

Rohan Mahendra Chaudhari

USCID : 6675-5653-85

1 Problem 1 - Principal Component Analysis

In the class we showed that PCA is finding the directions with the most variance. In this problem, you will show that PCA is in fact also minimizing reconstruction error in some sense.

1.1 Specifically, suppose we have a dataset $x_1, \dots, x_N \in R^D$ with zero mean, and we would like to compress it into a one-dimensional dataset $c_1, \dots, c_N \in R$. To reconstruct the dataset (approximately), we also keep a direction vector $v \in R^D$ with unit norm (i.e. $\|v\|_2 = 1$) so that the reconstructed dataset is $c_1 v, \dots, c_N v \in R^D$.

Solution

1. Consider the hint of fixing v

For any fixed v , we optimize over each c_n independently:

$$\begin{aligned} \operatorname{argmin}_{c_n} \|\mathbf{x}_n - c_n \mathbf{v}\|_2^2 &= \operatorname{argmin}_{c_n} (c_n^2 \|\mathbf{v}\|_2^2 - 2\mathbf{x}_n^T \mathbf{v} c_n + \|\mathbf{x}_n\|_2^2) \\ &= \operatorname{argmin}_{c_n} (c_n^2 - 2\mathbf{x}_n^T \mathbf{v} c_n), \text{ since } \|\mathbf{v}\|_2^2 = 1 \\ &= \operatorname{argmin}_{c_n} (c_n - \mathbf{x}_n^T \mathbf{v})^2 \\ &= \mathbf{x}_n^T \mathbf{v} \end{aligned}$$

2. Substitute $c = \mathbf{x}_n^T \mathbf{v}$ into $\operatorname{argmin}_{c_1, \dots, c_N, v: \|v\|_2=1} \sum_{n=1}^N \|\mathbf{x}_n - c_n \mathbf{v}\|_2^2$, we see that v is the solution of,

$$\begin{aligned} \operatorname{argmin}_{v: \|v\|_2=1} \sum_{n=1}^N \|\mathbf{x}_n - \mathbf{v} \mathbf{x}_n^T \mathbf{v}\|_2^2 &= \operatorname{argmin}_{v: \|v\|_2=1} \sum_{n=1}^N (\|\mathbf{x}_n\|_2^2 - 2(\mathbf{x}_n^T \mathbf{v})^2 + (\mathbf{x}_n^T \mathbf{v})^2 \|\mathbf{v}_n\|_2^2) \\ &= \operatorname{argmin}_{v: \|v\|_2=1} \sum_{n=1}^N (-(\mathbf{x}_n^T \mathbf{v})^2) \\ &= \operatorname{argmin}_{v: \|v\|_2=1} -\sum_{n=1}^N (\mathbf{v}^T \mathbf{x}_n \mathbf{x}_n^T \mathbf{v}) \\ &= \operatorname{argmin}_{v: \|v\|_2=1} \mathbf{v}^T (\sum_{n=1}^N \mathbf{x}_n \mathbf{x}_n^T) \mathbf{v} \end{aligned}$$

As discussed in the lecture, the above mentioned equation is the top eigenvector of the covariance matrix of the dataset and can thus be concluded as the first principal component.

1.2 Next, you are asked to generalize the same idea to an arbitrary compression dimension $p < D$. Specifically, we would like to compress the same zero-mean dataset into a p -dimensional dataset $\mathbf{c}_1, \dots, \mathbf{c}_N \in \mathbb{R}^p$. To reconstruct the dataset (approximately), we also keep p orthogonal direction vectors $\mathbf{v}_1, \dots, \mathbf{v}_p \in \mathbb{R}^D$ with unit norm. For notational convenience, we stack these vectors together as a matrix $\mathbf{V} \in \mathbb{R}^{D \times p}$ whose j -th column is \mathbf{v}_j .

Solution

1. The reconstructed dataset is $\mathbf{V}\mathbf{c}_1, \dots, \mathbf{V}\mathbf{c}_N$. Thus the optimization problem is

$$\operatorname{argmin}_{\mathbf{c}_1, \dots, \mathbf{c}_N \in \mathbb{R}^p, \mathbf{V} \in \mathbb{R}^{D \times p}: \mathbf{V}^T \mathbf{V} = \mathbf{I}} \sum_{n=1}^N \|\mathbf{x}_n - \mathbf{V}\mathbf{c}_n\|_2^2$$

Here, \mathbf{I} is a identity matrix of order $p \times p$.

