

# Probabilistic Latent Semantic Analysis

Seungjin Choi

Department of Computer Science and Engineering  
Pohang University of Science and Technology  
77 Cheongam-ro, Nam-gu, Pohang 37673, Korea  
[seungjin@postech.ac.kr](mailto:seungjin@postech.ac.kr)

# Outline

- ▶ Singular value decomposition
- ▶ Latent semantic analysis (a.k.a. latent semantic indexing)
- ▶ Probabilistic latent semantic analysis (a.k.a. probabilistic latent semantic indexing)

# Range Space and Null Space

## Definition

For any  $\mathbf{A} \in \mathbb{R}^{m \times n}$ , its **range space**  $\mathcal{R}(\mathbf{A})$  and **null space**  $\mathcal{N}(\mathbf{A})$  are defined as follows:

$$\begin{aligned}\mathcal{R}(\mathbf{A}) &= \{\mathbf{y} \in \mathbb{R}^m \mid \mathbf{y} = \mathbf{A}\mathbf{x} \text{ for some } \mathbf{x} \in \mathbb{R}^n\} \\ \mathcal{N}(\mathbf{A}) &= \{\mathbf{x} \in \mathbb{R}^n \mid \mathbf{A}\mathbf{x} = \mathbf{0}\}.\end{aligned}$$

**Examples:** What are range space and null space of  $\mathbf{A}$  and  $\mathbf{B}$ ?

$$\mathbf{A} = \begin{bmatrix} 0 & 0 \\ 0 & 0 \end{bmatrix} \quad \mathbf{B} = \begin{bmatrix} 1 & 0 \\ 0 & 0 \end{bmatrix}$$

# Linear Algebraic Equations

A linear algebraic equation has the form

$$\mathbf{Ax} = \mathbf{b}, \tag{1}$$

where  $\mathbf{A} \in \mathbb{R}^{m \times n}$ ,  $\mathbf{x} \in \mathbb{R}^n$ , and  $\mathbf{b} \in \mathbb{R}^m$ .

## Theorem

*A solution of  $\mathbf{Ax} = \mathbf{b}$  exists if and only if  $\mathbf{b} \in \mathcal{R}(\mathbf{A})$ .*

## Theorem

*Let  $\mathbf{x}^\circ$  be a particular solution to (1), then  $\mathbf{x} = \mathbf{x}^\circ + \mathcal{N}(\mathbf{A}) = \{\mathbf{x}^\circ + \mathbf{w} \mid \mathbf{w} \in \mathcal{N}(\mathbf{A})\}$  is the general solution.*

**Remark:** Note that  $\mathcal{N}(\mathbf{A})$  is the set of homogeneous solutions to (1). Therefore the set of all solutions is the sum of a particular solution and the set of homogeneous solutions.

# Singular Value Decomposition

## Theorem (SVD)

Let  $\mathbf{A} \in \mathbb{R}^{m \times n}$ . Then there exist orthogonal matrices

$$\begin{aligned}\mathbf{U} &= [\mathbf{u}_1, \mathbf{u}_2, \dots, \mathbf{u}_m] \in \mathbb{R}^{m \times m} \\ \mathbf{V} &= [\mathbf{v}_1, \mathbf{v}_2, \dots, \mathbf{v}_n] \in \mathbb{R}^{n \times n}\end{aligned}$$

such that

$$\mathbf{U}^T \mathbf{A} \mathbf{V} = \text{diag}(\sigma_1, \sigma_2, \dots, \sigma_p) \in \mathbb{R}^{m \times n} \quad (2)$$

where  $p = \min\{m, n\}$  and  $\sigma_1 \geq \sigma_2 \geq \dots \geq \sigma_p \geq 0$ . The scalars  $\{\sigma_i \in \mathbb{R}\}$  are called *singular values*.

