

Introduction to Probabilistic Topic Models

Ko, Youngjoong

**Computer Engineering
Dong-A University**

Intelligent System Laboratory, Dong-A University

Contents

- ❖ Introduction
- ❖ Generative Models
- ❖ Probabilistic Topic Models
- ❖ Algorithm for Extracting Topics
- ❖ Polysemy with Topics
- ❖ Computing Similarities



Introduction

❖ Why topic model?

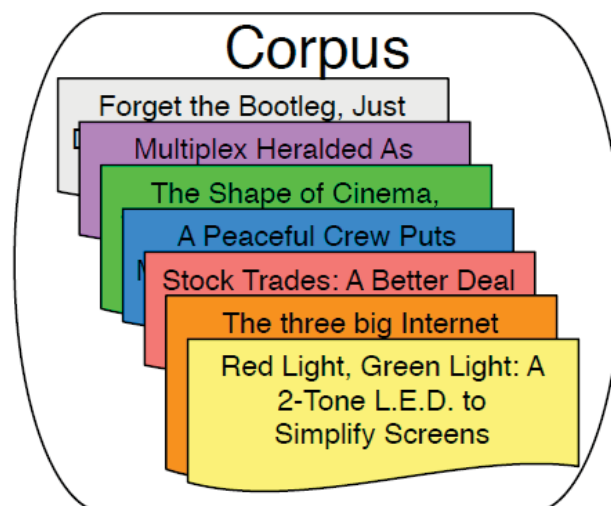
- Suppose you have a huge number of documents
- You want to know what's going on
- Don't have time to read them (e.g. every New York Times article from the 50's)
- Topic models offer a way to get a corpus-level view of major themes
- Unsupervised



Introduction

❖ Conceptual Approach

From an **input corpus** → words to topics



Introduction

❖ Conceptual Approach

From an input corpus → **words to topics**

TOPIC 1

computer,
technology,
system,
service, site,
phone,
internet,
machine

TOPIC 2

sell, sale,
store, product,
business,
advertising,
market,
consumer

TOPIC 3

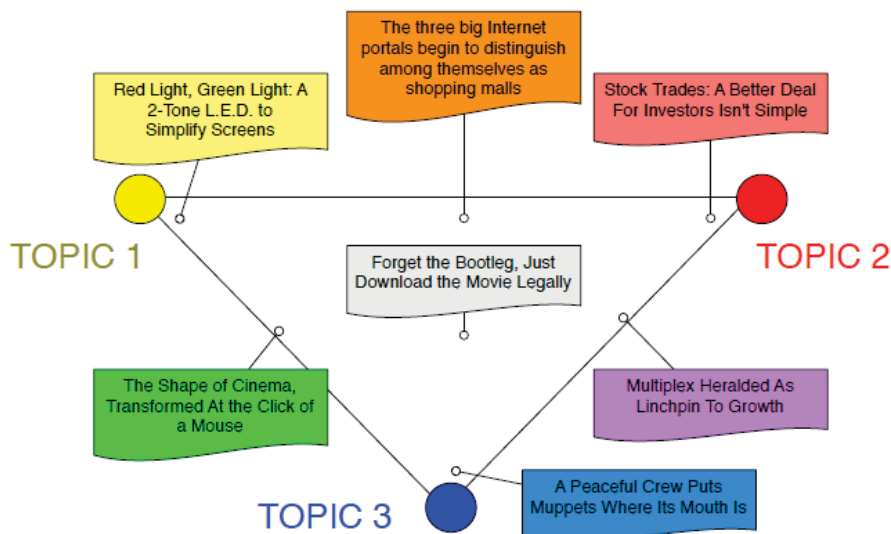
play, film,
movie, theater,
production,
star, director,
stage



Introduction

❖ Conceptual Approach

➤ For each document, what topics are expressed by that document?



Introduction

❖ Topics from Science

human	evolution	disease	computer
genome	evolutionary	host	models
dna	species	bacteria	information
genetic	organisms	diseases	data
genes	life	resistance	computers
sequence	origin	bacterial	system
gene	biology	new	network
molecular	groups	strains	systems
sequencing	phylogenetic	control	model
map	living	infectious	parallel
information	diversity	malaria	methods
genetics	group	parasite	networks
mapping	new	parasites	software
project	two	united	new
sequences	common	tuberculosis	simulations



Introduction

❖ Latent Semantic Analysis (LSA)

- People illustrate that applying a statistical method such as LSA to large databases can yield insight into human cognition
- Three claims
 - semantic information can be derived from a word-document co-occurrence matrix
 - dimensionality reduction is an essential part of this derivation
 - words and documents can be represented as points in Euclidean space
- A different approach
 - that is consistent with the first two of these claims,
 - but differs in the third, describing a class of statistical models in which the semantic properties of words and documents are expressed in terms of probabilistic topics.



Introduction

❖ Topic models

- Based upon the idea that documents are mixtures of topics, where a topic is a probability distribution over words.
- A topic model is a generative model for documents
 - It specifies a simple probabilistic procedure by which documents can be generated
 - To make a new document, one chooses a distribution over topics.
 - Then, for each word in that document, one chooses a topic at random according to this distribution, and draws a word from that topic.
- Standard statistical techniques can be used to invert this process
 - Four example topics from the TASA corpus, a collection of over 37,000 text passages
 - The sixteen words that have the highest probability under each topic.
 - Four topics relate to **drug use**, **colors**, **memory and the mind**, and **doctor visits**.
 - Documents with different content can be generated by choosing different distributions over topics.
 - Ex)
 - ✓ by giving equal probability to the first two topics, one could construct a document about a person that has taken too many drugs, and how that affected color perception.



