

# Dimensionality reduction

EE219: Large Scale Data Mining

Professor Roychowdhury

Jan 17, 2018

# Summary

- ▶ Review
  - ▶ Binary classifier
  - ▶ Covariance matrix
- ▶ Unsupervised binary classifier
- ▶ PCA
  - ▶ Eigenvalue and eigenvector
  - ▶ Approximation
  - ▶ pick  $k$
- ▶ SVD
  - ▶ SVD approximation
  - ▶ Term-Document matrix

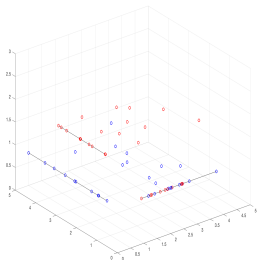
## Review: binary classifier

- ▶ In the linear binary classifier, the input is  $x_1, \dots, x_n \in \mathbb{R}^d$  and the corresponding output is  $y_1, \dots, y_n$ . We want to find a function  $f: \mathbb{R}^d \rightarrow \mathbb{R}$ . It can be viewed as a projection operation.
- ▶ Without loss of generality, we replace  $x_i$  with  $x_i - \bar{x}$ , where  $\bar{x} = \frac{1}{n} \sum_{i=1}^n x_i$  is the center of data.
- ▶ We want to pick a  $w \in \mathbb{R}^d$ , then  $y_i = w^T x_i$ , and 
$$\bar{y}_i = \frac{1}{n} \sum_{i=1}^n w^T x_i = w^T \frac{1}{n} \sum_{i=1}^n x_i = 0$$
- ▶ 
$$\hat{\sigma}_y = \frac{1}{n} \sum_{i=1}^n y_i^2 = \frac{1}{n} \sum_{i=1}^n (w^T x_i)(w^T x_i) = w^T \left( \frac{1}{n} \sum_{i=1}^n x_i x_i^T \right) w$$
- ▶ Define  $R = \frac{1}{n} \sum_{i=1}^n x_i x_i^T$ , then  $R = \text{Cov}(X) = E[XX^T]$ , where 
$$R_{k\ell} = \frac{1}{n} \sum_{i=1}^n x_i(k) x_i(\ell)$$
- ▶ For the original  $x_1, \dots, x_n$  samples before shifting,  $R_{k\ell} = \text{Cov}(X(k), X(\ell)) = E[(X(k) - E[X(k)])(X(\ell) - E[X(\ell)])]$

# Unsupervised binary classifier

If the output or labels  $y_i$  are not given, our aim is to find a  $w$  that maximize  $\hat{\sigma}_y(w) = w^T R w$ .

This is also called Principal Component Analysis(PCA)



As shown in the picture, two projections bring different output distributions. Left one can distinguish data better than the right one.

# PCA

- ▶  $\hat{\sigma}_y(cw) = c^2 \hat{\sigma}_y(w)$ , so if picking  $c \rightarrow \infty$ , you can get unbounded result. Without loss of generality, we add constraint  $\|w\|_2 = 1$  to the optimization problem.
- ▶  $\max_w : w^T R w = \max_w : \frac{w^T R w}{w^T w} = \lambda_{\max}$   
s.t.  $\|w\|_2 = 1$
- ▶  $\lambda_{\max}$  is the largest eigenvalue of  $R$ .
- ▶ How to find the second largest eigenvalue and corresponding eigenvector?
- ▶ How to find  $k$  largest eigenvalues and corresponding eigenvectors? How to pick  $k$ ?

