# Machine Learning (CS 6140)
# Homework 3 Solutions

Instructor: Ehsan Elhamifar

Due Date: November 27, 2017, 11:45am

**1. PCA Objective Value.** Show that for the optimal solution of fitting a d-dimensional subspace to data using PCA, the objective function value is equal to $\sum_{i \geq d+1} \sigma_i^2$. Here, $\sigma_i$ denotes the $i$-th largest singular value of the mean subtracted data matrix, $\bar{Y} = \begin{bmatrix} y_1 - \bar{y} & \cdots & y_N - \bar{y} \end{bmatrix}$, where $\bar{y} = \frac{1}{N} \sum_{i=1}^{N} y_i$.

**Solution:** Denote the mean subtracted data by $\bar{Y}$. Take the SVD of it to be $\bar{Y} = U\Sigma V^T$, where $U^T U = I_N$ and $V^T V = I_D$. Denote the d-dimensional data by $X$ and let $U_d = U(:, 1:d)$. From class, we have $X = U_d^\top \bar{Y} = \Sigma(1:d,:)V^T$. The PCA objective function can be written as

$$||\bar{Y} - U_d X||_F^2 = ||\bar{Y} - U(1:d,:)\Sigma(1:d,:)V^T||_F^2 = ||\sum_{i=1}^{D} \sigma_i u_i v_i^T - \sum_{i=1}^{d} \sigma_i u_i v_i^T||_F^2$$

$$= ||\sum_{i=d+1}^{D} \sigma_i u_i v_i^T||_F^2$$

$$= \sum_{i=d+1}^{D} \sigma_i^2$$

To get the last equality, we used the fact that $u_i^\top u_j = 0$ for every $i \neq j$ and

$$||\sum_{i=d+1}^{D} \sigma_i u_i v_i^T||_F^2 = \operatorname{tr}\left( (\sum_{i=d+1}^{D} \sigma_i u_i v_i^T)^\top (\sum_{i=d+1}^{D} \sigma_i u_i v_i^T) \right)$$

$$= \operatorname{tr}\left( \sum_{i=d+1}^{D} \sigma_i^2 v_i v_i^\top \right) = \operatorname{tr}\left( \sum_{i=d+1}^{D} \sigma_i^2 v_i^\top v_i \right) \tag{1}$$

$$= \operatorname{tr}(\sum_{i=d+1}^{D} \sigma_i^2) = \sum_{i=d+1}^{D} \sigma_i^2.$$

**2. Locally Linear Embedding.** As discussed in the class, the first step of LLE involves solving the optimization

$$\min_{w_i} \frac{1}{2} ||Y_i w_i||_2^2 \quad \text{s.t.} \quad \mathbf{1}^\top w_i = 1,$$

where $Y_i = \begin{bmatrix} \cdots y_j - y_i \cdots \end{bmatrix}$ for all $j$'s that are neighbors of point $j$ and $\mathbf{1}$ is a vector of all 1's of appropriate dimension. a) Prove that the solution of this optimization is given by $w_i =$

$\frac{(\boldsymbol{Y}_i^\top \boldsymbol{Y}_i)^{-1}\mathbf{1}}{\mathbf{1}^\top(\boldsymbol{Y}_i^\top \boldsymbol{Y}_i)^{-1}\mathbf{1}}$. b) Assume that we have $N$ data points in the $D$-dimensional ambient space. We set the number of nearest neighbors in LLE to be $K$. We try to find a $d$-dimensional representation of the data. What is the computational cost of step 1 and step 2 of LLE, in terms of big $O$ notation (for example, $O(D^p N^q)$)? Explain your answer in details.

**Solution:**

**Part A:** We form the Lagrangian function as

$$L(\boldsymbol{w}_i, \lambda) = \frac{1}{2}\|\boldsymbol{Y}_i\boldsymbol{w}_i\|\|_2^2 - \beta(\mathbf{1}^\top \boldsymbol{w}_i - 1).$$

We set the partial derivative of $L(w_i, \lambda)$ w.r.t. $w_i$ and $\beta$ to zero respectively.

$$\begin{cases} \dfrac{\partial L}{\partial \boldsymbol{w}_i} = \boldsymbol{Y}_i^\top \boldsymbol{Y}_i \boldsymbol{w}_i - \beta\mathbf{1} = 0 \\ \dfrac{\partial L}{\partial \beta} = \mathbf{1}^\top \boldsymbol{w}_i - 1 \quad\quad = 0 \end{cases}$$

$$\Rightarrow \begin{cases} \boldsymbol{w}_i = \beta(\boldsymbol{Y}_i^\top \boldsymbol{Y}_i)^{-1}\mathbf{1} \quad \text{①} \\ \mathbf{1}^\top \boldsymbol{w}_i - 1 = 0 \quad\quad \text{②} \end{cases}$$

