## Intelligent Agents Topic Analysis: pLSI and LDA

## Ralf Möller Universität zu Lübeck Institut für Informationssysteme



**IM FOCUS DAS LEBEN** 

## Summary and Agenda

- IR Agents
  - Task/goal: Information retrieval
  - Agents visits document repositories and returns doc recommendations
  - Means:
    - Vector space (bag-of-words)
      - Dimension reduction (LSI)

Non-standard Databases and Data Mining

- Probability based retrieval (binary)
  - Language models
- Today: Language models with dimension reduction
  - Probabilistic Latent Semantic Indexing (pLSI)
  - Latent Dirichlet Allocation (LDA): Topic Models
- Soon: What agents can take with them
  - What agents can leave at the repository (win-win)



#### Acknowledgments

Ramesh M. Nallapati presentation on Generative Topic Models for Community Analysis & Sina Miran presentation on Probabilistic Latent Semantic Indexing (PLSI) & David M. Blei presentation on Probabilistic Topic Models



#### **Topic Models**

- Statistical methods that analyze the words of texts in order to:
  - Discover the themes that run through them (topics)
  - How those themes are connected to each other



## **Topic Modeling Scenario**



- Each topic is a distribution over words
- Each document is a mixture of corpus-wide topics
- Each word is drawn from one of those topics



## **Topic Modeling Scenario**



- In reality, we only observe the documents
- The other structures are hidden variables
- Topic modeling algorithms infer these variables from data



#### **Plate Notation**

INIVERSITÄT ZU LÜBECK INSTITUT FÜR INFORMATIONSSYSTEME

- Naïve Bayes Model: Compact representation
  - C = topic/class (name for a word distribution)
  - N = number of words in document
  - W<sub>i</sub> one specific word in corpus
  - M documents, W now words in documents







gene

genetic 0.01

0.04 0.02

0 01

0.02

#### Generative vs. Descriptive Models

- Generative models: Learn P(x, y)
  - Tasks:
    - Predict (infer) new data
    - Transform P(x,y) into P(y | x) for classification
  - Advantages
    - Assumptions and model are explicit
    - Use well-known algorithms
- Descriptive models: Learn P(y | x)
  - Task: Classify data
  - Advantages

/ERSITÄT ZU LÜBECK

- Fewer parameters to learn
- Better performance for classification

Input: Bayesian network

 $X = \{X_1, \dots, X_N\}, N- #nodes, T - # samples$ Output: T samples

Process nodes in topological order – first process the ancestors of a node, then the node itself:

- 1. For t = 0 to T
- $2. \quad For i = 0 to N$
- 3.  $X_i \leftarrow \text{sample } x_i^t \text{ from } P(x_i | pa_i)$

M. Henrion, "Propagating uncertainty in Bayesian networks by probabilistic logic sampling", Uncertainty in AI, pp. = 149-163, **1988** 



What does it mean to sample  $x_i^t$  from  $P(X_i | pa_i)$ ?

- Assume  $D(X_i) = \{0,1\}$
- Assume  $P(X_i | pa_i) = (0.3, 0.7)$



Draw a random number **r** from [0,1]
 If **r** falls in [0,0.3], set X<sub>i</sub> = 0
 If **r** falls in [0.3,1], set X<sub>i</sub>=1



#### Forward Sampling (Example)



Evidence :  $X_3 = 0$ 

// generate sample k 1. Sample  $x_1$  from  $P(x_1)$ 2. Sample  $x_2$  from  $P(x_2 | x_1)$ 3. Sample  $x_3$  from  $P(x_3 | x_1)$ 4. If  $x_3 \neq 0$ , reject sample and start from 1, otherwise 5. sample  $x_4$  from  $P(x_4 | x_2, x_3)$ 

Rejection sampling (rather inefficient)



#### Earlier Topic Models: Topics Known



fifth, an, of, futures, the, an, incorporated, a, a, the, inflation, most, dollars, quarter, in, is, mass

thrift, did, eighty, said, hard, 'm, july, bullish

that, or, limited, the

Automatically generated sentences from a unigram model



#### **Multinomial Naïve Bayes**





- How to specify Domain(C)?
  - Domain(C) =  $\{1, 2, ..., k\}$  or
  - Domain(C) =  $\{0, 1\}^k$
- How to specify  $P(c_d)$ ?
  - Define a table

	P(C)	
1	$p_1$	
К	$p_K$	

- or use parameterized distribution  $\pi = (p_1, ..., p_K)$ 

•  $P(C=c|\pi) = \prod_{k=1}^{K} \pi_k^{z_k}$ 

13

#### **Recap: Binomial Distribution**

- Describes the number of successes in a series of independent trials with two possible outcomes "success" or "no success"
- n = #trials
  - p = #successful trials / n
- Description of frequency of having exactly k successful trials as a function

$$\mathsf{B}_{\mathsf{p,n}}(\mathsf{k}) = \binom{n}{k} p^k (1-p)^{n-k}$$

- Es gilt:  $\sum_{i=0}^{n} B_{p,n}(i) = 1$
- If n=1: Bernoulli distribution



# $\binom{n}{k} = \frac{n!}{k!(n)}$

# Multinomial Distribution Mult(n | $\pi$ )

- Generalization of binomial distribution
  - K possible outcomes instead of 2
  - Probability mass function
    - n = number of trials
    - $x_j \in \{0, 1\}$  a count for how often class j occurs  $\sum_{i=1}^{n} \sum_{j=1}^{n} \sum_{i=1}^{n} \sum_{i=1}^{n} \sum_{i=1}^{n} \sum_{i=1}^{n} \sum_{i=1}^{n} \sum_{j=1}^{n} \sum_{i=1}^{n} \sum_{i=1}^{n$