$$\underline{\text{SVD } \mathbf{A} = \mathbf{U} \mathbf{\Sigma} \mathbf{V}^T}$$

# Low-Rank Approximation

## Theorem (Eckart and Young, 1936)

Suppose that  $\mathbf{A}$  is an  $m \times n$  matrix of rank  $r$ , with singular values  $\sigma_1 \geq \sigma_2 \geq \dots \geq \sigma_r$ . and SVD  $\mathbf{A} = \mathbf{U}\mathbf{\Sigma}\mathbf{V}^T$ . Then the best rank- $k$  approximation of  $\mathbf{A}$  is  $\mathbf{A}_k = \mathbf{U}_1\mathbf{\Sigma}_1\mathbf{V}_1^T$ , where

$$\begin{aligned}\mathbf{U} &= [\mathbf{U}_1, \mathbf{U}_2], \\ \mathbf{V} &= [\mathbf{V}_1, \mathbf{V}_2], \\ \mathbf{\Sigma}_1 &= \text{diag}(\sigma_1, \dots, \sigma_k).\end{aligned}$$

That is,

$$\|\mathbf{A} - \mathbf{A}_k\|_F = \min \left\{ \|\mathbf{A} - \hat{\mathbf{A}}\|_F \mid \hat{\mathbf{A}} \in \mathbb{R}^{m \times n}, \text{rank}(\hat{\mathbf{A}}) = k \right\}.$$

# Proof of Eckart-Young Theorem

Note that

$$\begin{aligned}\|\mathbf{A} - \hat{\mathbf{A}}\|_F &= \|\mathbf{U}\mathbf{\Sigma}\mathbf{V}^T - \hat{\mathbf{A}}\|_F \\ &= \|\mathbf{\Sigma} - \underbrace{\mathbf{U}^T \hat{\mathbf{A}} \mathbf{V}}_{\mathbf{N}}\|_F.\end{aligned}$$

Direct calculations give

$$\begin{aligned}\|\mathbf{\Sigma} - \mathbf{N}\|_F^2 &= \sum_{i,j} |\Sigma_{i,j} - N_{i,j}|^2 \\ &= \sum_{i=1}^r |\sigma_i - N_{i,i}|^2 + \sum_{i>r} |N_{i,i}|^2 + \sum_{i \neq j} |N_{i,j}|^2,\end{aligned}$$

which is minimal when all the non-diagonal entries of  $\mathbf{N}$  are equal to zero and so are all  $N_{i,i}$  for  $i > r$ . The minimum of  $\sum_{i=1}^r |\sigma_i - N_{i,i}|^2$  is attained when  $\sigma_i = N_{i,i}$  for  $i = 1, \dots, k$  and all other  $N_{i,j}$  are zero.

# Least Squares Problem

Consider a matrix  $\mathbf{A} \in \mathbb{R}^{m \times n}$  with  $\text{Rank}(\mathbf{A}) = r$ . Then

$$\begin{aligned}\mathcal{R}(\mathbf{A}) &= \mathcal{R}(\mathbf{U}_1) = \text{sp}[\mathbf{u}_1, \mathbf{u}_2, \dots, \mathbf{u}_r] \\ \mathcal{N}(\mathbf{A}) &= \text{sp}[\mathbf{v}_{r+1}, \mathbf{v}_{r+2}, \dots, \mathbf{v}_n] = \mathcal{R}(\mathbf{V}_2).\end{aligned}$$

The least squares (LS) problem is as follows:

$$\min_{\mathbf{x} \in \mathbb{R}^n} \|\mathbf{b} - \mathbf{Ax}\|^2. \quad (3)$$

Setting  $\nabla_{\mathbf{x}} \|\mathbf{b} - \mathbf{Ax}\|^2 = 0$  leads to

$$\mathbf{A}^T \mathbf{Ax} = \mathbf{A}^T \mathbf{b}. \quad (4)$$



## Remarks

- ▶ From SVD, one can see that the particular solution to (4) is  $\mathbf{x}^\circ = \mathbf{V}_1 \boldsymbol{\Sigma}^{-1} \mathbf{U}_1^T \mathbf{b}$ .
- ▶ The complete solution is  $\mathbf{x} = \mathbf{x}^\circ + \mathcal{N}(\mathbf{A}^T \mathbf{A})$ . Note that  $\mathcal{N}(\mathbf{A}^T \mathbf{A}) = \mathcal{N}(\mathbf{A}) = \text{sp}[\mathbf{V}_2]$ . Hence  $\mathbf{x} = \mathbf{V}_1 \boldsymbol{\Sigma}^{-1} \mathbf{U}_1^T \mathbf{b} + \mathbf{V}_2 \mathbf{w}$  for some  $\mathbf{w} \in \mathbb{R}^r$ .
- ▶ Since  $\mathbf{V}_1 \boldsymbol{\Sigma}^{-1} \mathbf{U}_1^T \mathbf{b} \perp \mathbf{V}_2 \mathbf{w}$  for all  $\mathbf{w}$ ,  $\mathbf{x}^\circ = \mathbf{V}_1 \boldsymbol{\Sigma}^{-1} \mathbf{U}_1^T \mathbf{b}$  is the **MMSE solution** for which  $\|\mathbf{x}\|$  is the smallest among all solutions.
- ▶ The **pseudo-inverse** of  $\mathbf{A}$  is

