x}_n \in \mathbb{R}^D$. We will model them as having a K -dimensional latent continuous representation with $K \ll D$. One way to do this is to imagine a procedure with the following steps:

1. Generate N latent data $\{\mathbf{z}_n\}_{n=1}^N$ from a K -dimensional spherical Gaussian, i.e., $\mathbf{z}_n \sim \mathcal{N}(0, \mathbb{I}_K)$.
2. Linearly transform the data with a (not necessarily orthonormal) matrix $\mathbf{W}^{D \times K}$ into $\mathbf{W}\mathbf{z}_n$.
3. Add a bit of independent Gaussian noise to this quantity to get $\mathbf{x}_n = \mathbf{W}\mathbf{z}_n + \boldsymbol{\eta}_n$ where $\boldsymbol{\eta}_n \sim \mathcal{N}(\boldsymbol{\mu}, \boldsymbol{\Lambda})$ and $\boldsymbol{\Lambda}$ is a diagonal positive definite matrix.

- A. Write down the marginal distribution of \mathbf{x}_n .
- B. Conditioned on $\boldsymbol{\mu}$, \mathbf{W} , $\boldsymbol{\Lambda}$, and \mathbf{x}_n , what is the distribution over \mathbf{z}_n ?
- C. Write out the “complete data log likelihood” for the n th datum — the log of the joint probability of *both* \mathbf{x}_n and \mathbf{z}_n , conditioned on $\boldsymbol{\mu}$, \mathbf{W} , and $\boldsymbol{\Lambda}$.
- D. Denote the mean and covariance from part (B) as \mathbf{m}_n and \mathbf{S}_n , respectively. Compute the expectation of part (C) under this mean and covariance. This quantity is sometimes called the “expected complete data log likelihood”. The *expectation maximization* (EM) algorithm consists of alternating between maximizing this quantity summed over all the data, with respect to \mathbf{W} and $\boldsymbol{\Lambda}$, and updating \mathbf{m}_n and \mathbf{S}_n via the computation from part (B), across all of the data. The EM algorithm is one way to find maximum likelihood estimate in models with latent variables like \mathbf{z}_n .

Problem 3 (20pts)

Imagine that you have a matrix $\mathbf{X} \in \mathbb{R}^{N \times D}$ that you’d like to model as the (approximate) product of a tall-and-thin matrix $\mathbf{U} \in \mathbb{R}^{N \times K}$ and a short-and-fat matrix \mathbf{V}^T where $\mathbf{V} \in \mathbb{R}^{D \times K}$. Here $K \ll N, D$. We want to find a \mathbf{U} and \mathbf{V} so that $\mathbf{U}\mathbf{V}^T$ is close to \mathbf{X} . There are various ways to measure when two matrices are “close” but we’ll choose the squared *Frobenius norm* in this case, which is essentially the sum of the squares across all entries:

$$\|\mathbf{A}\|_{\text{Fro}}^2 = \sum_{i,j} A_{i,j}^2 = \text{trace}(\mathbf{A}^T \mathbf{A}). \quad (1)$$

Thus we could write our loss function for learning \mathbf{U} and \mathbf{V} as

$$L(\mathbf{U}, \mathbf{V}) = \frac{1}{2} \|\mathbf{X} - \mathbf{U}\mathbf{V}^T\|_{\text{Fro}}^2, \quad (2)$$

and minimize this to find good \mathbf{U} and \mathbf{V} .

- A. Would the solution of this minimization (the pair of matrices \mathbf{U} and \mathbf{V}) be unique? Why or why not?
- B. One approach to minimizing this loss function is to use coordinate descent and alternate between fixing \mathbf{V} and minimizing with respect to \mathbf{U} , and fixing \mathbf{U} and minimizing with respect to \mathbf{V} . The alternating updates in this case have a nice and familiar form, with the caveat that you’ll need to think about the “gradient” with respect to a matrix. Just do the obvious thing and find matrices of the same size as \mathbf{U} and \mathbf{V} containing the partial derivatives for each matrix element. Then set these to zero and solve for the iterative updates.

Problem 4 (20pts)

Take a greyscale image of your choice and load it into Python as a matrix where the entries in the matrix are the pixel intensities. Compute the SVD of the image matrix.

A. Square the singular values, sort them from largest to smallest and plot their cumulative sum as a function of latent dimension K . Is there interesting structure in the curve, e.g., an elbow?

B. Choose 8 different low-rank truncations and create images of each of the reconstructions. How does the visual quality of the reconstruction relate to the plot from part A?

Problem 5 (35pts)

Go out and grab an image data set like:

- CIFAR-10 or CIFAR-100:
<http://www.cs.toronto.edu/~kriz/cifar.html>
- MNIST Handwritten Digits:
<http://yann.lecun.com/exdb/mnist/>
- Small NORB (toys):
<http://www.cs.nyu.edu/~ylclab/data/norb-v1.0-small/>
- Street View Housing Numbers:
<http://ufldl.stanford.edu/housenumbers/>
- STL-10:
<http://cs.stanford.edu/~acoates/stl10/>
- Labeled Faces in the Wild:
<http://vis-www.cs.umass.edu/lfw/>

Figure out how to load it into your environment and turn it into a set of vectors. Center the data (make it zero mean) and compute the covariance matrix. Compute the eigenvalue decomposition of the covariance using `numpy.linalg` or `scipy.linalg`.

A. Produce a Scree plot, similar to the problem above, showing the cumulative sum of eigenvalues (when they are ordered from largest to smallest). Is there interesting structure in this plot indicating that the data are low dimensional?

B. Take the top 16 eigenvectors and put them back into image space, probably by rescaling them to be in $[0, 1]$, reshaping, and then using `imshow`. Produce a figure with these images as subplots.

C. Take 1000 of the data or so and project them onto the top six eigenvectors. Produce three visualizations, each showing a two-dimensional representation arising from PC1 vs. PC2, PC3 vs. PC4, and PC5 vs. PC6. (These are essentially what I showed in lecture.) To make the visualizations, I suggest looking into the `matplotlib.figure` function which lets you place an image anywhere you want in a figure.

Changelog

- 15 April 2019 – Initial version.