2. Here, we do the optimization independently for each \mathbf{c}_n ,

$$\begin{aligned} \operatorname{argmin}_{\mathbf{c}_n} \sum_{n=1}^N \|\mathbf{x}_n - \mathbf{V}\mathbf{c}_n\|_2^2 &= \operatorname{argmin}_{\mathbf{c}_n} (\|\mathbf{x}_n\|_2^2 - 2\mathbf{x}_n^T \mathbf{V}\mathbf{c}_n + \mathbf{c}_n^T \mathbf{V}^T \mathbf{V} \mathbf{c}_n) \\ &= \operatorname{argmin}_{\mathbf{c}_n} (\mathbf{c}_n^T \mathbf{c}_n - 2\mathbf{x}_n^T \mathbf{V}\mathbf{c}_n) \end{aligned}$$

Thus, set the gradient

$$2\mathbf{c}_n - 2\mathbf{V}^T \mathbf{x}_n = 0$$

Thus,

$$\mathbf{c}_n = \mathbf{v}^T \mathbf{x}_n$$

3. Substitute $\mathbf{c}_n = \mathbf{v}^T \mathbf{x}_n$,

$$\begin{aligned} \operatorname{argmin}_{\mathbf{V}: \mathbf{V}^T \mathbf{V} = \mathbf{I}} \sum_{n=1}^N \|\mathbf{x}_n - \mathbf{V}\mathbf{V}^T \mathbf{x}_n\|_2^2 &= \operatorname{argmin}_{\mathbf{V}: \mathbf{V}^T \mathbf{V} = \mathbf{I}} \sum_{n=1}^N (\|\mathbf{x}_n\|_2^2 - 2\mathbf{x}_n^T \mathbf{V}\mathbf{V}^T \mathbf{x}_n + \mathbf{x}_n^T \mathbf{V}\mathbf{V}^T \mathbf{V}\mathbf{V}^T \mathbf{x}_n) \\ &= \operatorname{argmin}_{\mathbf{V}: \mathbf{V}^T \mathbf{V} = \mathbf{I}} \sum_{n=1}^N (-2\mathbf{x}_n^T \mathbf{V}\mathbf{V}^T \mathbf{x}_n + \mathbf{x}_n^T \mathbf{V}\mathbf{V}^T \mathbf{x}_n) \\ &= \operatorname{argmin}_{\mathbf{V}: \mathbf{V}^T \mathbf{V} = \mathbf{I}} \sum_{n=1}^N (\mathbf{x}_n^T \mathbf{V}\mathbf{V}^T \mathbf{x}_n) \\ &= \operatorname{argmin}_{\mathbf{V}: \mathbf{V}^T \mathbf{V} = \mathbf{I}} \sum_{n=1}^N (\mathbf{x}_n^T (\sum_{j=1}^p \mathbf{v}_j \mathbf{v}_j^T) \mathbf{x}_n) \\ &= \operatorname{argmin}_{\mathbf{V}: \mathbf{V}^T \mathbf{V} = \mathbf{I}} \sum_{j=1}^p (\mathbf{v}_j^T (\sum_{n=1}^N \mathbf{x}_n \mathbf{x}_n^T) \mathbf{v}_j) \end{aligned}$$

To solve the last problem, we first find \mathbf{v}_1 to maximize $\mathbf{v}_1^T (\sum_{n=1}^N \mathbf{x}_n \mathbf{x}_n^T) \mathbf{v}_1$ with the constraint that $\|\mathbf{v}_1\|_2^2 = 1$ is the top eigenvector of the covariance matrix.

Next we find \mathbf{v}_2 to maximize $\mathbf{v}_2^T (\sum_{n=1}^N \mathbf{x}_n \mathbf{x}_n^T) \mathbf{v}_2$ with the constraint that \mathbf{v}_2 is orthogonal to \mathbf{v}_1 and $\|\mathbf{v}_2\|_2^2 = 1$. This will be the second eigenvector.

Thus, similarly \mathbf{v}_j will be the j -th eigenvector or the j -th principle component.