$$\mathbf{A}^\dagger = \mathbf{V}_1 \boldsymbol{\Sigma}^{-1} \mathbf{U}_1^T. \quad (5)$$

- ▶ The condition number of  $\mathbf{A}$  is given by  $\frac{\sigma_1}{\sigma_r}$ .

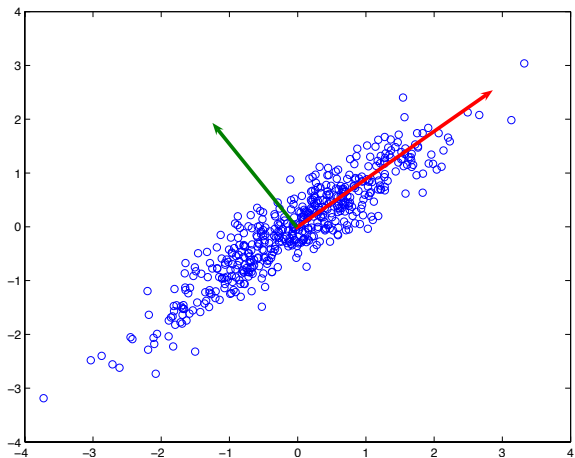
# PCA

- ▶ Principal component analysis (PCA) is a well-established technique for dimension reduction. Its applications include data compression, image processing, data visualization, exploratory data analysis, pattern recognition, and time series prediction.
- ▶ The most common derivation of PCA is in terms of a orthogonal projection which maximizes the variance in the projected space. Given a set of  $m$ -dimensional observation vector,  $\{\mathbf{x}_t\}$ , the PCA aims at finding a orthogonal linear projection  $\mathbf{y} = \mathbf{W}\mathbf{x}$  such that the variance of  $\mathbf{y} \in \mathbb{R}^q$  ( $q < d$ ) is maximized. The  $i$ th element of  $\mathbf{y}$  is called the  *$i$ th principal component*.
- ▶ Alternatively, the PCA provides an orthogonal linear projection which minimize the squared reconstruction error  $\sum_t \|\mathbf{x}_t - \hat{\mathbf{x}}_t\|^2$ . Thus the PCA is an optimal linear encoding in MS sense.

# PCA and SVD

- ▶ It was shown that the  $i$ th row vector of  $\mathbf{W}$  denoted by  $\mathbf{w}_i^T$  corresponds to the normalized eigenvector associated with the  $i$ th largest eigenvalue of the covariance matrix  $\mathbf{R}_x = E\{\mathbf{x}\mathbf{x}^T\}$
- ▶ Principal components can be found by SVD, linear neural networks, or probabilistic methods.
- ▶ The SVD of  $\mathbf{R}_x$  has the form  $\mathbf{R}_x = \mathbf{U}\mathbf{\Sigma}\mathbf{V}^T$ . Then we select  $n$  column vectors to construct  $\mathbf{U}_1 = [\mathbf{u}_1, \dots, \mathbf{u}_n]$ . The PCA transform leads to  $\mathbf{y} = \mathbf{U}_1^T \mathbf{x}$ .

# PCA: An Example



# Spectral Decomposition (Eigen-Decomposition)

Given a symmetric matrix  $\mathbf{C} \in \mathbb{R}^{m \times m}$ , its spectral decomposition is given by

$$\mathbf{C} = \lambda_1 \mathbf{u}_1 \mathbf{u}_1^T + \lambda_2 \mathbf{u}_2 \mathbf{u}_2^T + \cdots + \lambda_m \mathbf{u}_m \mathbf{u}_m^T,$$

where  $|\lambda_1| \geq |\lambda_2| \geq \cdots \geq |\lambda_m|$  are eigenvalues of  $\mathbf{C}$  and  $\mathbf{u}_i$  are associated eigenvectors.