$$\sum_{i=1}^k x_i = n$$

- $\mathbf{p}_{j} = \text{probability of class } j \text{ occurring}$  $Mult(x_{1}, \dots, x_{K}; p_{1}, \dots, p_{K}) = \frac{\Gamma(\sum_{i} x_{i} + 1)}{\prod_{i} \Gamma(x_{i} + 1)} \prod_{i=1}^{K} p_{i}^{x_{i}}$
- Here, the input to  $\Gamma(\cdot)$  is a positive integer, so  $\Gamma(n) = (n-1)!$
- If n=1: called categorial distribution ("multinoulli")
  - Often written  $Mult(.; p_1, ..., p_K)$  or  $Mult(. | p_1, ..., p_K)$
  - Generates a one-hot vector

# Sampling

- A variable value a can be sampled from a discrete distribution π = (p<sub>1</sub>, ..., p<sub>K</sub>)
- Notation: a ~ Mult(  $| \pi \rangle$
- Generate random number *x* from (0, 1]
- Find  $l \in \{1, 2, ..., k\}$  such that  $\sum_{i=1}^{l-1} p_i < x \le \sum_{i=1}^{l} p_i$

One-hot vector to be generated with position probability of indicator controlled by π

• Return  $(z_1, ..., z_K)$  such that  $z_l = 1$  and  $z_i = 0$  für  $i \neq l$ 



## **Multinomial with Matrices**

- Let  $\beta$  be a  $K \times V$  matrix (V vocabulary size), each row denotes a word distribution of a topic
- Select row k before applying multinomial:
  - Notation: Mult(.  $|\beta_k$ ) or Mult(.  $|\beta, k$ ) or Mult(.  $|k, \beta$ )





## Mixture of Unigrams: Known Topics



## Mixture of Unigrams: Unknown Topics



VERSITÄT ZU LÜBECK

TIONSSYSTEM

- Topics/classes are hidden
  - Joint probability of words and classes

$$\prod_{d=1}^{M} P(w_1, \dots, w_{N_d}, z_d \mid \beta, \pi) = \prod_{d=1}^{M} \pi_{z_d} \prod_{i=1}^{N_d} \beta_{z_k, w_i}$$

Sum over topics (K = number of topics)

$$\prod_{d=1}^{M} P(w_1, \dots, w_{N_d} | \beta, \pi) = \prod_{d=1}^{M} \sum_{k=1}^{K} \pi_{z_k} \prod_{i=1}^{N_d} \beta_{z_k, w_i}$$

Kamal Nigam, Andrew Kachites Mccallum, Sebastian Thrun & Tom Mitchell, Learning to Classify Text from Labeled and Unlabeled Documents, Proc. AAAI 98, Pages 792–799, **1998**.

Kamal Nigam, Andrew Kachites Mccallum, Sebastian Thrun & Tom Mitchell Text Classification from Labeled and Unlabeled Documents using EM Journal of Machine Learning volume 39, pages 103–134, **2000**.

$$\pi_{z_k} \coloneqq P(z_k | \pi)$$
  
$$\beta_{z_k, w_i} \coloneqq P(w_i | \beta, z_k)$$

• Learn parameters  $\pi$  and  $\beta$  $argmax_{\beta\pi}\prod_{d=1}^{M} P(w_1, ..., w_{N_d} | \beta, \pi)$ 

$$P(w_1, ..., w_{N_d} | \beta, \pi) = \sum_{k=1}^{K} \pi_{z_k} \prod_{i=1}^{N_d} \beta_{z_k, w_i}$$

Convex

• Use likelihood

$$\sum_{d=1}^{M} \log P(w_1, \dots, w_{N_d} | \beta, \pi) = \sum_{d=1}^{M} \log \sum_{k=1}^{K} \pi_{z_k} \prod_{i=1}^{N_d} \beta_{z_k, w_i}$$

Solve

IVERSITÄT ZU LÜBECK

INFORMATIONSSYSTEME

- $argmax_{\beta\pi}\sum_{d=1}^{M}\log\sum_{k=1}^{K}\pi_{z_{k}}\prod_{i=1}^{N_{d}}\beta_{z_{k},w_{i}}$
- Not a concave/convex function
- Note: a non-concave/non-convex function is not necessarily convex/concave
- Possibly no unique max, many saddle or turning points No easy way to find roots of derivative

Concave

#### Trick: Optimize Lower Bound









$$argmax_{\beta\pi}\sum_{d=1}^{M}\log\sum_{k=1}^{K}\pi_{z_{k}}\prod_{i=1}^{N_{d}}\beta_{z_{k},w_{i}}$$

- Optimize w.r.t. each document
- Derive lower bound

a, b

 $\log \sum_{i} \gamma_{i} x_{i} \geq \sum_{i} \gamma_{i} \log x_{i} \text{ where } \gamma_{i} \geq 0 \land \sum_{i} \gamma_{i} = 1 \qquad \text{Jensen's inequality} \\ \log(\mathbf{a} \cdot \mathbf{b}) \geq \mathbf{a} \cdot \log \mathbf{b}$ 

$$\log \sum_{i} x_{i} = \log \sum_{i} \gamma_{i} \frac{x_{i}}{\gamma_{i}} \ge \sum_{i} (\gamma_{i} \log x_{i} - \gamma_{i} \log \gamma_{i})$$
  
Entropy of  $\gamma_{di}$   
Sometimes  
called I(.)

$$\log \sum_{k=1}^{K} \pi_{z_{k}} \prod_{i=1}^{N_{d}} \beta_{z_{k},w_{i}} \geq \sum_{k=1}^{K} \left( \gamma_{k} \log(\pi_{z_{k}} \prod_{i=1}^{N_{d}} \beta_{z_{k},w_{i}}) \right) + H(\gamma)$$



IM FOCUS DAS LEBEN

Optimization problem for each document

 $argmax_{\beta\pi}\sum_{k=1}^{K}\left(\gamma_k\log(\pi_{z_k}\prod_{i=1}^{N_d}\beta_{z_k,w_i})\right)+H(\gamma)$ 

Convex? Concave?

