# Dimensionality reduction

## EE219: Large Scale Data Mining

Professor Roychowdhury

Jan 17, 2018

# Summary

- PCA
  - Eigenvalue and eigenvector
  - Approximation
  - pick k
- SVD
  - SVD approximation
  - Term-Document matrix

# Review: Basic Definitions

- We are given a set of feature vectors: $x_1, .. x_n \in \mathbb{R}^d$ and we want reduce the dimension of the data set to a single scalar.

- Without loss of generality, we replace $x_i$ with $x_i - \overline{x}$, where $\overline{x} = \frac{1}{n} \sum\limits_{i=1}^{n} x_i$ is the center of the given data set.

- We want to pick a $w \in \mathbb{R}^d$, and then use it to project each $x_i$ to get $y_i = w^T x_i$. Now each $y_i \in \mathbb{R}$ and is a scalar. Now the average value of $y_i$ is given as:
$$\overline{y_i} = \frac{1}{n} \sum\limits_{i=1}^{n} w^T x_i = w^T \frac{1}{n} \sum\limits_{i=1}^{n} x_i = 0$$

- $\hat{\sigma_y}^2 = \frac{1}{n} \sum\limits_{i=1}^{n} y_i^2 = \frac{1}{n} \sum\limits_{i=1}^{n} (w^T x_i)(w^T x_i) = w^T (\frac{1}{n} \sum\limits_{i=1}^{n} x_i x_i^T) w$

# Review: Covariance Matrix

- Define $R = \frac{1}{n} \sum_{i=1}^{n} x_i x_i^T$, then $R = Cov(X) = E[XX^T]$, where $R_{k\ell} = \frac{1}{n} \sum_{i=1}^{n} x_i(k) x_i(\ell)$
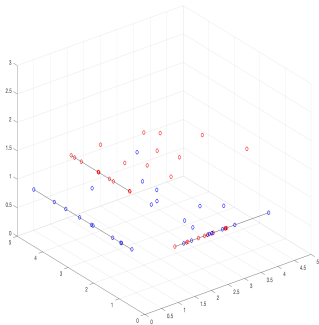
- For the original $x_1, ..x_n$ samples before centering, $R_{k\ell} = Cov(X(k), X(\ell)) = E[(X(k) - E[X(k)])(X(\ell) - E[X(\ell)])]$

# Projections and Clustering

Our aim is to find a w that maximizes $\hat{\sigma_y}(w) = w^T R w$.
This is also called Principal Component Analysis(PCA)



If the data was labeled (blue are data points that belong to one class, and red are data points that belong to a second class) then two different projections lead to two different distributions of blue and red classes: If the projection direction is parallel to the $X - Y$ plane, then the red and blue dots are very well separated. For a second projection (shown in the $Y - Z$ plane) the blue and red classes get all mixed. Thus the projection direction can lead to different clustering.

# PCA

- $\hat{\sigma_y}(cw) = c^2\hat{\sigma_y}(w)$, so if picking $c \to \infty$, you can get unbounded result. Without loss of generality, we add constraint $\|w\|_2 = 1$ to the optimization problem.

- $\max\limits_{\substack{w \\ s.t.\|w\|_2=1}} : w^T R w = \max\limits_{w} : \frac{w^T R w}{w^T w} = \lambda_{max}$

- $\lambda_{max}$ 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 $z \in \mathrm{C}^d$ is an eigenvector of an arbitrary matrix $R \in \mathbb{R}^{d \times d}$ if $Rz = \lambda z, \lambda \in \mathrm{C}$.

- ▶ If $R = R^T$ and real valued and $R$ is positive semidefinite (which covariance matrices always are), 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 $z_1^T z_2 = 0$

- ▶
  - ▶ $R[z_1...z_d] = [z_1...z_d] \begin{bmatrix} \lambda_1 & & \\ & \ddots & \\ & & \lambda_d \end{bmatrix}$ where $\lambda_1 \geq \lambda_2... \geq \lambda_d$

  - ▶ $RU = U\Lambda$, 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:

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

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

Example projections of $x_i \in \mathrm{R}^d$:

- $f : \mathrm{R}^d \to \mathrm{R}^3$

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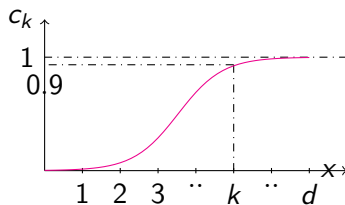
- $f : \mathrm{R}^d \to \mathrm{R}^k$

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

# PCA

How to pick k?

- Total variance post projection is $\sum\limits_{i=1}^{d} \lambda_i$

- Variance after projecting along the first k eigenvectors is $\sum\limits_{i=1}^{k} \lambda_i$

- The fraction $c_k = \dfrac{\sum\limits_{i=1}^{k} \lambda_i}{\sum\limits_{i=1}^{d} \lambda_i}$

# Generalization of eigenvalue decomposition

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

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

- $Y = U\Sigma V^T$

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

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

# SVD applications

When n,d have meanings, we can consider SVD.
For example, in the text analysis, we use $D_1, .. D_n$ to represent n documents and $T_1, .. T_d$ to represent all the words or terms shown in these documents and forgetting their orders.

$$Y = \begin{array}{c} \\ T_1 \\ \vdots \\ T_i \\ \vdots \\ T_d \end{array} \begin{array}{ccccc} D_1 & .. & D_j & .. & D_n \end{array} \\ \left( \begin{array}{ccccc} & & & & \\ & & & & \\ & & f_{ij} & & \\ & & & & \\ & & & & \end{array} \right)$$

is called Term/Document

Matrix, where $Y_{ij}$ represents the number of times ith 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,

$$\Sigma = \left[ \begin{array}{ccccc} \sigma_1 & 0 & .. & 0 & 0..0 \\ \ldots & & & & \ldots \\ 0 & 0 & .. & \sigma_d & 0..0 \end{array} \right]$$

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

$$\left[ \begin{array}{ccccccc} \sigma_1 & 0 & .. & 0 & 0 & 0..0 \\ \ldots & & & & \ldots & 0..0 \\ 0 & 0 & .. & \sigma_k & 0 & 0..0 \\ 0 & 0 & .. & 0 & 0 & 0..0 \end{array} \right]$$

# Term-Document matrix

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

- $(YY^T)_{i,j} = T_i^T T_j$ measures the cooccurrence of the jth and ith 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.