# Power Iteration

- ▶ The power iteration is a classical method which finds the largest eigenvector (associated with the largest eigenvalue) of a matrix  $\mathbf{C} \in \mathbb{R}^{m \times m}$ .
- ▶ Given a symmetric matrix  $\mathbf{C} \in \mathbb{R}^{m \times m}$  (hence its eigenvalues are real), the power iteration starts from a nonzero vector  $\mathbf{w}(0)$  and iteratively updates  $\mathbf{w}(t)$  by

$$\tilde{\mathbf{w}}(t+1) = \mathbf{C}\mathbf{w}(t), \quad (6)$$

$$\mathbf{w}(t+1) = \frac{\tilde{\mathbf{w}}(t+1)}{\|\tilde{\mathbf{w}}(t+1)\|_2}, \quad (7)$$

where  $\|\cdot\|_2$  represents Euclidean norm.

- ▶ Combining (6) and (7) leads to the updating rule which has the form

$$\mathbf{w}(t+1) = \mathbf{C}\mathbf{w}(t) \left[ \mathbf{w}^T(t)\mathbf{C}^2\mathbf{w}(t) \right]^{-\frac{1}{2}}. \quad (8)$$

- ▶ Assume that  $\mathbf{C}$  has an unique eigenvalue of maximum modulus  $\lambda_1$  associated with the leading eigenvector  $\mathbf{u}_1$ . Then the power iteration (8) leads  $\mathbf{w}(t)$  to converge to  $\mathbf{u}_1$ .

# Deflation

- ▶ Suppose that we are interested in computing eigenvectors of the data covariance matrix  $\mathbf{C} = \langle \mathbf{x}\mathbf{x}^T \rangle$ .
- ▶ The power iteration is applied to  $\mathbf{C}$  for extracting its first eigenvector. Our question arises, "How can we compute the second eigenvector of  $\mathbf{C}$  using the power iteration?"
- ▶ The **deflation** method is a common numerical technique for computing several eigenvalues and eigenvectors of  $\mathbf{C}$ .
- ▶ Assume that the first eigenvector is already computed. Then the output value can be deflated by the following transformation:

$$\tilde{\mathbf{x}} = (\mathbf{I} - \mathbf{u}_1\mathbf{u}_1^T) \mathbf{x}. \quad (9)$$

- ▶ One can easily see that  $\langle \tilde{\mathbf{x}}\tilde{\mathbf{x}}^T \rangle = \sum_{i=2}^m \lambda_i \mathbf{u}_i \mathbf{u}_i^T$ .
- ▶ The power iteration is applied to the deflated data, in order to extract the second eigenvector of  $\mathbf{C}$ .

# Term-Document Matrix

A term-document matrix  $\mathbf{X} \in \mathbb{R}^{D \times N}$  is a collection of vector space representations of documents, where rows are terms (words) and columns are documents

$$X_{ij} = t_{ij} \log \left( \frac{N}{idf_i} \right),$$

where  $t_{ij}$  is the term frequency of word  $i$  in document  $j$  and  $idf_i$  is the number of documents containing word  $i$ .

We write

$$\begin{aligned} \mathbf{X} &= [\mathbf{d}_1, \mathbf{d}_2, \dots, \mathbf{d}_N] \quad (\text{document vectors}) \\ &= \begin{bmatrix} \mathbf{t}_1^\top \\ \mathbf{t}_2^\top \\ \vdots \\ \mathbf{t}_D^\top \end{bmatrix} \quad (\text{term vectors}). \end{aligned}$$



# Latent Semantic Analysis

- ▶ A method for automatic indexing and retrieval, uncovering latent semantic structure of a term-document matrix.
- ▶ SVD of  $\mathbf{X} \in \mathbb{R}^{D \times N}$  is given by  $\mathbf{X} \approx \mathbf{U}\mathbf{\Sigma}\mathbf{V}^T$  where  $\mathbf{U} \in \mathbb{R}^{D \times K}$ ,  $\mathbf{\Sigma} \in \mathbb{R}^{K \times K}$ , and  $\mathbf{V} \in \mathbb{R}^{N \times K}$ .