Introduction

❖ An illustration of four (out of 300) topics

Topic 247	Topic 5	Topic 43	Topic 56																																																																																																																																								
<table border="1"><thead><tr><th>word</th><th>prob.</th></tr></thead><tbody><tr><td>DRUGS</td><td>.069</td></tr><tr><td>DRUG</td><td>.060</td></tr><tr><td>MEDICINE</td><td>.027</td></tr><tr><td>EFFECTS</td><td>.026</td></tr><tr><td>BODY</td><td>.023</td></tr><tr><td>MEDICINES</td><td>.019</td></tr><tr><td>PAIN</td><td>.016</td></tr><tr><td>PERSON</td><td>.016</td></tr><tr><td>MARIJUANA</td><td>.014</td></tr><tr><td>LABEL</td><td>.012</td></tr><tr><td>ALCOHOL</td><td>.012</td></tr><tr><td>DANGEROUS</td><td>.011</td></tr><tr><td>ABUSE</td><td>.009</td></tr><tr><td>EFFECT</td><td>.009</td></tr><tr><td>KNOWN</td><td>.008</td></tr><tr><td>PILLS</td><td>.008</td></tr></tbody></table>	word	prob.	DRUGS	.069	DRUG	.060	MEDICINE	.027	EFFECTS	.026	BODY	.023	MEDICINES	.019	PAIN	.016	PERSON	.016	MARIJUANA	.014	LABEL	.012	ALCOHOL	.012	DANGEROUS	.011	ABUSE	.009	EFFECT	.009	KNOWN	.008	PILLS	.008	<table border="1"><thead><tr><th>word</th><th>prob.</th></tr></thead><tbody><tr><td>RED</td><td>.202</td></tr><tr><td>BLUE</td><td>.099</td></tr><tr><td>GREEN</td><td>.096</td></tr><tr><td>YELLOW</td><td>.073</td></tr><tr><td>WHITE</td><td>.048</td></tr><tr><td>COLOR</td><td>.048</td></tr><tr><td>BRIGHT</td><td>.030</td></tr><tr><td>COLORS</td><td>.029</td></tr><tr><td>ORANGE</td><td>.027</td></tr><tr><td>BROWN</td><td>.027</td></tr><tr><td>PINK</td><td>.017</td></tr><tr><td>LOOK</td><td>.017</td></tr><tr><td>BLACK</td><td>.016</td></tr><tr><td>PURPLE</td><td>.015</td></tr><tr><td>CROSS</td><td>.011</td></tr><tr><td>COLORED</td><td>.009</td></tr></tbody></table>	word	prob.	RED	.202	BLUE	.099	GREEN	.096	YELLOW	.073	WHITE	.048	COLOR	.048	BRIGHT	.030	COLORS	.029	ORANGE	.027	BROWN	.027	PINK	.017	LOOK	.017	BLACK	.016	PURPLE	.015	CROSS	.011	COLORED	.009	<table border="1"><thead><tr><th>word</th><th>prob.</th></tr></thead><tbody><tr><td>MIND</td><td>.081</td></tr><tr><td>THOUGHT</td><td>.066</td></tr><tr><td>REMEMBER</td><td>.064</td></tr><tr><td>MEMORY</td><td>.037</td></tr><tr><td>THINKING</td><td>.030</td></tr><tr><td>PROFESSOR</td><td>.028</td></tr><tr><td>FELT</td><td>.025</td></tr><tr><td>REMEMBERED</td><td>.022</td></tr><tr><td>THOUGHTS</td><td>.020</td></tr><tr><td>FORGOTTEN</td><td>.020</td></tr><tr><td>MOMENT</td><td>.020</td></tr><tr><td>THINK</td><td>.019</td></tr><tr><td>THING</td><td>.016</td></tr><tr><td>WONDER</td><td>.014</td></tr><tr><td>FORGET</td><td>.012</td></tr><tr><td>RECALL</td><td>.012</td></tr></tbody></table>	word	prob.	MIND	.081	THOUGHT	.066	REMEMBER	.064	MEMORY	.037	THINKING	.030	PROFESSOR	.028	FELT	.025	REMEMBERED	.022	THOUGHTS	.020	FORGOTTEN	.020	MOMENT	.020	THINK	.019	THING	.016	WONDER	.014	FORGET	.012	RECALL	.012	<table border="1"><thead><tr><th>word</th><th>prob.</th></tr></thead><tbody><tr><td>DOCTOR</td><td>.074</td></tr><tr><td>DR</td><td>.063</td></tr><tr><td>PATIENT</td><td>.061</td></tr><tr><td>HOSPITAL</td><td>.049</td></tr><tr><td>CARE</td><td>.046</td></tr><tr><td>MEDICAL</td><td>.042</td></tr><tr><td>NURSE</td><td>.031</td></tr><tr><td>PATIENTS</td><td>.029</td></tr><tr><td>DOCTORS</td><td>.028</td></tr><tr><td>HEALTH</td><td>.025</td></tr><tr><td>MEDICINE</td><td>.017</td></tr><tr><td>NURSING</td><td>.017</td></tr><tr><td>DENTAL</td><td>.015</td></tr><tr><td>NURSES</td><td>.013</td></tr><tr><td>PHYSICIAN</td><td>.012</td></tr><tr><td>HOSPITALS</td><td>.011</td></tr></tbody></table>	word	prob.	DOCTOR	.074	DR	.063	PATIENT	.061	HOSPITAL	.049	CARE	.046	MEDICAL	.042	NURSE	.031	PATIENTS	.029	DOCTORS	.028	HEALTH	.025	MEDICINE	.017	NURSING	.017	DENTAL	.015	NURSES	.013	PHYSICIAN	.012	HOSPITALS	.011
word	prob.																																																																																																																																										
DRUGS	.069																																																																																																																																										
DRUG	.060																																																																																																																																										
MEDICINE	.027																																																																																																																																										
EFFECTS	.026																																																																																																																																										
BODY	.023																																																																																																																																										
MEDICINES	.019																																																																																																																																										
PAIN	.016																																																																																																																																										
PERSON	.016																																																																																																																																										
MARIJUANA	.014																																																																																																																																										
LABEL	.012																																																																																																																																										
ALCOHOL	.012																																																																																																																																										
DANGEROUS	.011																																																																																																																																										
ABUSE	.009																																																																																																																																										
EFFECT	.009																																																																																																																																										
KNOWN	.008																																																																																																																																										
PILLS	.008																																																																																																																																										
word	prob.																																																																																																																																										
RED	.202																																																																																																																																										
BLUE	.099																																																																																																																																										
GREEN	.096																																																																																																																																										
YELLOW	.073																																																																																																																																										
WHITE	.048																																																																																																																																										
COLOR	.048																																																																																																																																										
BRIGHT	.030																																																																																																																																										
COLORS	.029																																																																																																																																										
ORANGE	.027																																																																																																																																										
BROWN	.027																																																																																																																																										
PINK	.017																																																																																																																																										
LOOK	.017																																																																																																																																										
BLACK	.016																																																																																																																																										
PURPLE	.015																																																																																																																																										
CROSS	.011																																																																																																																																										
COLORED	.009																																																																																																																																										
word	prob.																																																																																																																																										
MIND	.081																																																																																																																																										
THOUGHT	.066																																																																																																																																										
REMEMBER	.064																																																																																																																																										
MEMORY	.037																																																																																																																																										
THINKING	.030																																																																																																																																										
PROFESSOR	.028																																																																																																																																										
FELT	.025																																																																																																																																										
REMEMBERED	.022																																																																																																																																										
THOUGHTS	.020																																																																																																																																										
FORGOTTEN	.020																																																																																																																																										
MOMENT	.020																																																																																																																																										
THINK	.019																																																																																																																																										
THING	.016																																																																																																																																										
WONDER	.014																																																																																																																																										
FORGET	.012																																																																																																																																										
RECALL	.012																																																																																																																																										
word	prob.																																																																																																																																										
DOCTOR	.074																																																																																																																																										
DR	.063																																																																																																																																										
PATIENT	.061																																																																																																																																										
HOSPITAL	.049																																																																																																																																										
CARE	.046																																																																																																																																										
MEDICAL	.042																																																																																																																																										
NURSE	.031																																																																																																																																										
PATIENTS	.029																																																																																																																																										
DOCTORS	.028																																																																																																																																										
HEALTH	.025																																																																																																																																										
MEDICINE	.017																																																																																																																																										
NURSING	.017																																																																																																																																										
DENTAL	.015																																																																																																																																										
NURSES	.013																																																																																																																																										
PHYSICIAN	.012																																																																																																																																										
HOSPITALS	.011																																																																																																																																										

- Representing the content of words and documents with probabilistic topics has one distinct advantage over spatial representation
 - Each topic is individually interpretable,
 - ✓ providing a probability distribution over words that picks out a coherent cluster of correlated terms.
 - ✓ the topics are typically as interpretable as the ones shown here.
 - This contrasts with the arbitrary axes of a spatial representation, and can be extremely useful in many applications



Introduction

❖ Why should you care?