# Eigenvalue and eigenvector

- ▶ A vector  $z \in \mathbb{C}^d$  is an eigenvector of an arbitrary matrix  $R \in \mathbb{R}^{d \times d}$  if  $Rz = \lambda z, \lambda \in \mathbb{C}$ .
- ▶ If  $R = R^T$  and real valued, then  $\lambda$  is real and  $\lambda \geq 0$ . In addition, if  $Rz_1 = \lambda_1 z_1, Rz_2 = \lambda_2 z_2$ , then  $z_1$  and  $z_2$  are orthogonal, or  $z_1^T z_2 = 0$
- ▶ ▶  $R[z_1 \dots z_d] = [z_1 \dots z_d] \begin{bmatrix} \lambda_1 & & \\ & \ddots & \\ & & \lambda_d \end{bmatrix}$  where  $\lambda_1 \lambda_2 \dots \geq \lambda_d$ 
  - ▶  $RU = U\Lambda$ , then  $R = U\Lambda U^T$  shows eigendecomposition of  $R$ , where  $UU^T = I, U^{-1} = U^T$
  - ▶  $w^* = z_1$  is the principal eigenvector corresponding to the largest eigenvalue  $\lambda_1$

# PCA

Use the previous properties to find the second largest eigenvalue:

$$\max_w : \frac{w^T R w}{w^T w} = \lambda_2$$

$$s.t. \|w\|_2 = 1, w^T z_1 = 0$$

For example:

$$\blacktriangleright f : \mathbb{R}^d \rightarrow \mathbb{R}^3$$

$$f(x_i) = \begin{bmatrix} \text{---} & z_1^T & \text{---} \\ \text{---} & z_2^T & \text{---} \\ \text{---} & z_3^T & \text{---} \end{bmatrix} \begin{bmatrix} x_i(1) \\ x_i(2) \\ \vdots \\ x_i(d) \end{bmatrix} = \begin{bmatrix} z_1^T x_i \\ z_2^T x_i \\ z_3^T x_i \end{bmatrix}$$

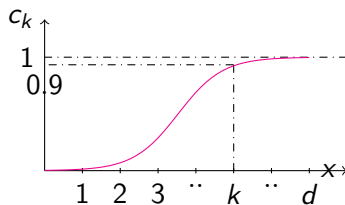
$$\blacktriangleright f: \mathbb{R}^d \rightarrow \mathbb{R}^k$$

$$f(x_i) = \begin{bmatrix} \text{---} & z_1^T & \text{---} \\ \text{---} & z_2^T & \text{---} \\ & \vdots & \\ \text{---} & z_k^T & \text{---} \end{bmatrix} \begin{bmatrix} x_i(1) \\ x_i(2) \\ \vdots \\ x_i(d) \end{bmatrix} = \begin{bmatrix} z_1^T x_i \\ \vdots \\ z_k^T x_i \end{bmatrix}$$

# PCA

How to pick k?

- ▶ Total variance post projection is  $\sum_{i=1}^d \lambda_i$
- ▶ Variance after projecting along the first k eigenvectors is  $\sum_{i=1}^k \lambda_i$
- ▶ The fraction  $c_k = \frac{\sum_{i=1}^k \lambda_i}{\sum_{i=1}^d \lambda_i}$





## Generalization of eigenvalue decomposition

Given  $x_1, x_2, \dots, x_N \in \mathbb{R}^d$ ,  $R = \frac{1}{n} \sum_{i=1}^n x_i x_i^T$ . Let  $Y = \begin{bmatrix} | & \vdots & | \\ x_1 & \vdots & x_n \\ | & \vdots & | \end{bmatrix}$

Then  $R = \frac{1}{n} Y Y^T$ . Instead of dealing with  $Y Y^T$ , we can analyze  $Y$  directly by singular value decomposition.