$$\begin{matrix} & & N & & K & & K & & N \\ & & \left[ \right. & & \left[ \right. & & \left[ \right. & & \left[ \right. \\ & & \mathbf{X} & \approx & \mathbf{U} & \mathbf{\Sigma} & \mathbf{V}^T & & \\ & & \left. \right] & & \left. \right] & & \left. \right] & & \left. \right] \\ D & & & & D & & K & & \end{matrix}$$

- ▶ Comparing documents:

$$\mathbf{\Sigma}^{-1}\mathbf{U}^T\mathbf{X} = \mathbf{\Sigma}^{-1}\mathbf{U}^T\mathbf{U}\mathbf{\Sigma}\mathbf{V}^T = \mathbf{V}^T = [\hat{\mathbf{d}}_1, \dots, \hat{\mathbf{d}}_N].$$

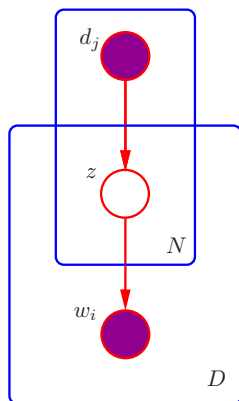
- ▶ Comparing terms:

$$\mathbf{X}\mathbf{V}\mathbf{\Sigma}^{-1} = \mathbf{U}\mathbf{\Sigma}\mathbf{V}^T\mathbf{V}\mathbf{\Sigma}^{-1} = \mathbf{U} = [\hat{\mathbf{t}}_1, \dots, \hat{\mathbf{t}}_D]^T.$$

# Probabilistic Latent Semantic Analysis

- ▶ The key idea in latent semantic analysis is to map high-dimensional count vectors (co-occurrence data, dyadic data) to a lower-dimensional representation in a so-called **latent semantic space**.
- ▶ The goal of LSA is to find a data mapping which reveals semantical relations between the entries of interest.
- ▶ Probabilistic latent semantic analysis (PLSA) is a probabilistic variant of LSA:
  - ▶ Has a sound statistical foundation;
  - ▶ Defines a proper **generative model** of the data.
- ▶ **Aspect model**: A latent variable model for co-occurrence data which associates an unobserved class  $z \in \{z_1, \dots, z_K\}$  with each occurrence of a word  $w \in \{w_1, \dots, w_D\}$  in a document  $d \in \{d_1, \dots, d_N\}$ .

# PLSA: Graphical Representation

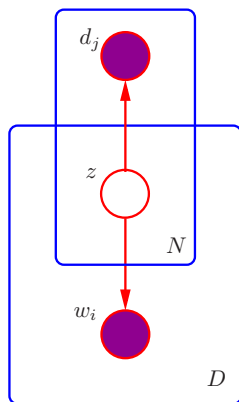


- ▶ A document  $d_j$  and a term (word)  $w_i$  are conditionally independent given an unobserved topic  $z$ :

$$p(w_i, d_j) = p(d_j) \left[ \sum_z p(w_i|z)p(z|d_j) \right].$$

- ▶ Generation process:
  - ▶ Select a document  $d_j$  with probability  $p(d_j)$ .
  - ▶ Pick a latent class (topic) with probability  $p(z|d_j)$ .
  - ▶ Generate a word  $w_i$  with probability  $p(w_i|z)$ .

# PLSA: Symmetric Parameterization



$$\begin{aligned} p(w_i, d_j) &= \sum_z p(w_i, d_j, z) \\ &= \sum_z p(w_i, d_j | z) p(z) \\ &= \sum_z p(w_i | z) p(d_j | z) p(z). \end{aligned}$$

# Model Fitting: EM Algorithm

- ▶ **Dyadic data  $\mathbf{X}$** : Entries  $X_{ij}$  are made for **dyads**  $(w_i, d_j)$  which refer to a domain with two sets of objects,  $\mathcal{W} = \{w_1, \dots, w_D\}$  and  $\mathcal{D} = \{d_1, \dots, d_N\}$ .
- ▶ **Complete-data likelihood**:

$$p(\mathbf{X}, z) = \prod_i \prod_j p(w_i, d_j, z)^{C_{ij}} = [p(w_i|z)p(d_j|z)p(z)]^{C_{ij}},$$

where  $C_{ij}$  are the empirical counts for dyads  $(w_i, d_j)$  and  $X_{ij} = C_{ij} / \sum_i \sum_j C_{ij}$ .

