

PCA

From $[U, \Sigma, V] = \text{SVD}(\text{Sigma})$ we get: $U = \begin{bmatrix} | & | & & | \\ u^1 & u^2 & \dots & u^n \\ | & | & & | \end{bmatrix} \in \mathbb{R}^{n \times n}$
取 k 个

$X \in \mathbb{R}^n \xrightarrow{\text{reduce}} Z \in \mathbb{R}^k$

$$Z = \begin{bmatrix} | & | & & | \\ u^1 & u^2 & \dots & u^k \\ | & | & & | \end{bmatrix}_{n \times k}$$

$$Z = Z^T \cdot X = \begin{bmatrix} -(u^1)^T \\ -(u^2)^T \\ \vdots \\ -(u^k)^T \end{bmatrix} \cdot X_{n \times 1}$$

$$\text{Sigma} = \frac{1}{m} \sum_{i=1}^m x^{(i)} \cdot (x^{(i)})^T = \frac{1}{m} \cdot X^T \cdot X$$

$$X = \begin{bmatrix} -(x^{(1)})^T \\ -(x^{(2)})^T \\ \vdots \\ -(x^{(m)})^T \end{bmatrix}$$

\Downarrow
 ~~$Z \in \mathbb{R}^{k \times 1}$~~
 $Z \in \mathbb{R}^{k \times 1}$

$$U_{\text{reduce}} = U(:, 1:k)$$

$$Z = U_{\text{reduce}}^T \cdot X$$

mathematics.

* Our goal: retain as much as possible info. after dimensional reduction

* info: $S_X = \text{Sigma} = \frac{1}{n} \cdot X^T \cdot X$... Covariance matrix

* Solve.

* What is a point? $X \in \mathbb{R}^D$, $Z \in \mathbb{R}^M$, $M < D$

linear dimensionality reduction

$$X = [x_1, \dots, x_n]^T \quad X \in \mathbb{R}^{N \times D}$$

$$Z = [z_1, \dots, z_n]^T \quad Z \in \mathbb{R}^{n \times m}$$

$$Z = U \cdot X$$

* problem: $\max_U U^T \cdot S_X \cdot U \quad \text{s.t.} \quad U^T \cdot U = I$

$$\text{Prop } S_Z = \frac{1}{n} \cdot Z^T \cdot Z \in \mathbb{R}^{m \times m}$$

$$\Leftrightarrow \max S_z = \max \frac{1}{N} \cdot z^T \cdot z \Leftrightarrow \max \frac{1}{N} \cdot (Xu)^T (Xu) \Leftrightarrow \max u^T X^T X u$$

$$\frac{\partial L}{\partial X} = 0 \quad \text{用 Lagrain.} \quad L(u, \lambda) = u^T S_X u + \lambda (I - u^T u) \quad \text{①} \quad \max u^T \cdot S_X \cdot u$$

$$\text{s.t. } u^T \cdot u = I$$

$$\text{EP } S_X \cdot u = \lambda u \quad \star$$

$$\Downarrow$$

$$S_X u_i = \lambda_i u_i$$

$$\star \text{ SVD} \quad \text{Total var} = \sum_{i=1}^N \lambda_i$$

$$\text{Retained var} = \sum_{i=1}^D \lambda_i \quad \Rightarrow \quad \% \text{ info} = \frac{\sum_{i=1}^D \lambda_i}{\sum_{i=1}^N \lambda_i} \quad 90\% - 95\%$$

$$\text{取 } \lambda_1 > \lambda_2 > \dots > \lambda_D \quad \text{取 top } M \quad \Rightarrow \quad U = [u_1 \dots u_M]^T$$

kernel PCA

$\phi: R^p \rightarrow R^m$
feature mapping

$x \in R^D$
 $\phi(x) \in R^m$
 $m < \infty$

max var retained we get

$$S_\phi \cdot u = \lambda \cdot u$$

$$\frac{1}{N} \cdot \phi^T \phi u = \lambda u$$

$$u = \frac{1}{\lambda N} \cdot \phi^T \phi u$$

2 properties of kernel:

① symmetry: $k(x_m, x_n) = k(x_n, x_m)$

② PSD: $\sum_m \sum_n V_m V_n \cdot k(x_m, x_n) \geq 0$

$$\phi^T(x_m) \phi^T(x_n) = k(x_m, x_n)$$

$$z = \phi u = \phi \cdot \frac{1}{\lambda n} \cdot \phi^T \phi u = \frac{1}{\lambda n} (\phi \phi^T) (\phi u)$$

$$z = \frac{1}{\lambda n} \cdot k z \Rightarrow k z = (\lambda n) z \quad \text{if } k z_i = \lambda' z_i \quad \text{just diag } k.$$