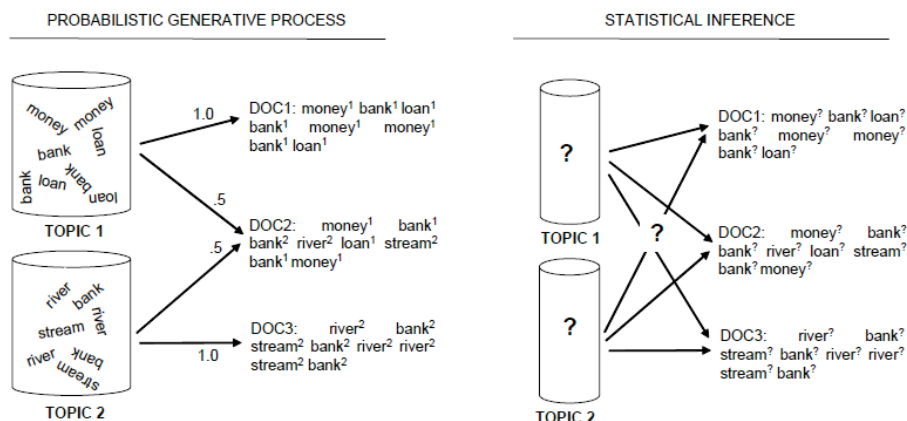
- Neat way to explore/understand corpus collections
- NLP Applications
 - POS Tagging [Toutanova and Johnson 2008]
 - Word Sense Disambiguation [Boyd-Graber et al. 2007]
 - Word Sense Induction [Brody and Lapata 2009]
 - Discourse Segmentation [Purver et al. 2006]
- Psychology [Griths et al. 2007b]: word meaning, polysemy
- Inference is (relatively) simple



Generative Models

❖ A Generative Model for Documents

- Based on simple probabilistic sampling rules that describe how words in documents might be generated on the basis of latent (random) variables
 - The goal is to find the best set of latent variables that can explain the observed data (i.e., observed words in documents), assuming that the model actually generated the data.
 - Ex) Illustration of the topic modeling approach in two distinct ways:



Generative Models

❖ A Generative Model for Documents

➤ Ex) Illustration of the topic modeling approach in two distinct ways:

▪ In the generative model (Left panel)

- ✓ With two topics (money and rivers)
- ✓ Bags containing different distributions over words
- ✓ Different documents can be produced by picking words from a topic depending on the weight given to the topic.

- ✓ topic models to capture polysemy (eg. Bank)
- ✓ there is no notion of mutual exclusivity that restricts words to be part of one topic only.

- ✓ Bag-of-words assumption : common to many statistical models of language with LSA

▪ In the statistical inference (Right Panel)

- ✓ Given the observed words in a set of documents, we would like to know what topic model is most likely to have generated the data.
- ✓ This involves inferring the probability distribution over words associated with each topic, the distribution over topics for each document, and, often, the topic responsible for generating each word



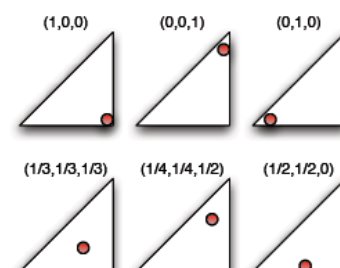
Generative Models

❖ A Generative Model for Documents

- How your data came to be
- Sequence of Probabilistic Steps
- Posterior Inference

❖ Multinomial Distribution

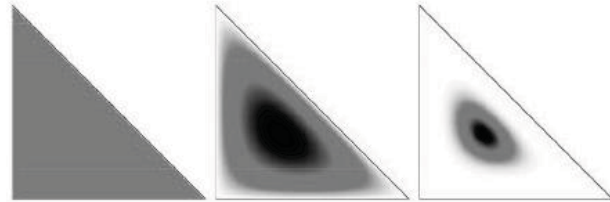
- Distribution over discrete outcomes
- Represented by non-negative vector that sums to one
- Picture representation
- Come from a Dirichlet distribution



Generative Models

❖ Dirichlet Distribution

$$P(\mathbf{p} | \alpha \mathbf{m}) = \frac{\Gamma(\sum_k \alpha m_k)}{\prod_k \Gamma(\alpha m_k)} \prod_k p_k^{\alpha m_k - 1}$$



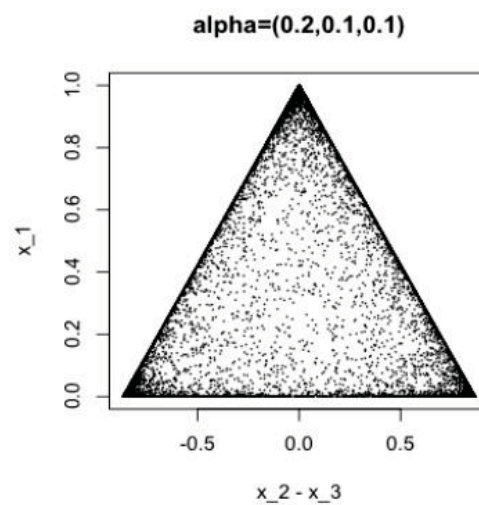
$\alpha = 3, \mathbf{m} = (\frac{1}{3}, \frac{1}{3}, \frac{1}{3})$ $\alpha = 6, \mathbf{m} = (\frac{1}{3}, \frac{1}{3}, \frac{1}{3})$ $\alpha = 30, \mathbf{m} = (\frac{1}{3}, \frac{1}{3}, \frac{1}{3})$

$\alpha = 14, \mathbf{m} = (\frac{1}{7}, \frac{5}{7}, \frac{1}{7})$ $\alpha = 14, \mathbf{m} = (\frac{1}{7}, \frac{1}{7}, \frac{5}{7})$ $\alpha = 2.7, \mathbf{m} = (\frac{1}{3}, \frac{1}{3}, \frac{1}{3})$



Generative Models

❖ Dirichlet Distribution



Generative Models

❖ Dirichlet Distribution

- If $\phi \sim \text{Dir}(\alpha)$, $\mathbf{w} \sim \text{Mult}(\phi)$, and $n_k = |\{w_i : w_i = k\}|$ then

$$p(\phi|\alpha, \mathbf{w}) \propto p(\mathbf{w}|\phi)p(\phi|\alpha) \quad (1)$$

$$\propto \prod_k \phi^{n_k} \prod_k \phi^{\alpha_k-1} \quad (2)$$

$$\propto \prod_k \phi^{\alpha_k+n_k-1} \quad (3)$$

- Conjugacy: this **posterior** has the same form as the **prior**



Probabilistic Topic Models

❖ The Fundamental Idea

- A document is a mixture of topics
- To introduce notation
 - $P(z)$: the distribution over topics z in a particular document
 - $P(w|z)$: the probability distribution over words w given topic z

➤ Generative Process

- Each word w_i in a document (where the index refers to the i -th word token) is generated by first sampling a topic from the topic distribution,
- then choosing a word from the topic-word distribution

$$P(w_i) = \sum_{j=1}^T P(w_i | z_i = j) P(z_i = j)$$

- ✓ $P(z_i = j)$: the probability that the j -th topic was sampled for the i -th word token
- ✓ $P(w_i | z_i = j)$ as the probability of word w_i under topic j .
- ✓ T : the number of topics.
- To simplify notation,
 - ✓ $\phi(j) = P(w | z=j)$: the multinomial distribution over words for topic j
 - ✓ $\theta(d) = P(z)$: the multinomial distribution over topics for document d



Probabilistic Topic Models

❖ The Fundamental Idea

➤ Generative Process

- the text collection : D documents
- each document d : N_d word tokens
- N : the total number of word tokens (i.e., $N = \sum N_d$).
- The parameters ϕ and θ indicate which words are important for which topic and which topics are important for a particular document, respectively.