Plug $\boldsymbol{w}_i$ in ① into ②. We obtain

$$\Rightarrow \mathbf{1}^\top(\boldsymbol{Y}_i^\top \boldsymbol{Y}_i)^{-1}\mathbf{1}\,\beta = 1$$

$$\Rightarrow \beta = [\mathbf{1}^\top(\boldsymbol{Y}_i^\top \boldsymbol{Y}_i)^{-1}\mathbf{1}]^{-1}$$

Plug the above $\beta$ back into ①. We obtain

$$\boldsymbol{w}_i = \beta(\boldsymbol{Y}_i^\top \boldsymbol{Y}_i)^{-1}\mathbf{1} = \frac{(\boldsymbol{Y}_i^\top \boldsymbol{Y}_i)^{-1}\mathbf{1}}{\mathbf{1}^\top(\boldsymbol{Y}_i^\top \boldsymbol{Y}_i)^{-1}\mathbf{1}}.$$

**Part B:**
Step 1) Find $\{w_{ij}\}_{j\in N_i}$ for $i = 1, ..., N$:
    1. Computing nearest neighbors scales (in the worst case) as $O(DN^2)$, or linearly in the input dimensionality, $D$, and quadratically in the number of data points, $N$.
    2. Computing the reconstruction weights scales as $O(N(k^3 + k^2D))$; this is the number of operations required to solve $\boldsymbol{w}_i$ for each data point.

Step 2) Given $\{w_{ij}\}_{j\in N_i}$, recover $\{x_i \in \mathcal{R}^d\}_{i=1}^N$:
    Computing the bottom eigenvectors scales as $O(dN^2)$, linearly in the number of embedding dimensions, $d$, and quadratically in the number of data points, $N$.

The overall complexity of standard LLE is $O(DN^2 + k^3N + k^2DN)$.

N : number of training data points
D : input dimension
k : number of nearest neighbors
d : output dimension

**3. Laplacian Matrix.** The Laplacian matrix of a graph on $N$ nodes is defined by $\boldsymbol{L} = \boldsymbol{D} - \boldsymbol{W}$, where $\boldsymbol{W}$ is the matrix of weights of the graph, where $W_{ij}$ is the weight between nodes $i$ and $j$. The degree matrix $\boldsymbol{D}$ is also defined as a diagonal matrix whose $i$-th diagonal entry is the sum of the weights of the nodes connected to node $i$, i.e., $D_{ii} = \sum_{j=1}^{N} W_{ij}$. a) Let $\boldsymbol{x} = \begin{bmatrix} x_1 & \cdots & x_N \end{bmatrix}^{\top}$. Compute a closed-form expression for $\boldsymbol{x}^{\top}\boldsymbol{L}\boldsymbol{x}$ that only depends on $W_{ij}$'s. b) Show that the Laplacian matrix is a positive semi-definite matrix. c) Is $L$ an invertible matrix? Prove. d) How many zero singular values does $\boldsymbol{L}$ have?

**Solution:**

**Part A:**

$$\boldsymbol{L} = \boldsymbol{D} - \boldsymbol{W} = \begin{bmatrix} D_{11} & -W_{12} & \dots & -W_{1N} \\ \dots & \dots & \dots & \dots \\ -W_{N1} & -W_{N2} & \dots & D_{NN} \end{bmatrix}$$

$$\boldsymbol{x}^T \boldsymbol{L} \boldsymbol{x} = \boldsymbol{x}^T (\boldsymbol{D} - \boldsymbol{W}) \boldsymbol{x}$$

$$= \sum_{i=1}^{N} D_{ii} x_i^2 - \sum_{i=1}^{N} \sum_{j=1}^{N} W_{ij} x_i x_j$$

$$= \frac{1}{2} \left( \sum_{i=1}^{N} D_{ii} x_i^2 - 2 \sum_{i=1}^{N} \sum_{j=1}^{N} W_{ij} x_i x_j + \sum_{j=1}^{N} D_{jj} x_j^2 \right)$$

$$= \frac{1}{2} \left( \sum_{i=1}^{N} \sum_{j=1}^{N} W_{ij} x_i^2 - 2 \sum_{i=1}^{N} \sum_{j=1}^{N} W_{ij} x_i x_j + \sum_{j=1}^{N} \sum_{i=1}^{N} W_{ij} x_j^2 \right)$$

$$= \frac{1}{2} \sum_{i=1}^{N} \sum_{j=1}^{N} W_{ij} (x_i - x_j)^2$$

**Part B:** $\boldsymbol{x}^T \boldsymbol{L} \boldsymbol{x} = \frac{1}{2} \sum_{i=1}^{N} \sum_{j=1}^{N} W_{ij} (x_i - x_j)^2 \geq 0, \forall \boldsymbol{x} \neq \boldsymbol{0}$. Hence, $\boldsymbol{L}$ is a positive semi-definite matrix.

