

Assignment #5 (Dimensionality Reduction and Clustering)

Instructor: Beilun Wang

Name: Haorui Li, ID: 61518407

Problem Description:

Problem 1: Principle component analysis

Consider n training sample points $\mathbf{x}_1, \dots, \mathbf{x}_n \in \mathbb{R}^p$ such that the sample mean of \mathbf{x} is zero: $\frac{1}{n} \sum_{i=1}^n \mathbf{x}_i = \mathbf{0}$.

(1) PCA: Maximizing the variance

Consider projection vector $\mathbf{u} \in \mathbb{R}^p$, where $\|\mathbf{u}\|_2^2 = 1$. We would like to maximize the sample variance $\tilde{V}[\mathbf{u}^\top \mathbf{x}]$, which is the sample variance of \mathbf{x} projected onto \mathbf{u} . The sample variance of n samples of a variable z is $\tilde{V}[z] = \frac{1}{n} \sum_{i=1}^n z_i^2$, where $\frac{1}{n} \sum_{i=1}^n z_i = 0$. We can write an optimization problem to maximize the sample variance

$$\max_{\mathbf{u}} \tilde{V}[\mathbf{u}^\top \mathbf{x}] \quad (1.1)$$

Reformulate maximum problem (1.1) to the following optimization problem, and define the $p \times p$ matrix Σ :

$$\max_{\mathbf{u}} \mathbf{u}^\top \Sigma \mathbf{u} \quad (1.2)$$

(2) PCA: Minimizing the reconstruction error

Consider the same variables as those defined in question (1). Instead of maximizing the projected variance, we now seek to minimize the reconstruction error. In other words, we would like to minimize the difference between \mathbf{x} and the reconstructed $\mathbf{u}\mathbf{u}^\top \mathbf{x}$. Note that $\mathbf{u}\mathbf{u}^\top \mathbf{x}$ first projects \mathbf{x} onto \mathbf{u} , and then projects this scalar back into the p -dimensional space along the \mathbf{u} axis. Minimizing the reconstruction error can be written as the following optimization problem.

$$\min_{\mathbf{u}} \frac{1}{n} \sum_{i=1}^n \|\mathbf{x}_i - \mathbf{u}\mathbf{u}^\top \mathbf{x}_i\|_2^2 \quad (1.3)$$

Reformulate minimum problem (1.3) to the following optimization problem with the same Σ as in the question (1).

$$\max_{\mathbf{u}} \mathbf{u}^\top \Sigma \mathbf{u} \quad (1.4)$$

(3) Kernel PCA

Kernel PCA is a non-linear dimensionality reduction where a principle vector \mathbf{u}_j is computed as a linear combination of training sample points in the feature space

$$\mathbf{u}_j = \sum_{i=1}^n \alpha_{ij} \phi(\mathbf{x}_i).$$

Computing the principle component of a new point \mathbf{x} can then be done using kernel evaluations

$$z_j(\mathbf{x}) = \langle \mathbf{u}_j, \phi(\mathbf{x}) \rangle = \sum_{i=1}^n \alpha_{ij} \langle \phi(\mathbf{x}_i), \phi(\mathbf{x}) \rangle = \sum_{i=1}^n \alpha_{ij} k(\mathbf{x}_i, \mathbf{x}).$$

You will show that kernel PCA can be represented by a neural network. First we define a *kernel node*. A kernel node with a vector \mathbf{w}_i of incoming weights and an input vector \mathbf{x} computes the output $y = k(\mathbf{x}, \mathbf{w}_i)$ (The kernel function is already known). Show that, given the training sample points $\mathbf{x}_1, \dots, \mathbf{x}_n$, there exists a network with a single hidden layer consisting of several kernel nodes and the output of the network is the kernel principle components $z_1(\mathbf{x}), \dots, z_k(\mathbf{x})$ for a given input \mathbf{x} .

Specify the number of nodes in the input, output and hidden layers, and the weights of the edges in terms of $\alpha_{ij}, \mathbf{x}_1, \dots, \mathbf{x}_n$.

Problem 2: Clustering

Given data set $\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_n \in \mathbb{R}^p$, we want to partition them into k clusters c_1, c_2, \dots, c_k .