➤ From probabilistic Latent Semantic Indexing method (pLSI)

- Hofmann (1999; 2001) introduced the probabilistic topic approach to document modeling in his Probabilistic Latent Semantic Indexing method
- The pLSI model does not make any assumptions about how the mixture weights θ are generated, making it difficult to test the generalizability of the model to new documents.
- Blei et al. (2003) extended this model by introducing a Dirichlet prior on θ , calling the resulting generative model Latent Dirichlet Allocation (LDA)
- As a conjugate prior for the multinomial, the Dirichlet distribution is a convenient choice as prior, simplifying the problem of statistical inference



Probabilistic Topic Models

- The probability density of a T dimensional Dirichlet distribution over the multinomial distribution $p=(p_1, \dots, p_T)$

$$\text{Dir}(\alpha_1, \dots, \alpha_T) = \frac{\Gamma(\sum_j \alpha_j)}{\prod_j \Gamma(\alpha_j)} \prod_{j=1}^T p_j^{\alpha_j - 1}$$

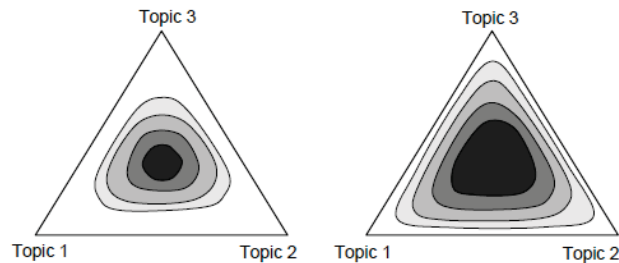
- parameters : $\alpha_1 \dots \alpha_T$
- each hyperparameter α_j can be interpreted as a prior observation count for the number of times topic j is sampled in a document, before having observed any actual words from that document.



Probabilistic Topic Models

➤ Dirichlet distribution over the multinomial distribution $p=(p_1, \dots, p_T)$

- Ex) Dirichlet distribution for three topics in a two-dimensional simplex (Left: $\alpha = 4$. Right: $\alpha = 2$)
 - ✓ a smoothed topic distribution, with the amount of smoothing determined by the α parameter



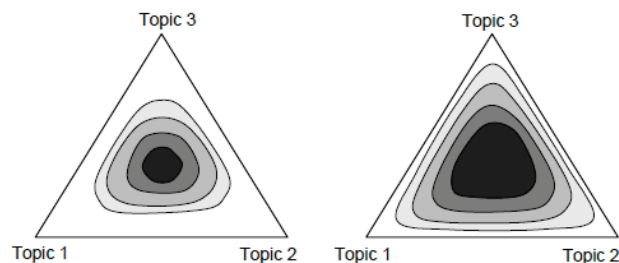
- Dirichlet prior on the topic distributions can be interpreted as forces on the topic combinations with higher α moving the topics away from the corners of the simplex, leading to more smoothing (compare the left and right panel).
- For $\alpha < 1$, the modes of the Dirichlet distribution are located at the corners of the simplex.
- In this regime (often used in practice), there is a bias towards sparsity, and the pressure is to pick topic distributions favoring just a few topics.



Probabilistic Topic Models

➤ Dirichlet distribution over the multinomial distribution $p=(p_1, \dots, p_T)$

- Ex) Dirichlet distribution for three topics in a two-dimensional simplex (Left: $\alpha = 4$. Right: $\alpha = 2$)
 - ✓ a smoothed topic distribution, with the amount of smoothing determined by the α parameter



- Dirichlet prior on the topic distributions can be interpreted as forces on the topic combinations with higher α moving the topics away from the corners of the simplex, leading to more smoothing (compare the left and right panel).
- For $\alpha < 1$, the modes of the Dirichlet distribution are located at the corners of the simplex.
- In this regime (often used in practice), there is a bias towards sparsity, and the pressure is to pick topic distributions favoring just a few topics.



Probabilistic Topic Models

➤ The hyperparameter β

- placing a symmetric Dirichlet(β) prior on ϕ
- interpreted as the prior observation count on the number of times words are sampled from a topic before any word from the corpus is observed
- This smooths the word distribution in every topic, with the amount of smoothing determined by β

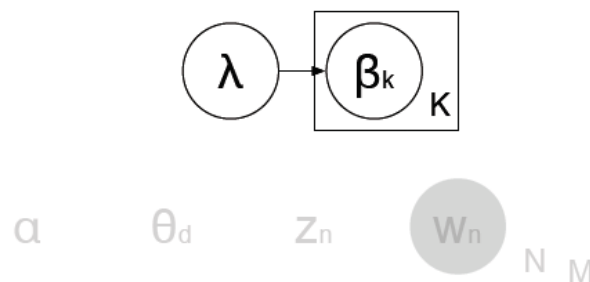
➤ Good choices for the hyperparameters α and β

- depend on number of topics and vocabulary size.
- From previous research, we have found $\alpha = 50/T$ and $\beta = 0.01$ to work well with many different text collections.



Probabilistic Topic Models

❖ Graphical Model for Generative Model Approach

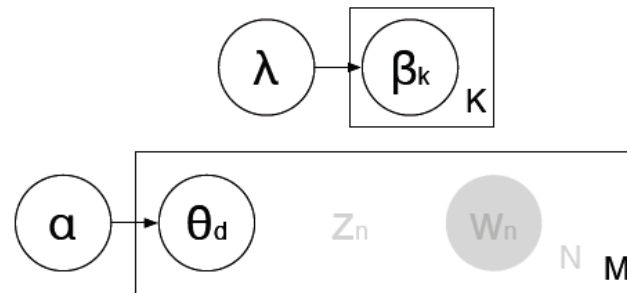


- For each topic $k \in \{1, \dots, K\}$, draw a multinomial distribution β_k from a Dirichlet distribution with parameter λ



Probabilistic Topic Models

❖ Graphical Model for Generative Model Approach

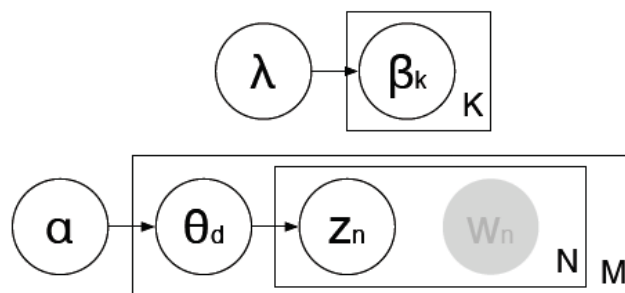


- For each topic $k \in \{1, \dots, K\}$, draw a multinomial distribution β_k from a Dirichlet distribution with parameter λ
- For each document $d \in \{1, \dots, M\}$, draw a multinomial distribution θ_d from a Dirichlet distribution with parameter α



Probabilistic Topic Models

❖ Graphical Model for Generative Model Approach

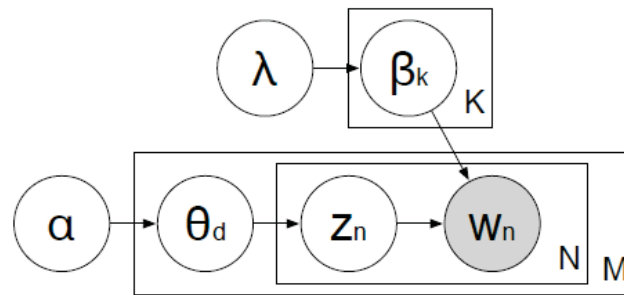


- For each topic $k \in \{1, \dots, K\}$, draw a multinomial distribution β_k from a Dirichlet distribution with parameter λ
- For each document $d \in \{1, \dots, M\}$, draw a multinomial distribution θ_d from a Dirichlet distribution with parameter α
- For each word position $n \in \{1, \dots, N\}$, select a hidden topic z_n from the multinomial distribution parameterized by θ .



Probabilistic Topic Models

❖ Graphical Model for Generative Model Approach

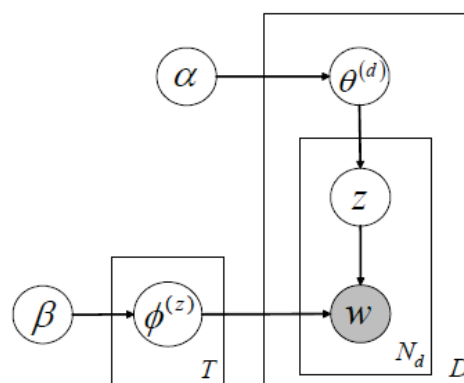


- For each topic $k \in \{1, \dots, K\}$, draw a multinomial distribution β_k from a Dirichlet distribution with parameter λ
- For each document $d \in \{1, \dots, M\}$, draw a multinomial distribution θ_d from a Dirichlet distribution with parameter α
- For each word position $n \in \{1, \dots, N\}$, select a hidden topic z_n from the multinomial distribution parameterized by θ .
- Choose the observed word w_n from the distribution β_{z_n} .