**Part C:** $\boldsymbol{L}$ is not an invertible matrix. Denote $\boldsymbol{1} = [1, ..., 1]^T$, where the dimension of $\boldsymbol{1}$ is $N \times 1$. We have

$$\boldsymbol{L}\boldsymbol{1} = \begin{bmatrix} D_{11} - \sum_{j=1}^{N} W_{1j} \\ \dots \\ D_{NN} - \sum_{j=1}^{N} W_{Nj} \end{bmatrix} = \boldsymbol{0}$$

Thus, $L$ has an eigenvalue of 0 and is not full rank, hence, it is not invertible.

**Part D:** The number of zero singular values equals to the number of connected components of the graph built using $W$.

**4. Neural Network.** Consider a neural net for a binary classification which has one hidden layer as shown in the figure below. We use a linear activation function $a(z) = cz$ at hidden units and a sigmoid activation function $a(z) = 1/(1 + e^{-z})$ at the output unit to learn the function for $P(y = 1|x, w)$ where $x = (x_1, x_2)$ and $w = (w_1, w_2, \ldots, w_9)$. A) What is the output $P(y = 1|x, w)$ from the above neural net? Express it in terms of $x_i, c$ and weights $w_i$. What is the final classification boundary? B) Draw a neural net with no hidden layer which is equivalent to the given neural net, and write weights $\tilde{w}$ of this new neural net in terms of $c$ and $w_i$. C) Is it true that any multi-layered neural net with linear activation functions at hidden layers can be represented as a neural net without any hidden layer? Explain your answer.
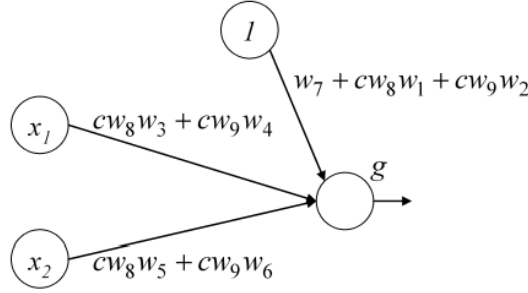


**Solution:**

**Part A:**

$$g(w_7 + w_8 h(w_1 + w_3 x_1 + w_5 x_2) + w_9 h(w_2 + w_4 x_1 + w_6 x_2))$$

$$= \frac{1}{1 + \exp(-(w_7 + cw_8 w_1 + cw_9 w_2 + (cw_8 w_3 + cw_9 w_4)x_1 + (cw_8 w_5 + cw_9 w_6)x_2))}$$

The classification boundary is :

$$w_7 + cw_8 w_1 + cw_9 w_2 + (cw_8 w_3 + cw_9 w_4)x_1 + (cw_8 w_5 + cw_9 w_6)x_2 = 0$$

**Part B:**



**Part C:** Yes. If linear activation functions are used for all the hidden units, output from hidden units will be written as linear combination of input features. Since these intermediate output serves as input for the final output layer, we can always find an equivalent neural net which does not have any hidden layer as seen in the example above.

**4. PCA Implementation.** Implement the PCA and the Kernel PCA algorithm. The input to the PCA algorithm must be the input data matrix $Y \in \mathbb{R}^{D \times N}$, where $D$ is the ambient dimension and $N$ is the number of points, as well as the dimension of the low-dimensional representation, $d$. The output of the PCA must be the basis of the low-dimensional subspace $U \in \mathbb{R}^{D \times d}$, the mean of the subspace $\mu \in \mathbb{R}^D$ and the low-dimensional representations $X \in \mathbb{R}^{d \times N}$. For KPCA, the input must be the Kernel matrix $K \in \mathbb{R}^{N \times N}$, where $N$ is the number of points, as well as the dimension of the low-dimensional representation, $d$. The output of KPCA must be the low-dimensional representations $X \in \mathbb{R}^{d \times N}$.

**5. Kmeans Implementation.** Implement the Kmeans algorithm. The input to the algorithm must be the data matrix $X \in \mathbb{R}^{D \times N}$, the desired number of clusters, $k$, as well as the number of repetitions of kmeans with different random initializations, $r$. The output of the algorithm must be the clustering of the data (indices of points in each group) for the best run of the kmeans, i.e., the run of kmeans that achieves the lowest cost function among all different initializations.

