

Review

* Matrix Factorization

How to approximate $X \approx VY$?

VQ: minimize $\|X - VY\|$ subject to $\begin{cases} Y_{\alpha n} \in \{0, 1\} \\ \sum_{\alpha} Y_{\alpha n} = 1 \end{cases}$

PCA: minimize $\|X - VY\|$ subject to $V^T V = I$

NMF: minimize $\|X - VY\|$ subject to $V_{i\alpha}, Y_{\alpha n} \geq 0$

* Exponential family PCA

Assume X_{ij} has exponential family distribution with natural parameter θ_{ij}

How to model X with $\theta = VY$?

Maximize log-likelihood:

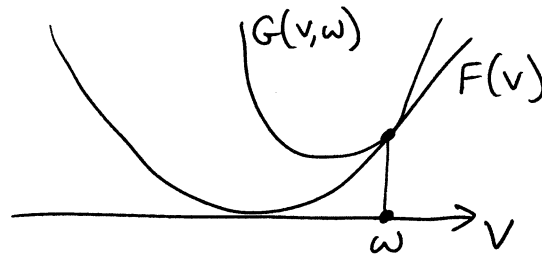
$$\begin{aligned} \mathcal{L} &= \sum_{ij} \log P(X_{ij}) \\ &= \sum_{ij} [\log P_0(X_{ij}) + X_{ij}(VY)_{ij} - G((VY)_{ij})] \end{aligned}$$

* Auxiliary functions

- Properties: $G(v, v) = F(v)$ and

$$G(v, w) \geq F(v) \text{ for all } v, w$$

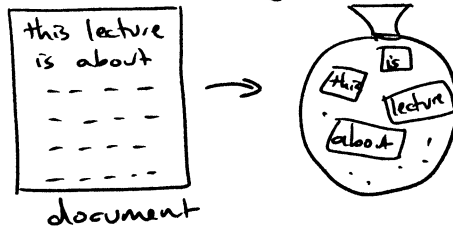
- Update: $v_{t+1} \leftarrow \operatorname{argmin}_v G(v, v_t)$



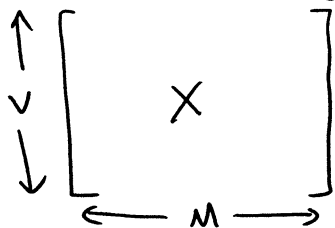
Document Modeling

* Bag-of-words model

- Represent document by vector of (raw or normalized) word counts.
- Ignore word ordering: words are "exchangeable"



- Corpus as term-by-document matrix



$V = \#$ words/tokens in vocabulary

$M = \#$ documents in corpus

$X_{ij} =$ raw or normalized count of word i in doc. j

$N_j = \sum_i X_{ij} = \#$ words in j^{th} document

- How to model X ?

* Matrix factorization: $X = VY$

- PCA known as "latent semantic indexing" in this context
- NMF can be applied to raw word counts
- But neither PCA/NMF provide probabilistic models of how documents are generated

Modeling Assumptions

* Notation / terminology:

Document = sequence of N words $\vec{w} = \{w_1, w_2, \dots, w_N\}$

Corpus = collection of M documents $D = \{\vec{w}_1, \vec{w}_2, \dots, \vec{w}_M\}$

* Bag-of-words / exchangeability assumption:

- summarize corpus by word-document counts X_{ij}

- data enters likelihood only via matrix X_{ij}

* What interesting models are there for $P(\vec{w})$ and $P(D)$?

(not n -grams, which model word ordering)

Early Work

* Unigram model

- draw words iid from single multinomial distribution
- document probability

$$\begin{aligned}P(\vec{w}) &= P(w_1, w_2, \dots, w_N) \\&= \prod_{n=1}^N P(w_n) \\&= \prod_{i=1}^V P(w=i)^{\text{count}(w=i)}\end{aligned}$$

- corpus probability

$$\begin{aligned}P(D) &= P(\vec{w}_1, \vec{w}_2, \dots, \vec{w}_M) \\&= \prod_{j=1}^M P(\vec{w}_j) \\&= \prod_{j=1}^M \prod_{i=1}^V P(w=i)^{x_{ij}} \leftarrow \text{word-document count}\end{aligned}$$

- Limitation.

No notion of "topic" to correlate words

* Mixture of Unigrams

- For each document:
 - Draw topic z from single multinomial distribution $P(z)$
 - Draw words iid from conditional multinomial distribution $P(w|z)$

- Document Probability

$$P(w_1, w_2, \dots, w_N) = \sum_{z=1}^K P(z) \prod_{n=1}^N P(w_n|z)$$

$K = \# \text{ topics (fixed)}$

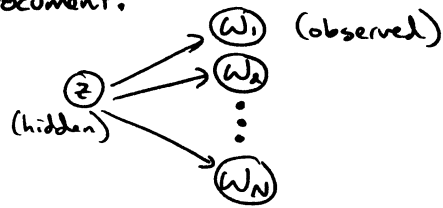
- Corpus Probability

$$\begin{aligned}P(D) &= \prod_{j=1}^M P(\vec{w}_j) \\&= \prod_{j=1}^M \left\{ \sum_z P(z) \prod_{i=1}^V P(w=i|z)^{x_{ij}} \right\}\end{aligned}$$