Probabilistic Topic Models

❖ Graphical Model



- the inner plate over z and w illustrates the repeated sampling of topics and words until N_d words have been generated for document d .
- The plate surrounding $\theta^{(d)}$ illustrates the sampling of a distribution over topics for each document d for a total of D documents.
- The plate surrounding $\phi^{(z)}$ illustrates the repeated sampling of word distributions for each topic z until T topics have been generated



Probabilistic Topic Models

❖ Graphical Model

➤ Plate Notation

- Probabilistic generative models with repeated sampling steps can be conveniently illustrated
- shaded variables : observed variable, unshaded variables : latent (i.e., unobserved) variables
- The variables ϕ and θ , as well as z (the assignment of word tokens to topics) are the three sets of latent variables that we would like to infer.
- To treat the hyperparameters α and β as constants in the model
- Arrows indicate conditional dependencies between variables
- plates (the boxes in the figure) refer to repetitions of sampling steps with the variable in the lower right corner referring to the number of samples.

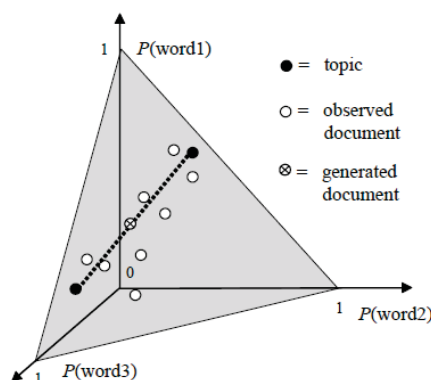


Probabilistic Topic Models

❖ Geometric Interpretation

➤ The probabilistic topic model has an elegant geometric interpretation

- With a vocabulary of W distinct words, a W dimensional space can be constructed where each axis represents the probability of observing a particular word type.
- The $W-1$ dimensional simplex represents all probability distributions over words
- the shaded region is the two-dimensional simplex that represents all probability distributions over three words
- the topics span a low-dimensional subsimplex and the projection of each document onto the low-dimensional subsimplex can be thought of as dimensionality reduction.

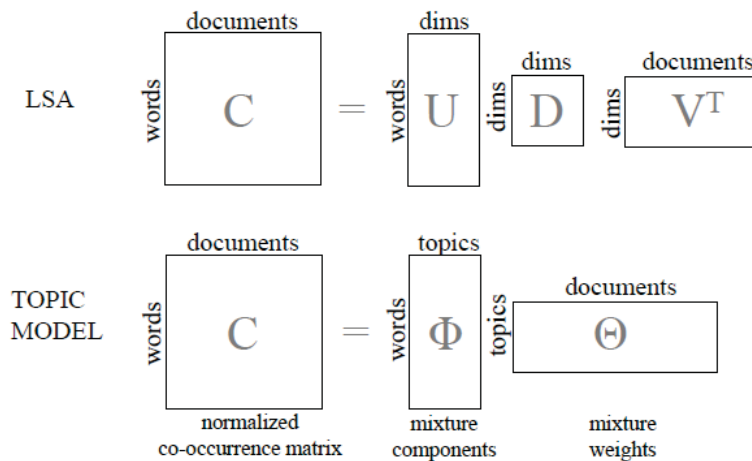


Probabilistic Topic Models

❖ Matrix Factorization Interpretation

➤ LSA vs. Topic Modeling

- In LSA, a word document co-occurrence matrix can be decomposed into three matrices: a matrix of word vectors, a diagonal matrix with singular values and a matrix with document vectors.
- In the topic model, the word-document co-occurrence matrix is split into two parts: a topic matrix Φ and a document matrix.



Probabilistic Topic Models

❖ Matrix Factorization Interpretation

➤ LSA vs. Topic Modeling

- To find a **low-dimensional representation** for the content of a set of documents
- In topic models, the word and document vectors of the two decomposed matrices are probability distributions with the accompanying constraint that the feature values are non-negative and sum up to one.
- In the LDA model, additional a priori constraints are placed on the word and topic distributions



Probabilistic Topic Models

❖ Topic Models: What's Important

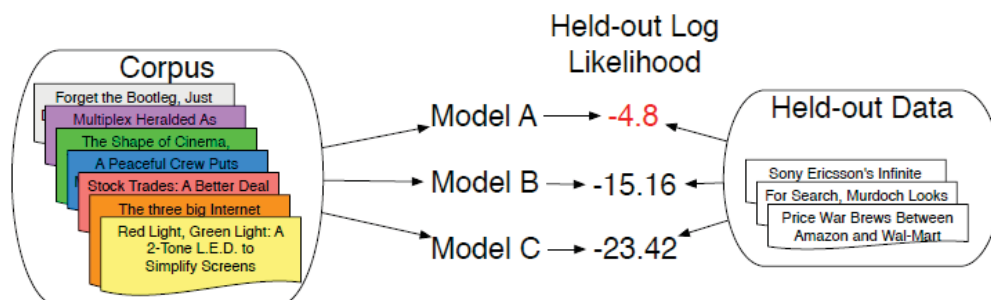
- Topic Models
 - Topics to words : multinomial distribution
 - Documents to topics : multinomial distribution
- Statistical structure inferred from data
- Have semantic coherence because of language use
- We use latent Dirichlet allocation (LDA) [Blei et al. 2003], a fully Bayesian version of pLSI [Hofmann 1999], probabilistic version of LSA [Landauer and Dumais 1997]



Probabilistic Topic Models

❖ Evaluation

➤ Likelihood



Measures predictive power, not what the topics are

$$P(\mathbf{w} | \mathbf{w}', \mathbf{z}', \alpha \mathbf{m}, \beta \mathbf{u}) = \sum_{\mathbf{z}} P(\mathbf{w}, \mathbf{z} | \mathbf{w}', \mathbf{z}', \alpha \mathbf{m}, \beta \mathbf{u})$$

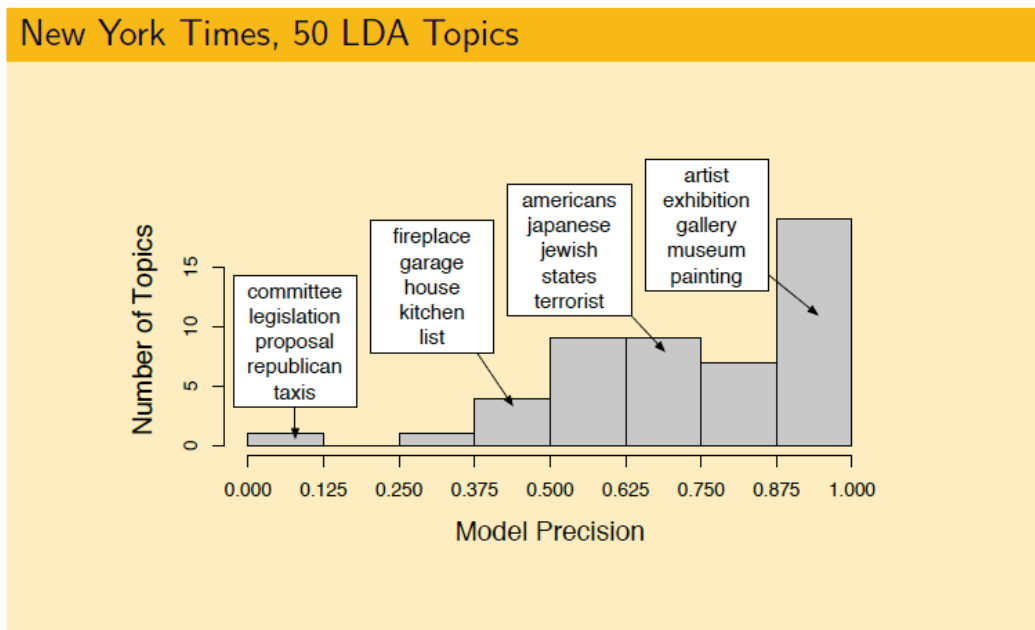
How you compute it is important too [?]