**6. Spectral Clustering Implementation.** Implement the Spectral Clustering algorithm. The input to the algorithm must be the similarity matrix $W \in Re^{N \times N}$, where $W_{ij}$ is the similarity between point $i$ and $j$, and the desired number of clusters, $k$. The output of the algorithm must be the clustering of data (indices of points in each group).

**7. Testing Algorithms on Data.**

**Part A**. Take the dataset1, which consists of 200 points in 40-dimensional space ($Y$ is $40 \times 200$ dimensional matrix).

i) Plot the data by using the first and second features (i.e., first and second rows) of $Y$. Can you see any structure in data? Explain.

ii) Plot the data by using the second and third features (i.e., first and third rows) of $Y$. Can you see any structure in data? Explain.

iii) Apply PCA to data and reduce the dimension of data to $d = 2$. Your output should be a $2 \times 200$ dimensional matrix. Plot the 2-dimensional representation of data using PCA. Can you see any structure in data? Explain.

iv) Apply Kmeans to the 40-dimensional data. Decide about $K$ based on the PCA plot of the data. Show the results by plotting the 2-dimensional data and indicating data in each cluster by a different color. Explain why Kmeans is or is not successful in recovering the true clustering/grouping of the data.

v) Apply Kmeans to the 2-dimensional data. Decide about $K$ based on the PCA plot of the data. Show the results by plotting the 2-dimensional data and indicating data in each cluster by a different color. Explain why Kmeans is or is not successful in recovering the true clustering/grouping of the data.

vi) Explain the similarities and differences of results between part iv and v? Looking at the plots, how much the results for applying Kmeans to 40-dimensional original data differ from the results of applying Kmeans to 2-dimensional representations of data? For general problems, is there any advantage of first applying PCA and then Kmeans to data?

**Part B**. Take the dataset2, which consists of 200 points in 40-dimensional space ($Y$ is $40 \times 200$ dimensional matrix).

i) Plot the data by using the first and second features (i.e., first and second rows) of $Y$. Can you see any structure in data? Explain.

ii) Plot the data by using the second and third features (i.e., first and third rows) of $Y$. Can you see any structure in data? Explain.

iii) Apply PCA to data and reduce the dimension of data to $d = 2$. Your output should be a $2 \times 200$ dimensional matrix. Plot the 2-dimensional representation of data using PCA. Can you see any structure in data? Explain.

iv) Apply Kmeans to the 2-dimensional data obtained by PCA. Show the results by plotting the 2-dimensional data and indicating data in each cluster by a different color. Explain why Kmeans is or is not successful in recovering the true clustering/grouping of the data.

v) Repeat part iv by applying Kernel PCA to data to reduce the dimension of original data to 2. For KPCA, use Gaussian kernel ($k_{ij} = e^{-\|y_i - y_j\|_2^2/(2\sigma^2)}$) and try several values of $\sigma$ (Hint: 2-dimensional PCA plots can help you to pick relatively good $\sigma$ values). Show the results (for the best $\sigma$ you could find) by plotting the 2-dimensional data and indicating data in each cluster by a different color. Explain why KPCA is or is not successful in recovering the true clustering/grouping of the data.

vi) Repeat part iv by applying Spectral to data to reduce the dimension of original data to 2. For spectral clustering, connect each point to its $K$ nearest neighbors using Euclidean distance between points and for any two points $i$ and $j$ connected, set the weights to be $w_{ij} = e^{-\|\boldsymbol{y}_i - \boldsymbol{y}_j\|_2^2/(2\sigma^2)}$. Try several values of $K$ and $\sigma$ (Hint: 2-dimensional PCA plots can help you to pick relatively good $K$ and $\sigma$ values). Show the results (for the best $K$ and $\sigma$ you could find) by plotting the 2-dimensional data and indicating data in each cluster by a different color. Explain why Spectral Clustering is or is not successful in recovering the true clustering/grouping of the data.

**Homework Submission Instructions:**
– Submission of Written Part: You must submit your written report in the class BEFORE CLASS STARTS. The written part, must contain the plots and results of running your code on the provided datasets.
–Submission of Code: You must submit your Python code (.py file) via email, BEFORE CLASS STARTS. For submitting your code, please send an email to me and CC both TAs.
  - The title of your email must be "CS6140: Code: HW3: Your First and Last Name".
  - You must attach a single zip file to your email that contains all python codes and a readme file on how to run your files.
  - The name of the zip file must be "HW3-Code: Your First and Last Name".