- ▶ **EM optimization**
  - ▶ E-step: Compute the expected complete-data log-likelihood:

$$\langle \mathcal{L}_c \rangle = \sum_i \sum_j \sum_k p(z_k | w_i, d_j) C_{ij} \log [p(w_i | z_k) p(d_j | z_k) p(z_k)].$$

- ▶ M-step: Re-estimate parameters  $p(w_i | z_k)$ ,  $p(d_j | z_k)$ ,  $p(z_k)$  which maximizes  $\langle \mathcal{L}_c \rangle$ .

## E-Step: Compute $p(z|w_i, d_j)$

Compute the posterior distribution over latent variables:

$$\begin{aligned} p(z_k|w_i, d_j) &= \frac{p(w_i, d_j|z_k)p(z_k)}{p(w_i, d_j)} \\ &= \frac{p(w_i|z_k)p(d_j|z_k)p(z_k)}{\sum_l p(w_i|z_l)p(d_j|z_l)p(z_l)}, \end{aligned}$$

where  $p(w_i|z_k)$ ,  $p(d_j|z_k)$ ,  $p(z_k)$  are estimated in the M-step.

## M-Step: Re-estimate Parameters

Re-estimate parameters:

$$\begin{aligned}p(w_i|z_k) &= \frac{\sum_j C_{ij} p(z_k|w_i, d_j)}{\sum_i \sum_j C_{ij} p(z_k|w_i, d_j)}, \\p(d_j|z_k) &= \frac{\sum_i C_{ij} p(z_k|w_i, d_j)}{\sum_i \sum_j C_{ij} p(z_k|w_i, d_j)}, \\p(z_k) &= \frac{\sum_i \sum_j C_{ij} p(z_k|w_i, d_j)}{\sum_i \sum_j C_{ij}}.\end{aligned}$$

Updating equations in the M-step are determined by solving

$$\frac{\partial}{\partial p(w_i|z_k)} \left[ \langle \mathcal{L}_c \rangle + \lambda(1 - \sum_i p(w_i|z_k)) \right] = 0,$$

$$\frac{\partial}{\partial p(d_j|z_k)} \left[ \langle \mathcal{L}_c \rangle + \lambda(1 - \sum_j p(d_j|z_k)) \right] = 0,$$

$$\frac{\partial}{\partial p(z_k)} \left[ \langle \mathcal{L}_c \rangle + \lambda(1 - \sum_l p(z_l)) \right] = 0,$$

for  $p(w_i|z_k)$ ,  $p(d_j|z_k)$ ,  $p(z_k)$ , respectively.



# Document Clustering by PLSA

- ▶ You are given parameters  $p(w_i|z_k)$ ,  $p(d_j|z_k)$ ,  $p(z_k)$  estimated by EM optimization.
- ▶ Compute

$$p(z_k|d_j) \propto p(d_j|z_k)p(z_k).$$

- ▶ Assign document  $d_j$  to cluster  $k^*$  if

$$k^* = \arg \max_k p(z_k|d_j).$$

## PLSA: Revisited

- ▶ The PLSA models each word in a document as a sample from a mixture model where the mixture components are multinomial random variables that can be viewed as representations of **topics**,

$$p(w, d) = \sum_z p(z)p(w|z)p(d|z) = p(d) \sum_z p(w|z)p(z|d).$$

- ▶ Each document is represented as a list of mixing proportions for mixture components and reduced to a probability distribution on a fixed set of topics. This distribution is the **reduced description** associated with the document.
- ▶ In PLSA, each document is represented as a list of numbers (the mixing proportions for topics) and there is no generative probabilistic model for these numbers. This leads to the following problems:
  - ▶ The number of parameters in the model grows linearly with the size of corpus, which leads to problems with overfitting.
  - ▶ It is not clear how to assign probability to a document outside of the training set.

## References

- ▶ S. Deerwester, S. T. Dumais, and R. Harshman, "Indexing by latent semantic analysis," *Journal of the American Society of Information Science*, vol. 41, no. 6, pp. 391-407, 1990.
- ▶ T. Hofmann, "Probabilistic latent semantic indexing," in *Proc. SIGIR-1999*.
- ▶ T. Hofmann, "Probabilistic latent semantic analysis," in *Proc. UAI-1999*.