2 Problem 2 - Hidden Markov Model

2.1 In the lecture, we discussed how to find the most likely hidden state path given only observations for the first $T_0 < T$ steps. In this problem, you need to generalize the algorithm to the case when you only observe data from an arbitrary subset of time steps. More concretely, for a given subset $M = \{1, \dots, T\}$, find

$$\arg \max_{z_{1:T}} P(Z_{1:T} = z_{1:T} | X_t = x_t, \forall t \in M)$$

Solution

Algorithm 1 Viterbi Algorithm with missing data

Input : Observations $\{x_t\}_{t \in M}$

Output : The most likely path z_1^*, \dots, z_T^*

Initialize : For each $s \in [S]$, compute $\delta_s(1) = \begin{cases} \pi_s b_{s,x_1} & 1 \in M \\ \pi_s & \text{else} \end{cases}$

for $t = 2, \dots, T$ **do**

for each $s \in [S]$ **do**

 Compute

$$\delta_s(t) = \begin{cases} b_{s,x_t} \max_{s'} a_{s',s} \delta_{s'}(t-1) & t \in M \\ \max_{s'} a_{s',s} \delta_{s'}(t-1) & \text{else} \end{cases}$$

$$\Delta_s(t) = \arg \max_{s'} a_{s',s} \delta_{s'}(t-1)$$

Backtracking: Let $z_T^* = \arg \max_s \delta_s(T)$, for $t = T, \dots, 2$, set $z_{t-1}^* = \Delta_{z_t^*}(t)$

2.2 (The next two questions are unrelated to the first one.) Suppose we observe a sequence of outcomes $x_1, \dots, x_{t-1}, x_{t+1}, \dots, x_T$ with the outcome at time t missing ($2 \leq t \leq T-1$). Derive the conditional probability of the state at time t being s , that is,

$$P(Z_t = s | X_{1:t-1} = x_{1:t-1}, X_{t+1:T} = x_{t+1:T})$$

You can use the proportional sign in your derivation. However, to test if you fully understand its meaning, you need to express your final answer **WITHOUT** using the proportional sign.

Solution

$$\begin{aligned} P(Z_t = s | X_{1:t-1} = x_{1:t-1}, X_{t+1:T} = x_{t+1:T}) &\propto P(Z_t = s, X_{1:t-1} = x_{1:t-1}, X_{t+1:T} = x_{t+1:T}) \\ &= P(X_{t+1:T} = x_{t+1:T} | Z_t = s, X_{1:t-1} = x_{1:t-1}) P(Z_t = s, X_{1:t-1} = x_{1:t-1}) \end{aligned}$$

$$\begin{aligned}
&= P(X_{t+1:T} = x_{t+1:T} | Z_t = s) \sum_{s'} P(Z_t = s, Z_{t-1} = s', X_{1:t-1} = x_{1:t-1}) \\
&= \beta_s(t) \sum_{s'} P(Z_t = s | Z_{t-1} = s') P(Z_{t-1} = s', X_{1:t-1} = x_{1:t-1}) \\
&= \beta_s(t) \sum_{s'} a_{s',s} \alpha_s(t-1)
\end{aligned}$$

2.3 Continuing from the last question, derive the conditional probability of the outcome at time t being o

Solution

$$\begin{aligned}
&P(X_t = o | X_{1:t-1} = x_{1:t-1}, X_{t+1:T} = x_{t+1:T}) \\
&= \sum_s P(X_t = o, Z_t = s | X_{1:t-1} = x_{1:t-1}, X_{t+1:T} = x_{t+1:T}) \\
&= \sum_s P(X_t = o | Z_t = s) P(Z_t = s | X_{1:t-1} = x_{1:t-1}, X_{t+1:T} = x_{t+1:T}) \\
&= \sum_s b_{s,o} \sum_{s'} P(Z_t = s, Z_{t-1} = s' | X_{1:t-1} = x_{1:t-1}, X_{t+1:T} = x_{t+1:T})
\end{aligned}$$

Thus, the above answer is written in terms of the equation in Q2.2

On substituting the answer of Q2.2,

$$= \sum_s b_{s,o} (\beta_s(t) \sum_{s'} a_{s',s} \alpha_s(t-1))$$