- We have introduced a new latent variable γ to approximate the original functional to be optimized
- Each document is assumed to be associated with a latent variable  $\gamma \in [0,1]^{K}$ ,  $\Sigma_{k} \gamma_{k} = 1$ independent of other random variables
- Can be seen as a class in the new space  $\gamma_k, \pi_{z_k}, \beta_{z_k, w_i}$



• New optimization problem:

$$argmax_{\beta\pi}\sum_{k=1}^{K} \left(\gamma_k \log(\pi_{z_k} \prod_{i=1}^{N_d} \beta_{z_k, w_i})\right) + H(\gamma)$$

- Solution: Expectation Maximization
  - Iterative algorithm to find local optimum
  - Guess values of  $\gamma_k$ ,  $\pi_{z_k}$ ,  $\beta_{z_k,w_i}$
  - Compute  $\gamma_k = P(\gamma_k | \pi_{z_k}, \beta_{z_k, w_i})$  according to model
  - Use Maximum-likelihood estimation of to optimize  $\pi_{Z_k}$ ,  $\beta_{Z_k,W_i}$
  - Until no further improvement
- Guaranteed to maximize a lower bound on the loglikelihood of the observed data
- Use  $\pi_{z_k}$ ,  $\beta_{z_k,w_i}$  to estimate  $P(z_k|\pi)$ ,  $P(w_i|\beta, z_k)$ , respectively



#### Graphical Idea of the EM Algorithm





- EM solution
  - E step (compute  $\gamma_k = P(\gamma_k | \pi_{z_k}, \beta_{z_k, w_i})$ )

$$\gamma_{k}^{(t+1)} = \frac{\gamma_{k}^{(t)} \pi_{Z_{k}}^{(t)} \prod_{i=1}^{N_{d}} \beta_{Z_{k},w_{i}}^{(t)}}{\sum_{j=1}^{K} \gamma_{Z_{dj}}^{(t)} \pi_{Z_{j}}^{(t)} \prod_{i=1}^{N_{d}} \beta_{Z_{j},w_{i}}^{(t)}} \qquad \qquad \text{Independence}$$
assumption

M step (maximum likelihood optimization: use frequencies)

$$\pi_{Z_k}^{(t+1)} = \frac{\sum_{d=1}^{M} \gamma_{dk}^{(t)}}{M} \qquad \qquad \beta_{Z_k, w_i}^{(t+1)} = \frac{\sum_{d=1}^{M} \gamma_{dk}^{(t)} n(d, w_i)}{\sum_{d=1}^{M} \gamma_{dk}^{(t)} \sum_{j=1}^{N_d} n(d, w_j)}$$

*n*(*d*, *w*<sub>*i*</sub>) number of times word *w*<sub>*i*</sub> occurs in document *d* 



## Back to Topic Modeling Scenario

- Documents are associated with a single topic
- Words do not depend on context
  - Bag-of-words model





## Probabilistic LSI



- Select a document d with probability P(d)
- For each word of d in the training set
  - Choose a topic z with probability  $P(z \mid d)$
  - Generate a word with probability P(w | z)

$$P(d, w_i) = P(d) \sum_{k=1}^{K} P(w_i | z_k) P(z_k | d)$$

• Documents can have multiple topics



Thomas Hofmann, Probabilistic Latent Semantic Indexing, Proceedings of the 22<sup>nd</sup> Annual International <u>SIGIR</u> Conference on Research and Development in Information Retrieval (SIGIR-99), **1999** 

- Joint probability for all documents, words  $\prod_{i=1}^{M} \prod_{j=1}^{N_d} P(d, w_i)^{n(d, w_i)}$
- Distribution for document d, word w<sub>i</sub>

 $\overline{d}=1$   $\overline{i}=1$ 

$$P(d, w_i) = P(d) \sum_{k=1}^{K} P(w_i | z_k) P(z_k | d)$$



• Reformulate  $P(z_k|d)$  with Bayes' Rule

$$P(d, w_i) = \sum_{k=1}^{K} P(d|z_k) P(z_k) P(w_i|z_k)$$





# pLSI: Learning Using EM

Model

$$\prod_{d=1}^{M} \prod_{i=1}^{N_d} P(d, w_i)^{n(d, w_i)} \qquad P(d, w_i) = \sum_{k=1}^{K} P(d|z_k) P(z_k) P(w_i|z_k)$$

• Likelihood

$$L = \sum_{d=1}^{M} \sum_{i=1}^{N_d} n(d, w_i) \log P(d, w_i) = \sum_{d=1}^{M} \sum_{i=1}^{N_d} n(d, w_i) \log \sum_{k=1}^{K} P(d|z_k) P(z_k) P(w_i|z_k)$$

- Parameters to learn (M step)
  - $P(d|z_k)$   $P(z_k)$   $P(w_i|z_k)$

 $P(z_k|d, w_i)$ 

• (E step)





