

## Latent Topic Models

Lecturer: Xiaojin Zhu

jerryzhu@cs.wisc.edu

Often a document can be thought of as a mixture of a small number of topics (e.g., a news article can be about 60% “finance”, 30% “politics”, and 10% “war”.) Latent topic models recover hidden topics from document collections. This is interesting because

- The topics can be viewed as *latent semantic concepts*. Psychologists use latent topic models to explain the concept space we operate in, which is called the latent semantic space.
- Each document can be represented at the topics level instead of the word level. This allows two documents which share no common words (one says “buy”, the other says “sell”) to be regarded as similar because they share the the same “finance” topic.

## 1 Latent Semantic Indexing (LSI) and Principal Component Analysis (PCA)

### 1.1 Singular Value Decomposition (SVD) and Eigen Decomposition of a Matrix

Let  $M_{v \times d}$  be a rectangular (i.e. not necessarily square nor symmetric) matrix. The singular value decomposition of  $M$  is

$$M_{v \times d} = U_{v \times m} S_{m \times m} V_{d \times m}^T, \quad (1)$$

where  $m = \min(v, d)$ ;  $U$  has columns  $u_1, \dots, u_m$  orthogonal to each other, so does  $V$ ;  $S$  is diagonal, where the diagonal elements are called the *singular values*.

Let  $M_{v \times v}$  be a square symmetric matrix. The eigen decomposition of  $M$  is

$$M_{v \times v} = \Phi_{v \times v} \Lambda_{v \times v} \Phi_{v \times v}^T, \quad (2)$$

where the columns of  $\Phi$  are the eigenvectors  $\phi_1, \dots, \phi_v$ , and  $\Lambda$  is a diagonal matrix with the eigenvalues on the diagonal.

### 1.2 Latent Semantic Indexing (LSI)

Let the document collection be represented by a  $v \times d$  matrix  $X = [x_1 | \dots | x_d]$ , where the columns are document BOW vectors,  $v$  is the vocabulary size,  $d$  is the number of documents. LSI is SVD applied to  $X$ :

$$X_{v \times d} = U_{v \times m} S_{m \times m} V_{d \times m}^T. \quad (3)$$

Furthermore, one only keeps the top  $k \ll m$  singular values. That is, let  $\hat{U}_{v \times k}$  be the first  $k$  columns of  $U$ ,  $\hat{S}_{k \times k}$  be the first  $k \times k$  submatrix of  $S$ ,  $\hat{V}_{d \times k}$  be the first  $k$  columns of  $V$ . Then

$$\hat{U}_{v \times k} \hat{S}_{k \times k} \hat{V}_{d \times k}^T$$

is the best rank- $k$  approximation to  $X$  in the least square sense. The  $k$  columns of  $\hat{U}_{v \times k}$  defines the new rotated and reduced coordinate system. The  $d$  columns of  $\hat{S}_{k \times k} \hat{V}_{d \times k}^T$  are the new coordinates of each document after dimensionality reduction.

### 1.3 Principal Component Analysis (PCA)

Let  $X$  be the document collection matrix above. Let  $\bar{x} = 1/d \sum_{i=1}^d x_i$  be the mean document vector. Now *center* each document by subtracting the mean  $y_i = x_i - \bar{x}$ . Let  $Y = [y_1 | \dots | y_d]$  be the centered matrix. Let  $\Sigma = 1/d Y Y^\top$  be the covariance matrix. The eigen-decomposition of  $\Sigma$

$$\Sigma = \Phi \Lambda \Phi^\top \quad (4)$$

is known as the Principal Component Analysis of  $X$ . Furthermore, one only keeps the top  $k \ll v$  eigenvalues. That is, let  $\hat{\Phi}_{v \times k}$  be the first  $k$  columns of  $\Phi$ ,  $\hat{\Lambda}_{k \times k}$  be the first  $k \times k$  submatrix of  $\Lambda$ . Then

$$\hat{\Phi}_{v \times k} \hat{\Lambda}_{k \times k} \hat{\Phi}_{v \times k}^\top$$

is the best rank- $k$  approximation to  $\Sigma$  in the least square sense. We approximate each  $y_i$  by

$$\hat{y}_i = \sum_{j=1}^k \alpha_{ij} \phi_j, \quad (5)$$

where  $\alpha_{ij} = \phi_j^\top y_i$ . This is dimensionality reduction. It turns out that such a  $\hat{Y}$  is the best rank- $k$  approximation to  $Y$ . To approximate  $x_i$ , one can use  $\hat{y}_i + \bar{x}$ .