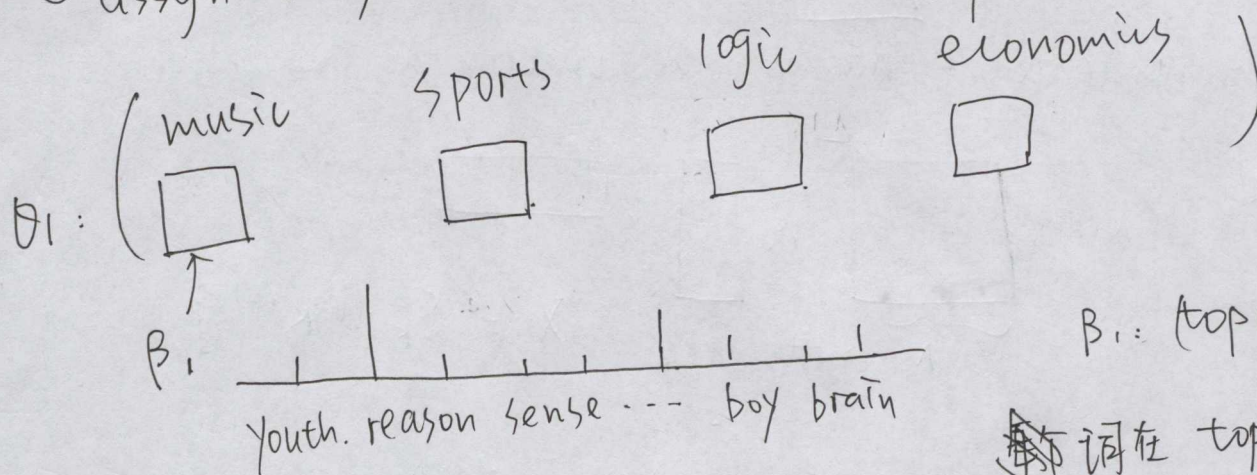
steps

- ① determine eigenvectors
- ② sort by eigenvectors
- ③ select top M vectors for z

Topic modeling

• a prob topic model:

- ① learn dist of words \rightarrow Topics 每个 topic 是一个 ~~词~~ 词集
- ② learn dist of topics for each doc \rightarrow 学习每个 doc 的 topic 分布
- ③ assign every word in a doc to a topic.



~~词~~ 词在 topic 为 music 上有一个 dist.

θ_1 : 一个 document. 含有若干 topics 的 dist.

LDA: 1) generate each topic $\beta_k \sim \text{Dirichlet}(\gamma)$, $k=1 \dots K$

2) generate dist on topic for each doc.

$\theta_d \sim \text{Dirichlet}(\alpha)$, $d=1 \dots D$

3) for the n -th word in the d -th doc:

a) allocate the word to a topic b) generate the word

Dirichlet dist: a continuous dist on discrete prob vector.

Let β_k be a prob vector and r a positive para vector.

$$p(\beta_k | r) = \frac{\Gamma(\sum_v r_v) \cdot \prod_{v=1}^V \beta_{k,v}^{r_v-1}}{\prod_{v=1}^V \Gamma(r_v)}$$

$r \uparrow$, var \downarrow .

对LDA来说, 会选得小的 γ, α . & focus on 很小的 subset of word topics.

LDA output ① topics. ② 每个doc的 topic dist.

Q: for a particular doc, what is $p(X_{dn}=i | \beta, \theta_d)$

A: $p(X_{dn}=i | \beta, \theta) = \sum_{k=1}^K p(X_{dn}=i, c_{dn}=k | \beta, \theta_d)$

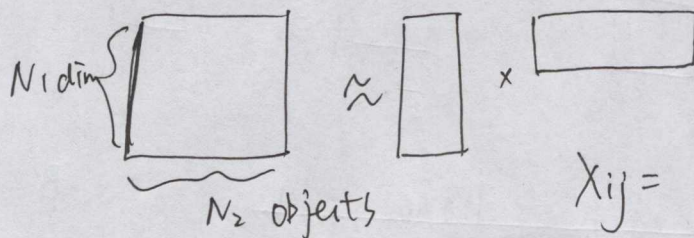
某个word 分配到第几个topic
 $c_{dn} \sim p(\theta_d)$
generate the word.

$X_{dn} \sim D(\beta_{c_{dn}})$

NMF

• LDA is an instance of non-negative matrix factorization.

• NMF 用处:



$$X_{ij} = \sum_k W_{ik} \cdot H_{kj}$$

• data

text: X_{ij} : # times word i appears in document j .

image: ~~X_{ij}~~ : put each vectorized $N \times M$ image of a face.
on a col of X

• obj function

$$\textcircled{1} \|X - WH\|^2 = \sum_i \sum_j (X_{ij} - (WH)_{ij})^2$$

$$\textcircled{2} D(X \| WH) = - \sum_i \sum_j [X_{ij} \ln (WH)_{ij} - (WH)_{ij}]$$

W, H non negative values