# pLSI: Learning Using EM

- EM solution
  - E step  $P(z_k|d, w_i) = \frac{P(z_k)P(d|z_k)P(w_i|z_k)}{\sum_{m=1}^{K} P(z_m)P(d|z_m)P(w_i|z_m)}$



– M step

$$P(w_{i}|z_{k}) = \frac{\sum_{d=1}^{M} n(d, w_{i}) P(z_{k}|d, w_{i})}{\sum_{d=1}^{M} \sum_{j=1}^{N_{d}} n(d, w_{j}) P(z_{k}|d, w_{j})}$$

$$P(d|z_{k}) = \frac{\sum_{i=1}^{N_{d}} n(d, w_{i}) P(z_{k}|d, w_{i})}{\sum_{d=1}^{M} \sum_{i=1}^{N_{d}} n(d, w_{i}) P(z_{k}|d, w_{i})}$$

$$P(z_{k}) = \frac{1}{R} \sum_{d=1}^{M} \sum_{i=1}^{N_{d}} n(d, w_{i}) P(z_{k}|d, w_{i}), R = \sum_{d=1}^{M} \sum_{i=1}^{N_{d}} n(d, w_{i}) P(z_{k}|d, w_{i}), R = \sum_{d=1}^{M} \sum_{i=1}^{N_{d}} n(d, w_{i}) P(z_{k}|d, w_{i}), R = \sum_{d=1}^{M} \sum_{i=1}^{N} n(d, w_{i}) P(z_{k}|d, w_{i}), R = \sum_{d=1}^{N} \sum_{i=1}^{N} n(d, w_{i}) P(z_{k}|d, w_{i}), R$$



## pLSI: Overview

- More realistic than mixture model
  - Documents can discuss multiple topics!
- Problems
  - Very many parameters
  - Danger of overfitting



# pLSI Testrun

- PLSI topics (TDT-1 corpus)
  - Approx. 7 million words, 15863 documents, K = 128

The two most probable topics that generate the term "flight" (left) and "love" (right).

List of most probable words per topic, with decreasing probability going down the list.

"plane"	"space shuttle"	"family"	"Hollywood"
$\mathbf{plane}$	space	home	film
airport	$\mathbf{shuttle}$	family	movie
$\operatorname{crash}$	mission	like	music
flight	astronauts	love	new
safety	launch	kids	$\mathbf{best}$
aircraft	station	$\operatorname{mother}$	hollywood
$\operatorname{air}$	crew	life	love
passenger	nasa	happy	actor
$\mathbf{board}$	$\mathbf{satellite}$	friends	entertainment
airline	$\operatorname{earth}$	$\operatorname{cnn}$	$\operatorname{star}$



#### Relation with LSI

$$P = U_k \Sigma_k V_k^T \qquad P(d, w_i) = \sum_{k=1}^K P(d|z_k) P(z_k) P(w_i|z_k)$$
$$U_k = \left( P(d|z_k) \right)_{d,k} \qquad \Sigma_k = \text{diag} \left( P(z_k) \right)_k \qquad V_k = \left( P(w_i|z_k) \right)_{i,k}$$



- Difference:
  - LSI: minimize Frobenius (L-2) norm
  - pLSI: log-likelihood of training data



## Intelligent Agents Topic Analysis: pLSI and LDA

## Ralf Möller Universität zu Lübeck Institut für Informationssysteme



**IM FOCUS DAS LEBEN** 

# pLSI with Multinomials




# Prior Distribution for Topic Mixture

- Goal: topic mixture proportions for each document could drawn from some distribution.
  - Distribution on multinomials (k-tuples of non-negative numbers that sum to one)
- The space of all of these multinomials can be interpreted geometrically as a (k-1)-*simplex*
  - K-1 independent values
  - Simplex = Generalization of a triangle to (k-1) dimensions
- Criteria for selecting our prior:
  - It needs to be defined for a (k-1)-simplex
  - Should have nice properties

ERSITÄT ZU LÜBECK STITUT FÜR INFORMATIONSSYSTEME





# LDA Model – Parameters



#### ← Proportions parameter

- (k-dimensional vector of real numbers)
- ← Per-document topic distribution (*k*-dimensional vector of probabilities summing up to 1)
- ← Per-word topic assignment (number from 1 to k)

#### ← Observed word

(number from 1 to v, where v is the number of words in the vocabulary)

← Word "prior" (v-dimensional)

### LDA Model





# Latent Dirichlet Allocation

- Document = mixture of topics (as in pLSI), but according to a Dirichlet prior
  - When we use a uniform Dirichlet prior, LDA= pLSI



# **Dirichlet Distributions**

$$p(\theta|\alpha) = \frac{\Gamma(\sum_{i} \alpha_{i})}{\prod_{i} \Gamma(\alpha_{i})} \prod_{i=1}^{K} \theta_{i}^{\alpha_{i}-1}$$

- Defined over a (k-1)-simplex
  - Takes K non-negative arguments which sum to one.
  - Consequently it is a natural distribution to use over multinomial distributions.