One major usage of PCA is in data visualization. Let  $k = 2$  or  $3$ . Each  $y_i$  (or  $x_i$ ) is approximated by  $\alpha_i$ , a  $k$ -dimension vector now suitable for plotting in 2D or 3D spaces.

Another major usage of PCA is to understand the few major axes ( $\phi_1, \dots, \phi_k$ ) of the dataset. For example, if  $X$  is a collection of aligned face images (each  $x$  vector is the image pixels in some scanning order), then  $\phi$ 's (each vector can be turned back into an image) tend to represent the major prototype faces, or major expressions. This is known as eigen-face.

Yet another usage of PCA is to extrapolate. Start with a new  $\alpha$  vector that is not in the dataset. One can generate its preimage  $\hat{x} = \sum_{j=1}^k \alpha_j \phi_j + \bar{x}$ . This can be used to create novel face images, for example.

A shortcoming of PCA is that it is an unsupervised learning algorithm. There is no guarantee that the  $\phi$ 's necessarily correspond to any particular goal (for classification).

### 1.4 LSI vs. PCA

LSI is equivalent to PCA, *only if  $X$  is already centered*. In this case,

$$X X^\top = d \Sigma \quad (6)$$

$$U S V^\top (U S V^\top)^\top = d \Phi \Lambda \Phi^\top \quad (7)$$

$$U S^2 U^\top = \Phi (d \Lambda) \Phi^\top, \quad (8)$$

That is, the LSI(SVD)  $U$  matrix is simply the eigenvector matrix in PCA, the singular values are the square root of eigenvalues in PCA (scaled by  $d$ ).

## 2 Probabilistic Latent Semantic Analysis (pLSA)

Recall Naive Bayes models a document collection by  $K$  topics (classes). Each topic is a multinomial over words, and each document is generated from a single topic:

$$p(w) = \sum_{k=1}^K p(z = k) p(w | z = k). \quad (9)$$

The parameters  $p(z = k), p(w | z = k)$  can be learned from labeled data (MLE or MAP), or with the EM algorithm if there is no class label (this corresponds to a mixture of multinomial model with  $K$  components). This is not a very flexible model, because it assumes ‘‘one document, one topic’’.

pLSA assumes that *each* document  $d$  (with word vector  $w$ ) is generated from all topics, with document-specific topic weights. The generative process of pLSA is the following. Given a fixed document collection with  $n$  documents, we represent it as a  $n \times V$  document-word matrix with entry  $c(d, w)$ , the count of word type  $w$  in document  $d$ . At each iteration, one picks a topic  $z = 1 \dots K \sim p(z)$ , then picks a document and a word type independent of each other, but both depends on the topic:  $d \sim p(d|z = k)$ ,  $w \sim p(w|z = k)$ . Generate (add one count of) word  $w$  to document  $d$ . Repeat until we generate the document-word matrix. Under this process, the probability of picking the cell  $(d, w)$  is

$$p(d, w) = \sum_{z=1}^K p(z)p(d|z)p(w|z). \quad (10)$$

The model parameters are  $\Theta = \{p(z), p(d|z), p(w|z)\}$ . We want to find the MLE to maximize the likelihood of the observed document-word matrix,

$$\max_{\Theta} \sum_{d=1}^n \sum_{w=1}^V c(d, w) \log p(d, w) \quad (11)$$

$$= \max_{\Theta} \sum_{d=1}^n \sum_{w=1}^V c(d, w) \log \left( \sum_{z=1}^K p(z)p(d|z)p(w|z) \right). \quad (12)$$

Note  $z$  is a hidden variable, and note the sum inside log. We can apply the EM algorithm:

$$\sum_{d,w} c(d, w) \log p(d, w) \quad (13)$$

$$= \sum_{d,w} c(d, w) \log \left( \sum_{z=1}^K p(z)p(d|z)p(w|z) \right) \quad (14)$$

$$= \sum_{d,w} c(d, w) \log \left( \sum_{z=1}^K p(z|d, w, \Theta^t) \frac{p(z)p(d|z)p(w|z)}{p(z|d, w, \Theta^t)} \right) \quad (15)$$

$$\geq \sum_{d,w} c(d, w) \sum_{z=1}^K p(z|d, w, \Theta^t) \left( \log \frac{p(z)p(d|z)p(w|z)}{p(z|d, w, \Theta^t)} \right) \quad (16)$$