Probabilistic Topic Models

❖ Evaluation

➤ Which Topics are Interpretable?



Algorithm for Extracting Topics

❖ Gibbs Sampling

➤ Inference

computer,
technology,
system,
service, site,
phone,
internet,
machine

sell, sale,
store, product,
business,
advertising,
market,
consumer

play, film,
movie, theater,
production,
star, director,
stage

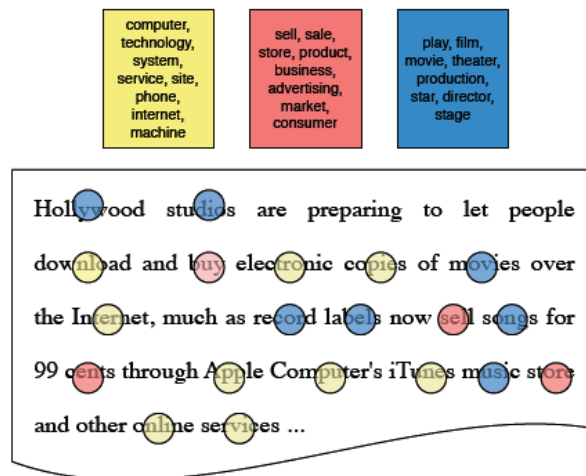
Hollywood studios are preparing to let people download and buy electronic copies of movies over the Internet, much as record labels now sell songs for 99 cents through Apple Computer's iTunes music store and other online services ...



Algorithm for Extracting Topics

❖ Gibbs Sampling

➤ Inference



And repeat, conditioning $z_{d,n}$ on all of the other assignments



Algorithm for Extracting Topics

❖ Gibbs Sampling

➤ A form of Markov chain Monte Carlo (MCMC)

- Hofmann (1999) used the expectation-maximization (EM) algorithm to obtain direct estimates of ϕ and θ ; suffers from problems involving local maxima.
- Markov chain Monte Carlo (MCMC) refers to a set of approximate iterative techniques designed to sample values from complex (often high-dimensional) distributions
- To simulate a high-dimensional distribution by sampling on lower-dimensional subsets of variables where each subset is conditioned on the value of all others.
- The sampling is done sequentially and proceeds until the sampled values approximate the target distribution.
- While the Gibbs procedure we will describe does not provide direct estimates of ϕ and θ , we will show how ϕ and θ can be approximated using posterior estimates of z .



Algorithm for Extracting Topics

❖ Gibbs Sampling

➤ Algorithm

- To consider each word token in the text collection in turn, and estimates the probability of assigning the current word token to each topic, conditioned on the topic assignments to all other word tokens.
- From this conditional distribution, a topic is sampled and stored as the new topic assignment for this word token.
- $P(z_i = j | \mathbf{z}_{-i}, w_i, d_i, \cdot)$,
 - ✓ $z_i = j$ represents the topic assignment of token i to topic j , \mathbf{z}_{-i} refers to the topic assignments of all other word tokens
 - ✓ “ \cdot ” refers to all other known or observed information such as all other word and document indices w_i and d_i , and hyperparameters α , and β .

$$P(z_i = j | \mathbf{z}_{-i}, w_i, d_i, \cdot) \propto \frac{C_{w_i j}^{WT} + \beta}{\sum_{w=1}^W C_{w j}^{WT} + W\beta} \frac{C_{d_i j}^{DT} + \alpha}{\sum_{t=1}^T C_{d_i t}^{DT} + T\alpha}$$



Algorithm for Extracting Topics

❖ Gibbs Sampling

➤ Algorithm

$$P(z_i = j | \mathbf{z}_{-i}, w_i, d_i, \cdot) \propto \frac{C_{w_i j}^{WT} + \beta}{\sum_{w=1}^W C_{w j}^{WT} + W\beta} \frac{C_{d_i j}^{DT} + \alpha}{\sum_{t=1}^T C_{d_i t}^{DT} + T\alpha}$$

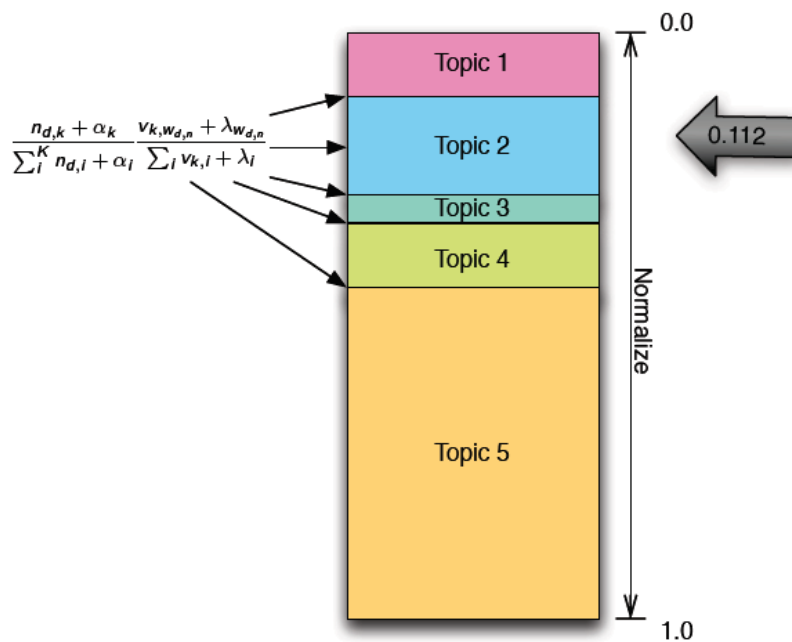
- C^{WT} and C^{DT} are matrices of counts with dimensions $W \times T$ and $D \times T$ respectively;
- $C_{w_j}^{WT}$ contains the number of times word w is assigned to topic j , not including the current instance i
- $C_{d_j}^{DT}$ contains the number of times topic j is assigned to some word token in document d , not including the current instance i .
- Note that this equation gives the unnormalized probability.
- The actual probability of assigning a word token to topic j is calculated by dividing the quantity in this equation for topic t by the sum over all topics T .
- The left part is the probability of word w under topic j
- The right part is the probability that topic j has under the current topic distribution for document d .
- Words are assigned to topics depending on how likely the word is for a topic, as well as how dominant a topic is in a document



Algorithm for Extracting Topics

❖ Gibbs Sampling

➤ How to sample from a distribution?