(1) Under the assumption of k -means clustering, the loss function could be written as

$$\min_{c_1, c_2, \dots, c_k} \sum_{i=1}^k \sum_{\mathbf{x} \in c_i} \|\mathbf{x} - \boldsymbol{\mu}_i\|^2 \quad (1.5)$$

where $\boldsymbol{\mu}_i = \frac{1}{|c_i|} \sum_{\mathbf{x} \in c_i} \mathbf{x}$ is the mean vector of each cluster. Here $|c_i|$ denotes the number of points in one cluster. Prove that the loss (1.5) is equivalent to the maximum problem below:

$$\max_{c_1, c_2, \dots, c_k} \sum_{i=1}^k \frac{1}{|c_i|} \sum_{\mathbf{x}, \mathbf{y} \in c_i} \mathbf{x}^\top \mathbf{y} \quad (1.6)$$

(Hint: prove by expansion.)

(2) Further show that the problem (1.6) is equivalent to the loss below:

$$\min_{c_1, c_2, \dots, c_k} \sum_{i=1}^k \frac{1}{|c_i|} \sum_{\mathbf{x}, \mathbf{y} \in c_i} \|\mathbf{x} - \mathbf{y}\|^2 \quad (1.7)$$

which represents the weighted within-class total distance.

(3) Reformulate problem (1.6) to the matrix form:

$$\max_{c_1, c_2, \dots, c_k} \text{tr}(\mathbf{H}^\top \mathbf{X} \mathbf{X}^\top \mathbf{H}) \quad (1.8)$$

where each entry \mathbf{H}_{ij} in $\mathbf{H} \in \mathbb{R}^{n \times k}$ is subject to

$$\mathbf{H}_{ij} = \begin{cases} \frac{1}{\sqrt{|c_j|}}, & \mathbf{x}_i \in c_j \\ 0, & \text{otherwise.} \end{cases} \quad (1.9)$$

and $\mathbf{X} = (\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_n)^\top$.

Answer:**Problem 1: Principle component analysis**

(1) Because:

$$\begin{aligned}
\tilde{V} [\mathbf{u}^T \mathbf{x}] &= \frac{1}{n} \sum_{i=1}^n \left[\mathbf{u}^T \mathbf{x}_i - \sum_{j=1}^n \mathbf{u}^T \mathbf{x}_j \right]^2 \\
&= \frac{1}{n} \sum_{i=1}^n \left[\mathbf{u}^T \mathbf{x}_i \mathbf{x}_i^T \mathbf{u} - 2 \mathbf{u}^T \mathbf{x}_i \sum_{j=1}^n \mathbf{x}_j^T \mathbf{u} + \left(\sum_{j=1}^n \mathbf{u}^T \mathbf{x}_j \right)^2 \right] \\
&= \frac{1}{n} \mathbf{u}^T \sum_{i=1}^n \left[\mathbf{x}_i \mathbf{x}_i^T - 2 \mathbf{x}_i \sum_{j=1}^n \mathbf{x}_j^T + \sum_{j=1}^n \sum_{k=1}^n \mathbf{x}_j \mathbf{x}_k^T \right] \mathbf{u} \\
&= \frac{1}{n} \mathbf{u}^T \left[\sum_{i=1}^n \mathbf{x}_i \mathbf{x}_i^T + (n-2) \sum_{i=1}^n \sum_{j=1}^n \mathbf{x}_i \mathbf{x}_j^T \right] \mathbf{u}
\end{aligned}$$

And $\sum_{i=1}^n \sum_{j=1}^n \mathbf{x}_i \mathbf{x}_j^T = \mathbf{O}$

So we can reformulate maximum problem (1.1) to the optimization problem:

$$\max_{\mathbf{u}} \mathbf{u}^T \Sigma \mathbf{u} \quad (2.10)$$

$$s.t. \mathbf{u} = \mathbf{x}^T \mathbf{v} \quad (2.11)$$

$$\mathbf{v}^T \mathbf{v} = 1 \quad (2.12)$$

and the sum matrix is:

$$\Sigma = \frac{1}{N} \sum_{i=1}^N \mathbf{x}_i \mathbf{x}_i^T$$

(2)

$$\begin{aligned}
&\min_{\mathbf{u}} \frac{1}{n} \sum_{i=1}^n \|\mathbf{x}_i - \mathbf{u} \mathbf{u}^T \mathbf{x}_i\|_2^2 \\
&= \min_{\mathbf{u}} \frac{1}{n} \sum_{i=1}^n (\mathbf{x}_i - \mathbf{u} \mathbf{u}^T \mathbf{x}_i)(\mathbf{x}_i - \mathbf{u} \mathbf{u}^T \mathbf{x}_i)^T \\
&= \min_{\mathbf{u}} \frac{1}{n} \sum_{i=1}^n (\mathbf{x}_i^T \mathbf{x}_i - 2 \mathbf{x}_i^T \mathbf{u} \mathbf{u}^T \mathbf{x}_i + \mathbf{x}_i^T \mathbf{u} \mathbf{u}^T \mathbf{x}_i)
\end{aligned}$$