- The Dirichlet parameter  $\alpha_i$  can be thought of as a prior count of the  $i^{\rm th}$  class

$$Dir(x_{1}, ..., x_{K}; p_{1}, ..., p_{K}) = \frac{\Gamma(\sum_{i} x_{i} + 1)}{\prod_{i} \Gamma(x_{i} + 1)} \prod_{i=1}^{K} p_{i}^{x_{i}}$$



#### **Dirichlet Distribution over a 2-Simplex**



A panel illustrating probability density functions of a few Dirichlet distributions over a 2-simplex, for the following  $\alpha$  vectors (clockwise, starting from the upper left corner): (1.3, 1.3, 1.3), (3,3,3), (7,7,7), (2,6,11), (14, 9, 5), (6,2,6). [Wikipedia]



JNIVERSITÄT ZU LÜBECK INSTITUT FÜR INFORMATIONSSYSTEME

# LDA Model – Plate Notation



- For each document d,
  - Generate  $\theta_d \sim \text{Dirichlet}(\alpha)$
  - For each position  $i = 1, ..., N_d$ 
    - Generate a topic  $z_i \sim Mult(\cdot | \theta_d)$
    - Generate a word  $w_i \sim Mult(\cdot | z_{i'}\beta)$

$$P(\beta, \theta, z_1, \dots, z_{N_d}, w_1, \dots, w_{N_d})$$
  
=  $\prod_{d=1}^{M} P(\theta_d | \alpha) \prod_{i=1}^{N_d} P(z_i | \theta_d) P(w_i | \beta, z_i)$ 

# Corpus-level Parameter $\alpha$ (uniform: $\alpha_{i} = \alpha_{j}$ )

- Let  $\alpha = 1$
- Per-document topic distribution: K = 10, D = 15





# Corpus-level Parameter *α*

•  $\alpha = 10$ 

•  $\alpha = 100$ 





# Corpus-level Parameter $\alpha$

•  $\alpha = 0.1$ 

•  $\alpha = 0.01$ 





# Back to Topic Modeling Scenario

#### What are the words' topics and word distribs of topics?

-  $P(\beta, \theta, \mathbf{z} | \mathbf{w}, \alpha)$ 





# Topic-specific Words: "Smoothed" LDA Model



- Give a different word distribution to each topic
  - β is K×V matrix (V vocabulary size), each row denotes word distribution of a topic
- For each document d
  - Choose  $\theta_d \sim \text{Dirichlet}(\alpha)$
  - Choose  $\beta_k \sim \text{Dirichlet}(\eta)$
  - For each position  $i = 1, ..., N_d$ 
    - Generate a topic  $z_k \sim Mult(\cdot \mid \theta_d)$
    - Generate a word  $w_i \sim Mult(\cdot | z_{k'}\beta_{zk})$



# But why does LDA actually work?

- Trade-off between two goals
  - 1. For each document, allocate its words to as few topics as possible
  - 2. For each topic, assign high probability to as few terms as possible
- These goals are at odds.
  - Putting a document in a single topic makes #2 hard:
     All of its words must have non-negligible probability under that topic
  - Putting very few words in each topic makes #1 hard:
     To cover a document's words, it must assign many topics to it
- Trading off these goals finds groups of tightly co-occurring words



#### Query Answering Problem (non-smoothed version)



To which topics does a given document belong?

$$P(\theta, \mathbf{z} | \mathbf{w}, \alpha, \beta) = \frac{P(\theta, \mathbf{z}, \mathbf{w} | \alpha, \beta)}{P(\mathbf{w} | \alpha, \beta)}$$

$$P(\theta, \mathbf{z}, \mathbf{w} | \alpha, \beta) = P(\theta | \alpha) \prod_{i=1}^{N} P(z_i | \theta) P(w_i | z_i, \beta)$$

$$P(\mathbf{w} | \alpha, \beta) = \int \sum_{k=1}^{K} P(\mathbf{w}, \theta, \mathbf{z} | \alpha, \beta) \ d\theta = \int \sum_{k=1}^{K} P(\theta | \alpha) \prod_{i=1}^{N} P(z_i | \theta) P(w_i | z_i, \beta) \ d\theta = \int \sum_{k=1}^{K} P(\theta | \alpha) \prod_{i=1}^{N} \sum_{k=1}^{K} \prod_{j=1}^{V} (\theta_k \beta_{kj})^{w_i^j} \ d\theta$$

This not only looks awkward, but is as well *computationally intractable* in general. Coupling between  $\theta$  and  $\beta_{ij}$ . Solution: *Approximations*.

 $p(\theta|\alpha) = \frac{\Gamma(\sum_{i} \alpha_{i})}{\prod_{i} \Gamma(\alpha_{i})} \prod_{i=1}^{K} \theta_{i}^{\alpha_{i}-1}$ 



# LDA Learning

- Parameter learning:
  - Variational Inference / EM
    - Numerical approximation using lower-bounds
    - Results in biased solutions
    - Convergence has numerical guarantees
  - Gibbs Sampling
    - Stochastic simulation
    - Unbiased solutions
    - Stochastic convergence

We have a lecture on Approximation Algorithms for Probabilistic Models !



D. Blei, A. Ng, and M. Jordan. Latent Dirichlet allocation. Journal of Machine Learning Research, 3:993-1022, January **2003** 

# Back to Agents

- Agents not only use models
- Agents *build* models that are appropriate to fulfil the agents' goals ...
  - ... or maximize the utilities derived from preference structures and goals
- Agents *derive approximation algorithms* for query answering on the models they find appropriate





# LDA Application: Reuters Data

- Setup
  - 100-topic LDA trained on a 16,000 documents corpus of news articles by Reuters
  - Some standard stop words removed
- Top-7 words from some of the P(w|z)

"Arts"	"Budgets"	"Children"	"Education"
new	million	children	school
film	tax	women	students
show	program	people	schools
music	budget	child	education
movie	billion	years	teachers
play	federal	families	high
musical	year	work	public



# LDA Application: Reuters Data

Result

Again: "Arts", "Budgets", "Children", "Education".