Algorithm for Extracting Topics

❖ Gibbs Sampling

➤ Implementation

Algorithm

- 1 For each iteration i :
 - 1 For each document d and word n currently assigned to z_{old} :
 - 1 Decrement $n_{d,z_{old}}$ and $v_{z_{old},w_{d,n}}$
 - 2 Sample $z_{new} = k$ with probability proportional to
$$\frac{n_{d,k} + \alpha_k}{\sum_i^K n_{d,i} + \alpha_i} \frac{v_{k,w_{d,n}} + \lambda_{w_{d,n}}}{\sum_i v_{k,i} + \lambda_i}$$
 - 3 Increment $n_{d,z_{new}}$ and $v_{z_{new},w_{d,n}}$



Algorithm for Extracting Topics

❖ Gibbs Sampling

➤ Algorithm

$$P(z_i = j | \mathbf{z}_{-i}, w_i, d_i, \cdot) \propto \frac{C_{w_i j}^{WT} + \beta}{\sum_{w=1}^W C_{w j}^{WT} + W\beta} \frac{C_{d_i j}^{DT} + \alpha}{\sum_{t=1}^T C_{d_i t}^{DT} + T\alpha}$$

- starts by assigning each word token to a random topic in [1..T].
- For each word token, the count matrices C^{WT} and C^{DT} are first decremented by one for the entries that correspond to the current topic assignment.
- Then, a new topic is sampled from the distribution in this equation and the count matrices C^{WT} and C^{DT} are incremented with the new topic assignment.
- Each Gibbs sample consists the set of topic assignments to all N word tokens in the corpus, achieved by a single pass through all documents.
- During the initial stage of the sampling process (also known as the **burnin period**), the Gibbs samples have to be discarded because they are poor estimates of the posterior.
- After the burnin period, the successive Gibbs samples start to approximate the target distribution (i.e., the posterior distribution over topic assignments)



Algorithm for Extracting Topics

❖ Gibbs Sampling

➤ Estimating parameters

$$\phi_i^{(j)} = \frac{C_{ij}^{WT} + \beta}{\sum_{k=1}^W C_{kj}^{WT} + W\beta} \quad \theta_j^{(d)} = \frac{C_{dj}^{DT} + \alpha}{\sum_{k=1}^T C_{dk}^{DT} + T\alpha}$$

➤ An Example

- The Gibbs sampling algorithm can be illustrated by generating artificial data from a known topic model and applying the algorithm to check whether it is able to infer the original generative structure

- Topic 1 gives equal probability to words MONEY, LOAN, and BANK

$$\phi_{MONEY}^{(1)} = \phi_{LOAN}^{(1)} = \phi_{BANK}^{(1)} = 1/3.$$

- Topic 2 gives equal probability to words RIVER, STREAM, and BANK

$$\phi_{RIVER}^{(2)} = \phi_{STREAM}^{(2)} = \phi_{BANK}^{(2)} = 1/3.$$

- 16 documents generated by arbitrarily mixing the two topics
- the color of the circles indicate the topic assignments (black=topic 1; white=topic 2).
- after 64 iterations: Topic 1 and Topic 2

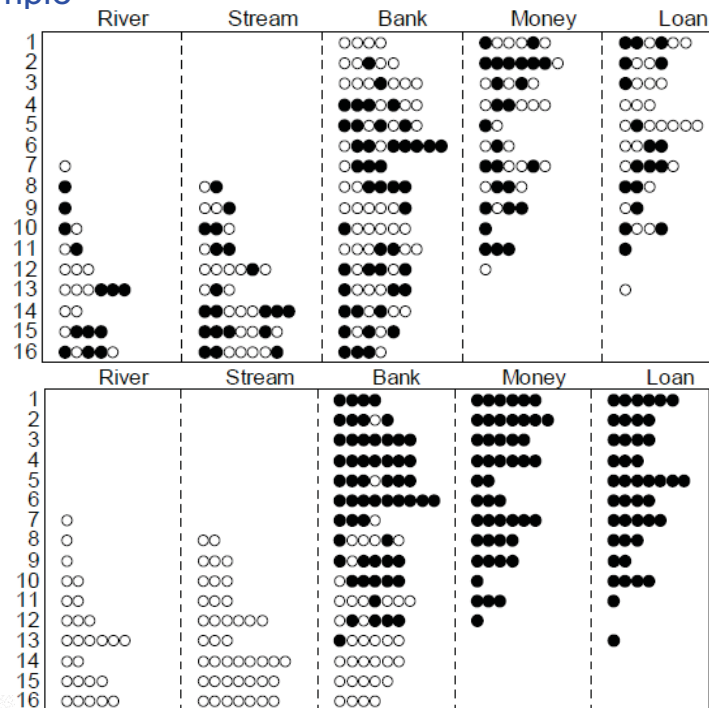
$$\phi_{MONEY}^{(1)} = .32, \quad \phi_{LOAN}^{(1)} = .29, \quad \phi_{BANK}^{(1)} = .39 \quad \phi_{RIVER}^{(2)} = .25, \quad \phi_{STREAM}^{(2)} = .4, \quad \phi_{BANK}^{(2)} = .35.$$



Algorithm for Extracting Topics

❖ Gibbs Sampling

➤ An Example



Algorithm for Extracting Topics

❖ Gibbs Sampling

➤ Exchangeability of topics

- when topics are used to calculate a statistic which is invariant to the ordering of the topics, it becomes possible and even important to average over different Gibbs samples
- Model averaging is likely to improve results because it allows sampling from multiple local modes of the posterior.



Algorithm for Extracting Topics

❖ Gibbs Sampling

- Hyperparameters: Sample them too (slice sampling)
- Initialization: Random
- Sampling: Until likelihood converges
- Lag / burn-in: Difference of opinion on this
- Number of chains: Should do more than one

➤ Available implementations

- Mallet (<http://mallet.cs.umass.edu>)
- LDAC (<http://www.cs.princeton.edu/blei/lda-c>)
- Topicmod (<http://code.google.com/p/topicmod>)



Polysemy with Topics

❖ Probabilistic topic models represent semantic ambiguity through uncertainty over topics.

➤ A 300 topic solution for the TASA corpus

- The word PLAY is given relatively high probability related to the different senses of the word (playing music, theater play, playing games).

Topic 77		Topic 82		Topic 166	
word	prob.	word	prob.	word	prob.
MUSIC	.090	LITERATURE	.031	PLAY	.136
DANCE	.034	POEM	.028	BALL	.129
SONG	.033	POETRY	.027	GAME	.065
PLAY	.030	POET	.020	PLAYING	.042
SING	.026	PLAYS	.019	HIT	.032
SINGING	.026	POEMS	.019	PLAYED	.031
BAND	.026	PLAY	.015	BASEBALL	.027
PLAYED	.023	LITERARY	.013	GAMES	.025
SANG	.022	WRITERS	.013	BAT	.019
SONGS	.021	DRAMA	.012	RUN	.019
DANCING	.020	WROTE	.012	THROW	.016
PIANO	.017	POETS	.011	BALLS	.015
PLAYING	.016	WRITER	.011	TENNIS	.011
RHYTHM	.015	SHAKESPEARE	.010	HOME	.010
ALBERT	.013	WRITTEN	.009	CATCH	.010
MUSICAL	.013	STAGE	.009	FIELD	.010



Polysemy with Topics

- Fragments of three documents are shown from TASA that use PLAY in three different senses.
- The presence of other less ambiguous words (e.g., MUSIC in the first document) builds up evidence for a particular topic in the document.
- When a word has uncertainty over topics, the topic distribution developed for the document context is the primary factor for disambiguating the word.

Document #29795

Bix beiderbecke, at age⁰⁶⁰ fifteen²⁰⁷, sat¹⁷⁴ on the slope⁰⁷¹ of a bluff⁰⁵⁵ overlooking⁰²⁷ the mississippi¹³⁷ river¹³⁷. He was listening⁰⁷⁷ to music⁰⁷⁷ coming⁰⁰⁹ from a passing⁰⁴³ riverboat. The music⁰⁷⁷ had already captured⁰⁰⁶ his heart¹⁵⁷, as well as his ear¹¹⁹. It was jazz⁰⁷⁷. Bix beiderbecke had already had music⁰⁷⁷ lessons⁰⁷⁷. He showed⁰⁰² promise¹³⁴ on the piano⁰⁷⁷, and his parents⁰³⁵ hoped²⁶⁸ he might consider¹¹⁸ becoming a concert⁰⁷⁷ pianist⁰⁷⁷. But bix was interested²⁶⁸ in another kind⁰⁵⁰ of music⁰⁷⁷. He wanted²⁶⁸ to play⁰⁷⁷ the comet. And he wanted²⁶⁸ to play⁰⁷⁷ jazz⁰⁷⁷ ...

