

## Lecture 19: Principal Component Analysis

April 2020

Lecturer: Steven Wu

Scribe: Steven Wu

We have thus far focused on supervised learning. The goal there is somewhat clear: we want to find a predictor  $\hat{f}$  based on a training data set  $\{(x_1, y_1), \dots, (x_n, y_n)\}$  such that  $\hat{f}(x)$  matches  $y$  on most instances  $(x, y)$ .

Now we will switch over to unsupervised learning, where the learner only observes a collection unlabeled examples  $x_1, \dots, x_n$ . There are no target labels. The objectives for unsupervised learning can be quite diverse. For example, we might want to encode the data in some compact representation (e.g., lower-dimensional space), recover *latent* structure (e.g., clusters, subspace) in the data, or extract features for subsequent supervised learning tasks. We will study PCA as our first example of unsupervised learning.

## Principal Component Analysis

**SVD** Let us first recall how *singular value decomposition (SVD)* works. In the lectures on linear regression, we study the “thin” version of SVD. Here we will state the “full” factorization of the SVD. Given any matrix  $M \in \mathbb{R}^{n \times d}$ , we want to factorize the matrix as  $M = USV^T$ , such that

- $r$  is the rank of the matrix  $M$ ;
- $U \in \mathbb{R}^{n \times n}$  is orthonormal, that is  $U^T U = I_n$ ;
- $V \in \mathbb{R}^{d \times d}$  is orthonormal, that is  $V^T V = I_d$ ;
- $S \in \mathbb{R}^{n \times d}$  is a rectangular diagonal matrix with entries  $S_{ii} = s_i$  and  $S_{ij} = 0$  for all  $i \neq j$ , where  $(s_1, \dots, s_r, 0, \dots, 0)$  such that  $s_1 \geq s_2 \geq s_3 \dots$  denote the ordered sequence of singular values.

We could also express the factorization as a sum  $M = \sum_{i=1}^r s_i u_i v_i^T$ , where each  $u_i$  is a column vector for  $U$  and each  $v_i$  is a column vector for  $V$ . Note that  $\{u_i\}_{i=1}^r$  spans the column space of  $M$  and  $\{v_i\}_{i=1}^r$  spans the row space of  $M$ .

**Truncated SVD** For any  $k \leq r$ , let  $U_k \in \mathbb{R}^{n \times k}$  be the matrix given by the first  $k$  columns of  $U$ ,  $V_k \in \mathbb{R}^{d \times k}$  be the matrix given by the first  $k$  columns of  $V$ , and  $S_k = \text{diag}(s_1, \dots, s_k)$ . This gives the following factorization restricted to the top- $k$  space:

$$M_k = U_k S_k V_k = \sum_{i=1}^k s_i u_i v_i^T$$

Note that when  $k = r$ , this recovers the thin SVD.

**PCA.** The problem of *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  $A$ ,

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

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

## Analysis

Why does the PCA method solve the problem in (1)? We will prove the following fact, which will be useful for analyzing the problem in (1).

**Fact 0.1.** Let  $X \in \mathbb{R}^{n \times d}$  and  $k \leq r = \text{rank}(X)$ .

$$\min_{M \in \mathbb{R}^{d \times d}, \text{rank}(M)=k} \|X - XM\|_F^2 = \min_{D \in \mathbb{R}^{k \times d}, E \in \mathbb{R}^{d \times k}} \|X - XED\|_F^2 = \min_{D \in \mathbb{R}^{d \times k}, D^\top D = I} \|X - XDD^\top\|_F^2$$

*Proof.* Note that feasible sets have the three minimization problems have the following relationship:

$$\begin{aligned} \{M \in \mathbb{R}^{d \times d} : \text{rank}(M) = k\} &\supseteq \{ED : E \in \mathbb{R}^{d \times k}, D \in \mathbb{R}^{k \times d}\} \\ &\supseteq \{DD^\top : D \in \mathbb{R}^{d \times k}, D^\top D = I_k\} \end{aligned}$$