Due to the x is coefficient and the problem can be rewritten as:

$$\begin{aligned}
&\min_{\mathbf{u}} \frac{1}{n} \sum_{i=1}^n \|\mathbf{x}_i - \mathbf{u} \mathbf{u}^T \mathbf{x}_i\|_2^2 \\
&= \max \sum_{i=1}^n \mathbf{x}_i^T \mathbf{u} \mathbf{u}^T \mathbf{x}_i \\
&= \max_{\mathbf{u}} \mathbf{u}^T \Sigma \mathbf{u}
\end{aligned}$$

(3) Fully-connected network:

Input: x node

Hidden layers: n kernel nodes

Output : k nodes and everyone connected to all the hidden layers nodes. The weight between the j_{th} output node and the i_{th} hidden layers node is α_{ij} **Problem 2: Clustering**

(1) Expansion formula we can get:

$$\begin{aligned}
& \min_{c_1, c_2, \dots, c_k} \sum_{i=1}^k \sum_{\mathbf{x} \in c_i} \|\mathbf{x} - \mu_i\|^2 \\
&= \min_{c_1, c_2, \dots, c_k} \sum_{i=1}^k \sum_{\mathbf{x} \in c_i} (\mathbf{x} - \mathbf{u}_i)^T (\mathbf{x} - \mathbf{u}_i) \\
&= \min_{c_1, c_2, \dots, c_k} \sum_{i=1}^k \sum_{\mathbf{x} \in c_i} \mathbf{x}^T \mathbf{x} - 2\mathbf{x}^T \mathbf{u}_i + \mathbf{u}_i^T \mathbf{u}_i \\
&= \max_{c_1, c_2, \dots, c_k} \sum_{i=1}^k \sum_{\mathbf{x} \in c_i} \mathbf{x}^T \mathbf{u}_i \\
&= \max_{c_1, c_2, \dots, c_k} \sum_{i=1}^k \frac{1}{|c_i|} \left(\sum_{\mathbf{x} \in c_i} \mathbf{x}^T \right) \left(\sum_{\mathbf{x} \in c_i} \mathbf{x} \right) \\
&= \max_{c_1, c_2, \dots, c_k} \sum_{i=1}^k \frac{1}{|c_i|} \sum_{\mathbf{x}, \mathbf{y} \in c_i} \mathbf{x}^T \mathbf{y}
\end{aligned}$$

(2)

Just do the reverse operation:

$$\begin{aligned}
& \max_{c_1, c_2, \dots, c_k} \sum_{i=1}^k \frac{1}{|c_i|} \sum_{\mathbf{x}, \mathbf{y} \in c_i} \mathbf{x}^T \mathbf{y} \\
&= \min_{c_1, c_2, \dots, c_k} \sum_{i=1}^k \frac{1}{|c_i|} \sum_{\mathbf{x}, \mathbf{y} \in c_i} \mathbf{x}^T \mathbf{x} - 2\mathbf{x}^T \mathbf{y} + \mathbf{y}^T \mathbf{y} \\
&= \min_{c_1, c_2, \dots, c_k} \sum_{i=1}^k \frac{1}{|c_i|} \sum_{\mathbf{x}, \mathbf{y} \in c_i} \|\mathbf{x} - \mathbf{y}\|^2
\end{aligned}$$

(3)

Like (2):

$$\begin{aligned}
& \max_{c_1, c_2, \dots, c_k} \sum_{i=1}^k \frac{1}{|c_i|} \sum_{\mathbf{x}, \mathbf{y} \in c_i} \mathbf{x}^T \mathbf{y} \\
&= \max_{c_1, c_2, \dots, c_k} \sum_{i=1}^k \frac{1}{|c_i|} \left(\sum_{\mathbf{x} \in c_i} \mathbf{x}^T \right) \left(\sum_{\mathbf{x} \in c_i} \mathbf{x} \right) \\
&= \max_{c_1, c_2, \dots, c_k} \sum_{i=1}^k \left(\sum_{\mathbf{x} \in c_i} \frac{1}{\sqrt{|c_i|}} \mathbf{x}^T \right) \left(\sum_{\mathbf{x} \in c_i} \frac{1}{\sqrt{|c_i|}} \mathbf{x} \right) \\
&= \max_{c_1, c_2, \dots, c_k} \sum_{i=1}^k (X^T \mathbf{h}_i)^T (X^T \mathbf{h}_i) \\
&= \max_{c_1, c_2, \dots, c_k} \text{tr} \left(\mathbf{H}^T \mathbf{X} \mathbf{X}^T \mathbf{H} \right)
\end{aligned}$$