The William Randolph Hearst Foundation will give \$1.25 million to Lincoln Center, Metropolitan Opera Co., New York Philharmonic and Juilliard School. "Our board felt that we had a real opportunity to make a mark on the future of the performing arts with these grants an act every bit as important as our traditional areas of support in health, medical research, education and the social services," Hearst Foundation President Randolph A. Hearst said Monday in announcing the grants.



# Measuring Performance

- Perplexity of a probability model
- Describe how well a probability distribution or probability model predicts a sample
  - q: Model of an unknown probability distribution p
     based on a training sample drawn from p
  - Evaluate q by asking how well it predicts a separate test sample  $x_1, \dots, x_N$  also drawn from p
  - Perplexity of q w.r.t. sample  $x_1, \ldots, x_N$  defined as

 $2^{-\frac{1}{N}\sum_{i=1}^{N}\log_2 q(x_i)}$ 

- A better model q will tend to assign higher probabilities to  $q(x_i)$ 
  - $\rightarrow$  lower perplexity ("less surprised by sample")



# Perplexity of Various Models





# Use of LDA

- A widely used topic model (Griffiths, Steyvers, 04)
- Complexity is an issue
- Use in IR:
  - Ad hoc retrieval (Wei and Croft, SIGIR 06: TREC benchmarks)
  - Improvements over traditional LM (e.g., LSI techniques)
  - But no consensus on whether there is any improvement over a relevance model, i.e., model with relevance feedback (relevance feedback part of the TREC tests)

T. Griffiths, M. Steyvers, Finding Scientific Topics. Proceedings of the National Academy of Sciences, 101 (suppl. 1), 5228-5235. **2004** 

Xing Wei and W. Bruce Croft. LDA-based document models for ad-hoc retrieval. In *Proceedings of the 29th annual international ACM SIGIR conference on Research and development in information retrieval* (SIGIR '06). ACM, New York, NY, USA, 178-185. **2006**.



TREC=Text REtrieval Conference

# Topic modelling

# Tomoharu Iwata



**IM FOCUS DAS LEBEN** 

# Social annotation services

- Delicious, Flickr, CiteULike, youtube, Last.fm, Technorati, Hatena
- Users can attach annotations freely to objects, and share the annotations.

delicious       assisted         social bookmarking       assisted         All your stuff in one place.       orgonic         Get to your bookmarks from any computer, anywhere.       orgonic         Image: State of the state	Jain Noor Sign In More	flickr ∏ <sup>™</sup> <sup>™</sup> ∏		Create Your Account Chy Likes a memory with your Yahed ID Share your photos.
Search the biggest collection of bookmarks in the universe	NOE MIRO 🔘	د By xillio 2,641 uploads in <u>the last</u> i	minude • 160,129 things tagged with morning • 2.5	SEARCH million things gestagged this month - Take the tour
See more Pegular kookmarks  Watch TV Online - Full Episodes save Watch TV Online - Full Episodes save Y tree online - Full Episodes save Het webdesign inventis Web Architects save	Popular Tags delign blog video software book	S. S		
Webdagen design and/or grid prit	nusic programming webdesign reference tutorial art	And George as a contact? Share & stay in touch	Upload & organize	
DevHub • Free Web Hosted Publishing Platform save vetolesign hosting development business publishing	howlo javascript tree linux web2.0 devolutionent	Crop, fix, edit	Explore	



# Derive content-unrelated annotations

- manufacturer of camera that took the photo – 'nikon', 'canon'
- when they were taken
  - '2008', 'november'
- remind the annotator
  - 'toread'
- qualities
  - 'great', '\*\*\*\*'
- ownership



29

# Proposed model

- generative model for contents (words) and annotations with relevance based on topic models
- infer relevance to the content for each annotation



# Latent Dirichlet allocation



[Blei et. al. Latent Dirichlet Allocation, JMLR2003]



**IM FOCUS DAS LEBEN** 

# **Correspondence LDA**



David M. Blei and Michael I. Jordan. Modeling annotated data. In Proceedings of the 26th annual international ACM SIGIR conference on Research and development in information retrieval (SIGIR '03). Association for Computing Machinery, New York, NY, USA, 127–134. **2003**.

**IM FOCUS DAS LEBEN** 

# Proposed model (Inference with Gibbs Sampling)



- N: #words, M: #annotations, D: #documents, K: #topics
- each annotation is associated with a latent variable r, r=1 if content-related, r=0 otherwise

INIVERSITÄT ZU LÜBECK INSTITUT FÜR INFORMATIONSSYSTEME

# **Topics in Delicious**

	unrelated	Topic1	Topic2	Topic3	Topic4	Topic5
	reference	money	video	opensource	food	windows
<b>^</b> \	web	finance	music	software	recipes	linux
Ч Ш	imported	economics	videos	programming	recipe	sysadmin
	design	business	fun	development	cooking	Windows
ō	internet	economy	entertainment	linux	Food	security
$\overline{\mathbf{v}}$	online	Finance	funny	tools	Recipes	computer
Ŧ.	cool	financial	movies	rails	baking	microsoft
<u>0</u>	toread	investing	media	ruby	health	network
	tools	bailout	Video	webdev	vegetarian	Linux
	blog	finances	film	rubyonrails	diy	ubuntu
		money	music	project	recipe	windows
Ŋ		financial	video	code	food	system
2		credit	link	server	recipes	microsoft
Ì		market	tv	ruby	make	linux
		economic	movie	rails	wine	software
		october	itunes	source	made	file
5		economy	film	file	add	server
6		banks	amazon	version	love	user
ř		government	play	files	eat	files
		bank	interview	development	good	ubuntu



