

Lecture 20: Kernel PCA

April 2020

Lecturer: Steven Wu

Scribe: Steven Wu

Principal Component Analysis

Principal component analysis aims to solve the following optimization problem: given as input a data matrix $X \in \mathbb{R}^{n \times d}$, find an encoder E and decoder D to minimize the following reconstruction error:

$$\min_{D \in \mathbb{R}^{k \times d}, E \in \mathbb{R}^{d \times k}} \|X - XED\|_F^2 \quad (1)$$

where $\|\cdot\|_F$ denotes the Frobenius norm: for any matrix A ,

$$\|A\|_F = \sqrt{\sum_{(i,j)} A_{ij}^2} = \sqrt{\text{tr}(A^T A)}$$

The PCA method solves the problem with the following procedure: compute $X = USV^T$, then return encoder $E = V_k$, decoder $D = V_k^T$, encoded data $XV_k = U_k S_k \in \mathbb{R}^{n \times k}$, and decoded data $XV_k V_k^T$. Note that $V_k V_k^T \in \mathbb{R}^{d \times d}$ performs orthogonal projection onto subspace spanned by V_k .

Last lecture, we showed that the optimization problem can be re-written as

$$\min_{D \in \mathbb{R}^{k \times d}, E \in \mathbb{R}^{d \times k}} \|X - XED\|_F^2 = \min_{D \in \mathbb{R}^{d \times k}, D^T D = I} \|X - XDD^T\|_F^2$$

We also showed that this new objective can be further decomposed

$$\|X - XDD^T\|_F^2 = \|X\|_F^2 - \|XD\|_F^2$$

This means,

$$\min_{D \in \mathbb{R}^{d \times k}, D^T D = I} \|X - XDD^T\|_F^2 \Leftrightarrow \max_{D \in \mathbb{R}^{d \times k}, D^T D = I} \|XD\|_F^2$$

Finally, the objective value of the maximization problem is singular values squared.

$$\max_{D \in \mathbb{R}^{d \times k}, D^T D = I} \|XD\|_F^2 = \|XV_k\|_F^2 = \sum_{j=1}^k s_j^2$$

where s_1, \dots, s_k are the top singular values of X .

Centered PCA. Typically, before running PCA, we replace each x_i with $x'_i = x_i - \bar{x}$, where $\bar{x} = \frac{1}{n} \sum_{i=1}^n x_i$. The objective then becomes

$$\|X'D\|_F^2 = \text{tr}((X'D)^\top(X'D)) = \sum_{i=1}^k (X'De_i)^\top(X'De_i)$$

Note that $\frac{1}{n}(X'De_i)^\top(X'De_i)$ corresponds to the variance on the i -th coordinate after the projection. Therefore, PCA is maximizing the resulting per-coordinate variances.

Power method. How to compute the SVD of $X \in \mathbb{R}^{n \times d}$? We can use the power method to first compute v_1 , u_1 and s_1 . The idea is to compute the top eigenvector and eigenvalue of the matrix $X^\top X$:

- Start with a random vector $y_0 \sim \mathcal{N}(0, I_d)$
- For $t = 1, \dots, T$:
 $y_t \leftarrow X^\top X y_{t-1}$
- $v_1 \leftarrow y_T / \|y_T\|_2$: the first column of V and also the top eigenvector of $X^\top X$
- $s_1 \leftarrow \|Xv_1\|_2$ top singular value
- $u_1 \leftarrow Xv_1 / s_1$

To compute the remainder of triplets (u_i, s_i, v_i) , repeat the same for the residual matrix $X - s_1 u_1 v_1^\top$. Note that we can also apply the power method to the matrix XX^\top for computing its top eigenvector, which is u_1 . This will be useful for the next kernel PCA method.

Kernel PCA

We can find the “high variance” directions in a richer feature space by first apply some feature mapping $\phi: \mathbb{R}^d \rightarrow \mathbb{R}^m$ and then runs PCA over the transformed data. Let $\Phi \in \mathbb{R}^{n \times m}$ such that each row of Φ is given by $\phi(x_i)$. Let $k(\cdot, \cdot)$ be the kernel such that $k(x, y) = \phi(x)^\top \phi(y)$. Kernel PCA then does the following:

- Compute the Gram matrix $G = \Phi\Phi^\top$ and the centered Gram matrix

$$\begin{aligned} \bar{G} &= (\Phi - E\Phi)(\Phi - E\Phi)^\top \\ &= \Phi\Phi^\top - E\Phi\Phi^\top - \Phi\Phi^\top E + E\Phi\Phi^\top E \\ &= G - EG - GE + EGE \end{aligned}$$

where $E \in \mathbb{R}^{n \times n}$ is the matrix with all entries of $1/n$.

- Find the top k eigenvectors of \bar{G} with normalization: call it $A \in \mathbb{R}^{n \times k}$

Original data



Data corrupted with Gaussian noise



Result after linear PCA



Result after kernel PCA, Gaussian kernel

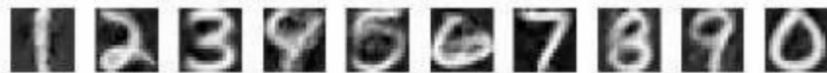


Figure 1: Denoising application of kernel PCA on the digits data set. Image from Haipeng Luo's lecture slide. Another application here.

- Construct the encoded dataset

$$(\Phi - E\Phi)(\Phi - E\Phi)^T A = \bar{G}A$$