Document #1883

There is a simple⁰⁵⁰ reason¹⁰⁶ why there are so few periods⁰⁷⁸ of really great theater⁰⁸² in our whole western⁰⁴⁶ world. Too many things³⁰⁰ have to come right at the very same time. The dramatists must have the right actors⁰⁸², the actors⁰⁸² must have the right playhouses, the playhouses must have the right audiences⁰⁸². We must remember²⁸⁸ that plays⁰⁸² exist¹⁴³ to be performed⁰⁷⁷, not merely⁰⁵⁰ to be read²⁵⁴ (even when you read²⁵⁴ a play⁰⁸² to yourself, try²⁸⁸ to perform⁰⁶² it, to put¹⁷⁴ it on a stage⁰⁷⁸, as you go along) as soon⁰²⁸ as a play⁰⁸² has to be performed⁰⁸², then some kind¹²⁶ of theatrical⁰⁸² ...

Document #21359

Jim²⁹⁶ has a game¹⁶⁶ book²⁵⁴. Jim²⁹⁶ reads²⁵⁴ the book²⁵⁴. Jim²⁹⁶ sees⁰⁸¹ a game¹⁶⁶ for one. Jim²⁹⁶ plays¹⁶⁶ the game¹⁶⁶. Jim²⁹⁶ likes⁰⁸¹ the game¹⁶⁶ for one. The game¹⁶⁶ book²⁵⁴ helps⁰⁸¹ jim²⁹⁶. Don¹⁸⁰ comes⁰⁴⁰ into the house⁰³⁸. Don¹⁸⁰ and jim²⁹⁶ read²⁵⁴ the game¹⁶⁶ book²⁵⁴. The boys⁰²⁰ see a game¹⁶⁶ for two. The two boys⁰²⁰ play¹⁶⁶ the game¹⁶⁶. The boys⁰²⁰ play¹⁶⁶ the game¹⁶⁶ for two. The boys⁰²⁰ like the game¹⁶⁶. Meg²⁸² comes⁰⁴⁰ into the house²⁸². Meg²⁸² and don¹⁸⁰ and jim²⁹⁶ read²⁵⁴ the book²⁵⁴. They see a game¹⁶⁶ for three. Meg²⁸² and don¹⁸⁰ and jim²⁹⁶ play¹⁶⁶ the game¹⁶⁶. They play¹⁶⁶ ...



Computing Similarities

❖ The similarity of words and documents

- Two words are similar to the extent that they appear in the same topics, and two documents are similar to the extent that the same topics appear in those documents

➤ Similarity between documents

- The similarity between documents d_1 and d_2 can be measured by the similarity between their corresponding topic distributions and $\theta^{(d_1)}$ and $\theta^{(d_2)}$.
- The Kullback Leibler (KL) divergence

$$D(p, q) = \sum_{j=1}^T p_j \log_2 \frac{p_j}{q_j} \quad KL(p, q) = \frac{1}{2} [D(p, q) + D(q, p)]$$

- ✓ This non-negative function is equal to zero when for all j , $p_j = q_j$.

- The symmetrized Jensen-Shannon (JS) divergence

$$JS(p, q) = \frac{1}{2} [D(p, (p+q)/2) + D(q, (p+q)/2)]$$

- ✓ measures similarity between p and q through the average of p and q
- ✓ two distributions p and q will be similar if they are similar to their average $(p+q)/2$

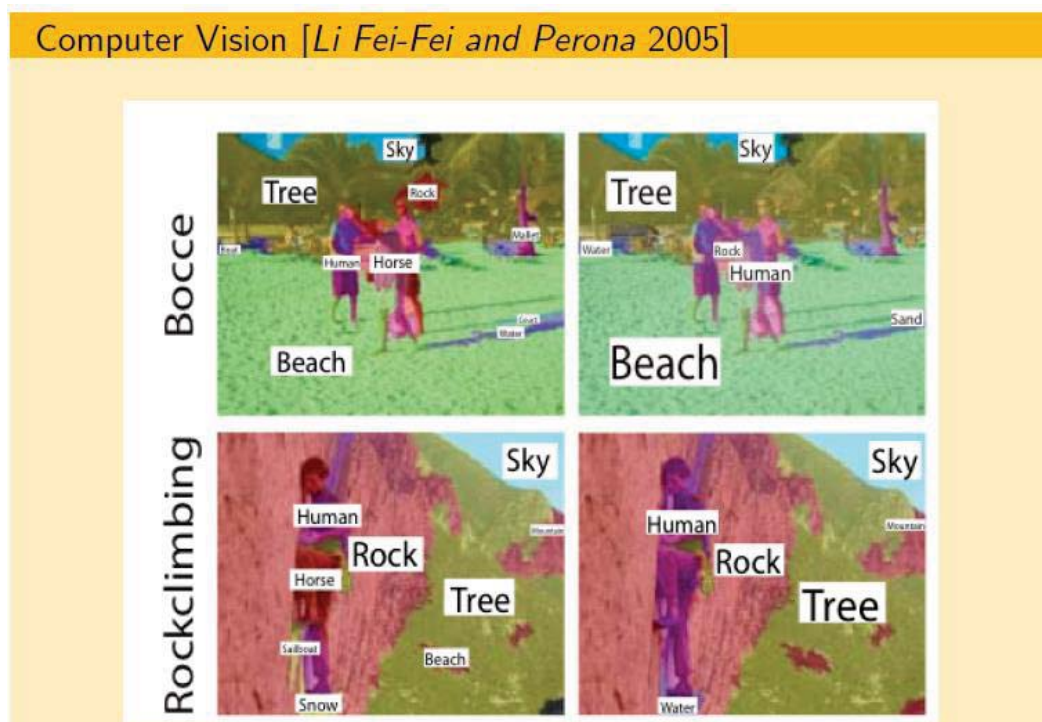


Applications

- What's a document?
- What's a word?
- What's your vocabulary?
- How do you evaluate?

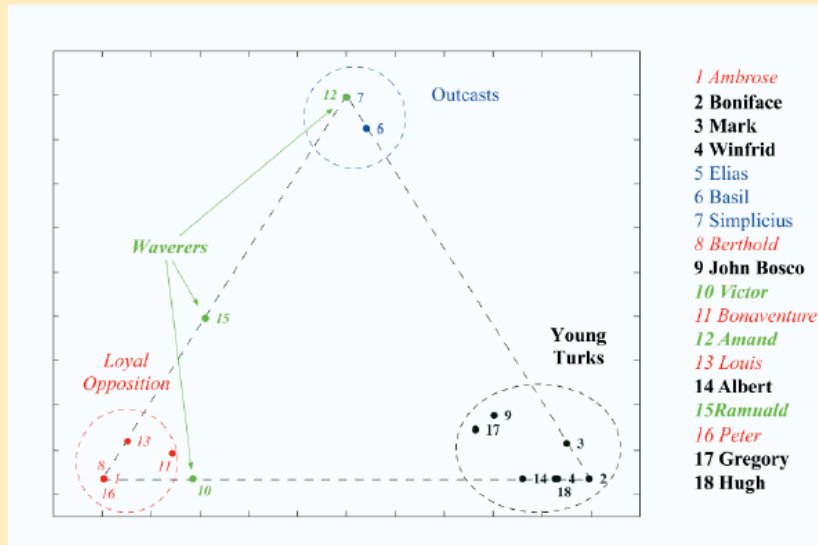


Applications



Applications

Social Networks [Airoldi et al. 2008]



Applications

Music [Hu and Saul 2009]



Thank you for your attention !!

My websites:

web.donga.ac.kr/yjko/

Islab.donga.ac.kr

Email address:

youngjoong.ko@gmail.com