# **Topics in Flickr**

	unrelated	Topic1	Topic2	Topic3	Topic4	Topic5
	2008	dance	sea	autumn	rock	beach
ດງ	nikon	bar	sunset	trees	house	travel
Ĩ	canon	dc	sky	tree	party	vacation
Б	white	digital	clouds	mountain	park	camping
ō	yellow	concert	mountains	fall	inn	landscape
t	red	bands	ocean	garden	coach	texas
at a	photo	music	panorama	bortescristian	creature	lake
	italy	washingtondc	south	geotagged	halloween	cameraphone
Y	california	dancing	ireland	mud	mallory	md
	color	work	oregon	natura	night	sun
prob						
able ir			6			
nage						



IM FOCUS DAS LEBEN

# Perplexity



The proposed method performed better than Corr-LDA <sub>37</sub> in the case of noisy social annotation data.



**IM FOCUS DAS LEBEN** 

# Generative Topic Models for Community Analysis

#### Ramesh Nallapati

http://www.cs.cmu.edu/~wcohen/10-802/lda-sep-18.ppt

&

Arthur Asuncion, Qiang Liu, Padhraic Smyth:

Statistical Approaches to Joint Modeling of Text and Network Data



# What if the corpus has network structure?



CORA citation network. Figure from [Chang, Blei, AISTATS 2009]



J. Chang, and D. Blei. Relational Topic Models for Document Networks. AISTATS, volume 5 of JMLR Proceedings, page 81-88. JMLR.org, **2009**.

# Outline



- Citation Modeling
  - with pLSI
  - with LDA
- Modeling influence of citations
- Relational Topic Models



# Hyperlink Modeling Using pLSI



- Select document d ~ Mult(.  $|\pi)$ 
  - For each position  $n = 1, ..., N_d$ 
    - Generate  $z_n \sim Mult(\cdot | \theta_d)$
    - Generate  $w_n \sim Mult(\cdot | \beta_{z_n})$
  - For each citation  $j = 1, ..., L_d$ 
    - Generate  $z_j \sim Mult(\cdot | \theta_d)$
    - Generate  $c_j \sim Mult(\cdot | \gamma_{Z_j})$

D. A. Cohn, Th. Hofmann, The Missing Link - A Probabilistic Model of Document Content and Hypertext Connectivity, In: Proc. NIPS, pp. 430-436, **2000** 

IM FOCUS DAS LEBEN 71

# Hyperlink Modeling Using pLSI





• pLSI likelihood

$$= \prod_{d=1}^{M} P(w_1, \dots, w_{N_d}, d | \theta, \beta, \pi)$$
$$= \prod_{d=1}^{M} \pi_d \left( \prod_{i=1}^{N_d} \sum_{k=1}^{K} \theta_{dk} \beta_{kw_n} \right)$$

New likelihood

$$\prod_{d=1}^{M} P(w_1, \dots, w_{N_d}, c_1, \dots, c_{L_d}, d | \theta, \beta, \gamma, \pi)$$

$$= \prod_{d=1}^{M} \pi_d \left( \prod_{i=1}^{N_d} \sum_{k=1}^{K} \theta_{dk} \beta_{kw_n} \right) \left( \prod_{j=1}^{L_d} \sum_{k=1}^{K} \theta_{dk} \gamma_{kc_j} \right)$$

• Learning using EM
## Hyperlink Modeling Using pLSI

- Heuristic
  - 0 <  $\alpha$  < 1 determines the relative importance of content and hyperlinks

$$\prod_{d=1}^{M} P(w_1, \dots, w_{N_d}, c_1, \dots, c_{L_d}, d | \theta, \beta, \gamma, \pi)$$
$$= \prod_{d=1}^{M} \pi_d \left( \prod_{i=1}^{N_d} \sum_{k=1}^{K} \theta_{dk} \beta_{kw_n} \right)^{\alpha} \left( \prod_{j=1}^{L_d} \sum_{k=1}^{K} \theta_{dk} \gamma_{kc_j} \right)^{1-\alpha}$$



## Hyperlink modeling using PLSA

- Experiments: Text Classification
- Datasets:
  - Web KB
    - 6000 CS dept web pages with hyperlinks
    - 6 Classes: faculty, course, student, staff, etc.
  - Cora
    - 2000 Machine learning abstracts with citations
    - 7 classes: sub-areas of machine learning
- Methodology:
  - Learn the model on complete data and obtain  $\theta_d$  for each document
  - Test documents classified into the label of the nearest neighbor in training set
  - Distance measured as cosine similarity in the  $\theta$  space
  - Measure the performance as a function of  $\alpha$