- Strength: words correlated by topics
- Weakness: one topic per document

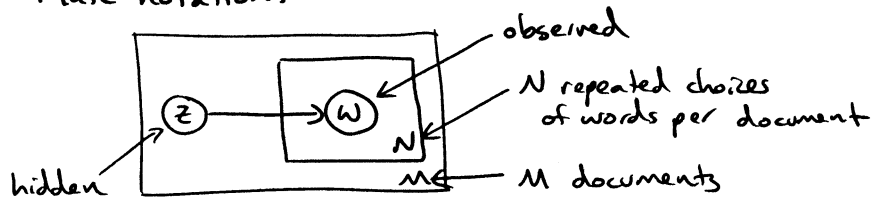
- Graphical model

Document:



Corpus: above is replicated M times

Plate notation:



* Probabilistic latent semantic indexing (PLSI)

- also known as an "aspect" model
- for each document d , and for each word in d :
 - draw a topic from multinomial distribution $P(z|d)$
 - draw a word from $P(w|z)$

- document probability

$$P(\vec{w}|d) = \prod_{n=1}^N \left[\sum_z P(z|d) P(w_n|z) \right]$$

- corpus probability

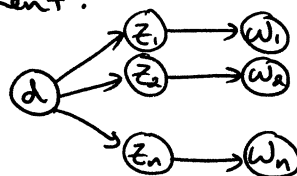
$$P(D) = \prod_{j=1}^M P(\vec{w}_j|d_j) = \prod_{j=1}^M \prod_{i=1}^V \left[\sum_z P(w=i|z) P(z|d_j) \right]^{x_{ij}}$$

log-likelihood of corpus:

$$\begin{aligned} \mathcal{L} &= \log P(D) \\ &= \sum_{ij} x_{ij} \log \left[\sum_z P(w=i|z) P(z|d_j) \right] \end{aligned}$$

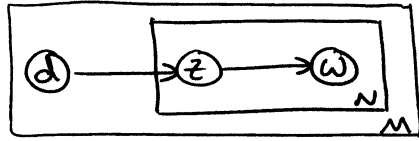
- Graphical model

document:



Corpus: replicate above M times

Plate notation:



M documents
 N repeated choices of
topic and word

- Strengths
 - multiple topics per document
 - EM algorithm can estimate $P(w|z)$ and $P(z|d_i)$ can discover "hidden" topics
 - Inference is tractable
- Weaknesses
 - Only assigns probabilities to documents in corpus. (doesn't generalize to new documents)
 - # parameters grows linearly in $M = \# \text{ docs}$
 \Rightarrow overfitting

Latent Dirichlet Allocation (LDA)

How to model distributions over topics?

* Dirichlet distribution

- a "distribution over distributions"

- let θ take values in simplex

$$\theta = (\theta_1, \theta_2, \dots, \theta_k) \text{ with } \theta_k \geq 0 \text{ and } \sum_{k=1}^k \theta_k = 1$$

- Exponential family distribution

$$P(\theta | \alpha) = \frac{\prod_{k=1}^k \alpha_k^{\theta_k}}{\prod_{k=1}^k \Gamma(\alpha_k)} \theta_1^{\alpha_1-1} \theta_2^{\alpha_2-1} \dots \theta_k^{\alpha_k-1}$$

parameters of Dirichlet distribution $\alpha_k > 0$ normalization constant

- Parameters

$$E[\theta_k] = \frac{\alpha_k}{\alpha_0} \text{ where } \alpha_0 = \sum_{k=1}^k \alpha_k$$

- Gamma function

$$\text{Definition! } \Gamma(x) = \int_0^{\infty} dt e^{-t} t^{x-1}$$

Properties: $\Gamma(x+1) = x \Gamma(x)$ prove using integration by parts

$$\Gamma(x+1) = x! \text{ for integers } x > 0$$

* Generative Model for LDA

- for each document \vec{w} in some corpus D :

1) Choose $N = \#$ words in document by some distribution (not modeled)

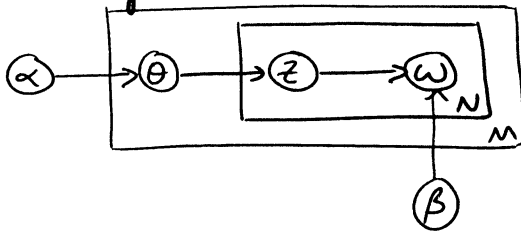
2) Choose θ from $P(\theta | \alpha)$

3) Do N times:

(a) choose $z_n \in \{1, 2, \dots, K\}$ from multinomial with weights θ

(b) choose w_n from multinomial $P(w_n | z_n, \beta)$ specifically:
 $P(w=j | z=K) = \beta_{Kj}$ matrix of parameters

- Graphical Model



- Modeling Parameters

α_k for $k=1 \dots K$ corpus/document - topic model

β_{kj} for $j=1 \dots V$ topic - word model

parameters = $K + KV$ no explicit linear dependence on M