Note Jensen's inequality involves  $p(z|d, w, \Theta^t)$ , which computes the probability of topics separately for each cell, under the current parameters  $\Theta^t$ . This is exactly the E-step. They can be computed as

$$p(z|d, w, \Theta^t) \propto p(z|\Theta^t)p(d|z)p(w|z, \Theta^t). \quad (17)$$

Maximizing (16) by setting the gradient to zero amounts to the M-step, which gives

$$p(z) \propto \sum_d \sum_w c(d, w) p(z|d, w, \Theta^t) \quad (18)$$

$$p(d|z) \propto \sum_w c(d, w) p(z|d, w, \Theta^t) \quad (19)$$

$$p(w|z) \propto \sum_d c(d, w) p(z|d, w, \Theta^t). \quad (20)$$

The E-step and M-step are repeated until convergence.

Once the model is trained, we can look at it in the following way:

- $p(w|z)$  are the topics. Each topic is defined by a word multinomial. Often people find that the topics seem to have distinct semantic meanings.
- From  $p(d|z)$  and  $p(z)$ , we can compute  $p(z|d) \propto p(d|z)p(z)$ .  $p(z|d)$  is the topic wights for document  $d$ .

One drawback of pLSA is that it is *transductive* in nature. That is, there is no easy way to handle a new document that is not already in the collection.

### 3 Latent Dirichlet Allocation (LDA)

LDA too assumes that each document is a mixture of multiple topics, and each document can have different topics weights. Unlike pLSA, LDA is a full generative model and readily generalizes to unseen documents. The LDA generative process is the following.

1. Sample  $K$  multinomial distributions (each of size  $V$ )  $\phi_{1:K}$  from a Dirichlet distribution  $Dir(\beta)$ , these are the topics. Note  $\beta$  is a parameter vector of length  $V$ .
2. For each document
  - (a) Sample a topic multinomial (of size  $K$ )  $\theta$  from a Dirichlet distribution  $Dir(\alpha)$ .
  - (b) For each word position
    - i. Sample topic index  $z \sim \theta$
    - ii. Sample a word from the topic  $w \sim \phi_z$

The observation is the document collection  $w_{1:n}$ . The parameters are  $\alpha$  and  $\beta$ . The other parameters  $z$ ,  $\phi$  and  $\theta$  are hidden variables that will be marginalized out.

The probability of a topic multinomial  $\phi$  drawn from  $Dir(\beta)$  is

$$p(\phi|\beta) = \frac{\Gamma(\sum_{i=1}^V \beta_i)}{\prod_{i=1}^V \Gamma(\beta_i)} \prod_{i=1}^V \phi_i^{\beta_i-1}. \quad (21)$$

The probability of drawing the  $K$  topic multinomials are

$$p(\phi_{1:K}|\beta) = \prod_{j=1}^K p(\phi_j|\beta). \quad (22)$$

Similarly,

$$p(\theta|\alpha) = \frac{\Gamma(\sum_{i=1}^K \alpha_i)}{\prod_{i=1}^K \Gamma(\alpha_i)} \prod_{i=1}^K \theta_i^{\alpha_i-1}. \quad (23)$$

Let  $p(z|\theta)$  and  $p(w|z, \phi) = p(w|\phi_z)$  be the corresponding multinomial probabilities. The probability of generating the words  $w_{1:N}$  for a single document, given  $\theta, \phi$  is

$$\prod_{n=1}^N \sum_{z_n=1}^K p(z_n|\theta) p(w_n|z_n, \phi), \quad (24)$$

where we marginalize out  $z$  for each word position. Putting things together, the probability of a single document, after marginalizing out all hidden variables, is

$$p(w|\alpha, \beta) = \int_{\phi_{1:K}} \int_{\theta} p(\phi_{1:K}|\beta) p(\theta|\alpha) \left( \prod_{n=1}^N \sum_{z_n=1}^K p(z_n|\theta) p(w_n|z_n, \phi) \right) d\theta d\phi_{1:K}. \quad (25)$$

Finally, the probability of a document collection  $w^{(1)}, \dots, w^{(M)}$  is

$$p(w^{(1)}, \dots, w^{(M)}|\alpha, \beta) = \prod_{d=1}^M p(w^d|\alpha, \beta). \quad (26)$$