#### **Overview on Evaluation Measures**

		True condition				
	Total population	Condition positive	Condition negative	$= \frac{\sum \text{ Condition positive}}{\sum \text{ Total population}}$		
Predicted condition	Predicted condition positive	True positive	False positive (Type I error)	Positive predictive value (PPV), Precision $= \frac{\Sigma \text{ True positive}}{\Sigma \text{ Test outcome positive}}$	False discovery rate (FDR) = $\frac{\Sigma \text{ False positive}}{\Sigma \text{ Test outcome positive}}$	
	Predicted condition negative	False negative (Type II error)	True negative	False omission rate (FOR) = $\frac{\Sigma \text{ False negative}}{\Sigma \text{ Test outcome negative}}$	$\frac{\text{Negative predictive value}}{(\text{NPV})} = \frac{\Sigma \text{ True negative}}{\Sigma \text{ Test outcome negative}}$	
	$\frac{\text{Accuracy (ACC)} =}{\sum \text{True positive} + \sum \text{True negative}}{\sum \text{Total population}}$	True positive rate (TPR), Sensitivity, Recall = $\frac{\Sigma \text{ True positive}}{\Sigma \text{ Condition positive}}$	False positive rate (FPR), Fall-out = $\frac{\Sigma \text{ False positive}}{\Sigma \text{ Condition negative}}$	Positive likelihood ratio (LR+) = $\frac{TPR}{FPR}$	Diagnostic odds ratio (DOR)	
		False negative rate (FNR), Miss rate = $\frac{\Sigma \text{ False negative}}{\Sigma \text{ Condition positive}}$	True negative rate (TNR), Specificity (SPC) = $\frac{\Sigma \text{ True negative}}{\Sigma \text{ Condition negative}}$	Negative likelihood ratio (LR–) = $\frac{FNR}{TNR}$	$=\frac{LR^{+}}{LR^{-}}$	



#### Hyperlink Modeling Using pLSI





# Hyperlink modeling using LDA



- For each document d,
  - Generate  $\theta_d \sim \text{Dirichlet}(\alpha)$
  - For each position  $i = 1, ..., N_d$ 
    - Generate a topic  $z_i \sim Mult(\cdot | \theta_d)$
    - Generate a word  $w_i \sim Mult (\cdot | \beta_{Z_n})$
  - For each citation  $j = 1, ..., L_c$ 
    - Generate  $z_i \sim Mult(\theta_d)$
    - Generate  $c_i \sim Mult (\cdot | \gamma_{Z_j})$
- Learning using variational EM, Gibbs sampling

E. Erosheva, S Fienberg, J. Lafferty, Mixed-membership models of scientific publications. Proc National Academy Science U S A. 2004 Apr 6;101 Suppl 1:5220-7. Epub **2004** Mar 12.

#### Link-pLSI-LDA: Topic Influence in Blogs





R. Nallapati, A. Ahmed, E. Xing, W.W. Cohen, Joint Latent Topic Models for Text and Citations, In: Proc. KDD, **2008**.

## Modeling Citation Influences - Copycat Model

 Each topic in a citing document is drawn from one of the topic mixtures of cited publications





L. Dietz, St. Bickel, and T. Scheffer, Unsupervised Prediction of Citation Influences, In: Proc. ICML **2007**.

### **Modeling Citation Influences**

 Citation influence model: Combination of LDA and Copycat model





L. Dietz, St. Bickel, and T. Scheffer, Unsupervised Prediction of Citation Influences, In: Proc. ICML **2007**.

#### **Modeling Citation Influences**

• Citation influence graph for LDA paper





## **Modeling Citation Influences**

• Words in LDA paper assigned to citations

Cited Title	Associated Words		
Probabilistic	text(0.04), latent(0.04),	0.49	
Latent Semantic	modeling(0.02), model(0.02),		
Indexing	indexing(0.01), $semantic(0.01)$ ,		
	document(0.01), collections(0.01)		
Modelling	dirichlet(0.02), mixture(0.02),	0.25	
heterogeneity	allocation(0.01), context(0.01),		
with and	variable(0.0135), bayes(0.01),		
without the	continuous(0.01), improves(0.01),		
Dirichlet process	model(0.01), $proportions(0.01)$		
Introduction to	variational(0.01), inference(0.01),	0.22	
Variational	algorithms(0.01), including(0.01),		
Methods for	each(0.01), we(0.01), via(0.01)		
Graphical			
Methods			



## Relational Topic Model (RTM) [ChangBlei 2009]

 Same setup as LDA, except now we have observed network information across documents



 $y_{d,d'} \sim \psi(y_{d,d'} | z_d, z_{d'}, \eta)$ 

"Link probability function"

Documents with similar topics are more likely to be linked.

J. Chang, and D. Blei. Relational Topic Models for Document Networks. AISTATS, volume 5 of JMLR Proceedings, page 81-88. JMLR.org, **2009**.

## Relational Topic Model (RTM) [ChangBlei 2009]



- For each document d
  - Draw topic proportions  $\theta_d | \alpha \sim Dir(\alpha)$
  - For each word  $w_{d,n}$ 
    - Draw assignment  $z_{d,n} | \theta_d \sim Mult(\theta_d)$
    - Draw word  $w_{d,n}|_{Z_{d,n}}, \beta_{1:K} \sim Mult(\beta_{Z_{d,n}})$
  - For each pair of documents *d*, *d*'
    - Draw binary link indicator  $y|z_d, z_{d'} \sim \psi(\cdot | z_d, z_{d'}, \eta)$



	# Docs	# Links	Ave. Doc- Length	Vocab-Size	Link Semantics
CORA	4,000	17,000	1,200	60,000	Paper citation (undirected)
Netflix Movies	10,000	43,000	640	38,000	Common actor/director
Enron (Undirected)	1,000	16,000	7,000	55,000	Communication between person i and person j
Enron (Directed)	2,000	21,000	3,500	55,000	Email from person i to person j



#### Conclusion

- Topic Modeling
- Relational topic modeling provides a useful start for combining text and network data in a single statistical framework

• RTM can improve over simpler approaches for link prediction

- Opportunities for future work:
  - Faster algorithms for larger data sets
  - Better understanding of non-edge modeling
  - Extended models