Since the three minimization problems are increasingly constrained, we have

$$\min_{M \in \mathbb{R}^{d \times d}, \text{rank}(M)=k} \|X - XM\|_F^2 \geq \min_{D \in \mathbb{R}^{k \times d}, E \in \mathbb{R}^{d \times k}} \|X - XED\|_F^2 \geq \min_{D \in \mathbb{R}^{d \times k}, D^\top D = I} \|X - XDD^\top\|_F^2$$

Thus, to establish the stated equality, it suffices to show

$$\min_{D \in \mathbb{R}^{k \times d}, D^\top D = I_k} \|X - XDD^\top\|_F^2 \leq \min_{M \in \mathbb{R}^{d \times d}, \text{rank}(M)=k} \|X - XM\|_F^2$$

Consider any  $M \in \mathbb{R}^{d \times d}$ . With slight abuse of notations, we will also write down the SVD of  $M$  as  $USV^\top$  with  $\text{rank}(M) \leq k$ . We can write

$$\begin{aligned} \|X - XM\|_F^2 &= \|X - XV_k V_k^\top + XV_k V_k^\top - XM\|_F^2 \\ &= \|X - XV_k V_k^\top\|_F^2 + \|XV_k V_k^\top - XM\|_F^2 + 2\text{tr}((X - XV_k V_k^\top)^\top (XV_k V_k^\top - XM)) \end{aligned}$$

If we can show that the trace term is zero, then the above implies

$$\|X - XM\|_F^2 \geq \|X - XV_k V_k^\top\|_F^2 \geq \min_{D \in \mathbb{R}^{d \times k}, D^\top D = I_k} \|X - XDD^\top\|_F^2$$

which will complete the proof.

Now we show that the trace term is zero.

$$\begin{aligned} & \text{tr}((X - XV_k V_k^\top)^\top (XV_k V_k^\top - XM)) \\ &= \text{tr}((I - V_k V_k^\top)^\top X^\top (X - XU_k S_k V_k^\top) V_k V_k^\top) \\ &= \text{tr}(X^\top (X - XU_k S_k V_k^\top) V_k V_k^\top (I - V_k V_k^\top)^\top) \end{aligned} \quad (\text{cyclic property of trace})$$

Let us focus on the last two factors:

$$(I - V_k V_k^\top)^\top V_k V_k^\top = \left( \sum_{j=1}^d v_j v_j^\top - \sum_{j=1}^k v_j v_j^\top \right) \sum_{j=1}^k v_j v_j^\top = 0$$

Thus,  $\text{tr}((X - XV_k V_k^\top)^\top (XV_k V_k^\top - XM)) = 0$ .  $\square$

The fact above suggests that it suffices to solve the following special case of the “encode-decode” problem with the same matrix  $D$ :

$$\min_{D \in \mathbb{R}^{d \times k}, D^\top D = I} \|X - XDD^\top\|_F^2 \quad (2)$$

Furthermore, since

$$\|XDD^\top\|_F^2 = \text{tr}((XDD^\top)^\top (XDD^\top)) = \text{tr}((XD)^\top (XDD^\top D)) = \text{tr}((XD)^\top (XD)) = \|XD\|_F^2$$

Then the objective can be further re-written as

$$\begin{aligned} \|X - XDD^\top\|_F^2 &= \|X\|_F^2 - 2\text{tr}((XDD^\top)^\top X) + \|XDD^\top\|_F^2 \\ &= \|X\|_F^2 - 2\text{tr}(D^\top X^\top X D) + \|XD\|_F^2 \\ &= \|X\|_F^2 - \|XD\|_F^2 \end{aligned}$$

Thus, (2) is equivalent to

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

Finally, it can be shown that

$$\max_{D \in \mathbb{R}^{d \times k}, D^\top D = I} \|XD\|_F^2 = \|XV_k\|_F^2 = \sum_{i=1}^k s_i^2$$

where  $V_k$  from the truncated SVD of  $X$ .

**PCA example.** In this famous paper, the authors performed PCA over the genome data of 1,387 Europeans, and show that the structure of the projected data looks remarkably like the geographic map of Europe. See the Figure below.