►  $Y = U \Sigma V^T$

$$= \underbrace{\begin{bmatrix} \mathbf{u}_1 & \dots & \mathbf{u}_r \end{bmatrix}}_{\text{Col } A} \underbrace{\begin{bmatrix} \dots & \mathbf{u}_m \end{bmatrix}}_{\text{Nul } A^T} \begin{bmatrix} \sigma_1 & 0 & \dots & 0 & 0 \dots 0 \\ \dots & & & & \\ 0 & 0 & \dots & \sigma_r & 0 \dots 0 \\ 0 & 0 & \dots & 0 & 0 \dots 0 \\ \dots & & & & \\ 0 & 0 & \dots & 0 & 0 \dots 0 \end{bmatrix} \begin{bmatrix} \mathbf{v}_1^T \\ \mathbf{v}_2^T \\ \dots \\ \mathbf{v}_r^T \\ \mathbf{v}_{r+1}^T \\ \dots \\ \mathbf{v}_n^T \end{bmatrix} \left. \begin{array}{l} \left. \begin{array}{l} \mathbf{v}_1^T \\ \mathbf{v}_2^T \\ \dots \\ \mathbf{v}_r^T \end{array} \right\} \text{Row } A \\ \left. \begin{array}{l} \mathbf{v}_{r+1}^T \\ \dots \\ \mathbf{v}_n^T \end{array} \right\} \text{Nul } A \end{array} \right\}$$

- $U U^T = I, V V^T = I$
- $Y Y^T = U (\Sigma \Sigma^T) U^T$ ,  $U$  are the eigen vectors of  $Y Y^T$
- $Y^T Y = V (\Sigma^T \Sigma) V$ ,  $V$  are the eigen vectors of  $Y^T Y$

# SVD applications

When  $n, d$  have meanings, we can consider SVD.

For example, in the text analysis, we use  $D_1, \dots, D_n$  to represent  $n$  documents and  $T_1, \dots, T_d$  to represent all the words or terms shown in these documents and forgetting their orders.

$$Y = \begin{matrix} & D_1 & \dots & D_j & \dots & D_n \\ \begin{matrix} T_1 \\ \vdots \\ T_i \\ \vdots \\ T_d \end{matrix} & \left( \begin{matrix} & & & & \\ & & & & \\ & & f_{ij} & & \\ & & & & \\ & & & & \end{matrix} \right) \end{matrix} \text{ is called Term/Document}$$

Matrix  $Y$ , where  $Y_{ij}$  represents the number of times  $i$ th term appears on the  $j$ th document.

$$Y = U_{d \times d} \Sigma_{d \times n} V_{n \times n}^T$$

## SVD application

For example, when  $d = 20k$ ,  $n = 100k$ , using SVD can reduce the dimension. In this case,

$$\Sigma = \left[ \begin{array}{cccc|ccc} \sigma_1 & 0 & .. & 0 & 0 & .. & 0 \\ \vdots & & & & \vdots & & \vdots \\ 0 & 0 & .. & \sigma_d & 0 & .. & 0 \end{array} \right]$$

► Approximate  $\hat{Y} = U\hat{\Sigma}V^T$ , where  $\hat{\Sigma} =$

$$\left[ \begin{array}{cccc|ccc} \sigma_1 & 0 & .. & 0 & 0 & .. & 0 \\ \vdots & & & & \vdots & & \vdots \\ 0 & 0 & .. & \sigma_k & 0 & .. & 0 \\ \hline 0 & 0 & .. & 0 & 0 & .. & 0 \end{array} \right]$$

## Term-Document matrix

$$YY^T_{d \times d} = \begin{matrix} & \begin{matrix} D_1 & .. & D_n \end{matrix} \\ \begin{matrix} T_1 \\ \vdots \\ T_i \\ \vdots \\ T_d \end{matrix} & \begin{bmatrix} & & \\ & & \\ - & - & \\ & & \end{bmatrix} \end{matrix} \begin{matrix} \begin{matrix} T_1 & .. & T_j & .. & T_n \end{matrix} \\ \begin{bmatrix} \\ \\ \vdots \\ \end{bmatrix} \end{matrix}$$

- ▶  $(YY^T)_{i,j} = T_i^T T_j$  measures the cooccurrence of the  $j$ th and  $i$ th term. This value measures how similarity they are in the document space. It can be used to cluster terms.
- ▶ Given  $D_j \in \mathbb{R}^d$ ,  $T_j \in \mathbb{R}^n$ , we can project the  $T_i$  and  $D_i$  to  $\mathbb{R}^k$ . This is Latent Semantic Analysis/Indexing. It will be further discussed in the following lecture.