Department of Computer Science
UNIVERSITY OF COLORADO **BOULDER**

Topic Models

Machine Learning: Jordan Boyd-Graber
University of Colorado Boulder
LECTURE 18

- Suppose you have a huge number of documents
- Want to know what's going on
- Can't read them all (e.g. every New York Times article from the 90's)
- Topic models offer a way to get a corpus-level view of major themes
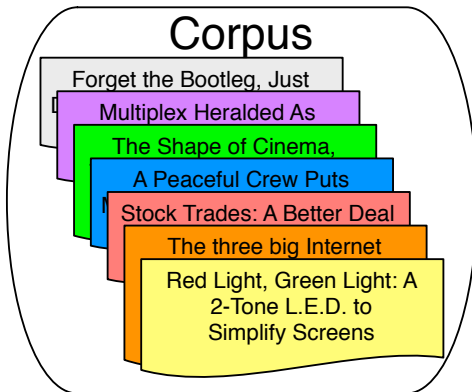
**Why topic models?**

- Suppose you have a huge number of documents
- Want to know what's going on
- Can't read them all (e.g. every New York Times article from the 90's)
- Topic models offer a way to get a corpus-level view of major themes
- Unsupervised

**Roadmap**

- What are topic models
- How to know if you have good topic model
- How to go from raw data to topics

## Conceptual Approach

From an **input corpus** and number of topics $K \rightarrow$ words to topics

## Conceptual Approach

From an input corpus and number of topics $K \rightarrow$ **words to topics**

## Conceptual Approach

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

## 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 |

**Why should you care?**

- Neat way to explore / understand corpus collections
  - E-discovery
  - Social media
  - Scientific data
- NLP Applications
  - Word Sense Disambiguation
  - Discourse Segmentation
  - Machine Translation
- Psychology: word meaning, polysemy
- Inference is (relatively) simple

## Matrix Factorization Approach



$$\begin{bmatrix} M \times K \end{bmatrix} \times \begin{bmatrix} K \times V \end{bmatrix} \approx \begin{bmatrix} M \times V \end{bmatrix}$$

Topic Assignment    Topics    Dataset

K Number of topics
M Number of documents
V Size of vocabulary

**Matrix Factorization Approach**



$$\begin{bmatrix} M \times K \end{bmatrix} \times \begin{bmatrix} K \times V \end{bmatrix} \approx \begin{bmatrix} M \times V \end{bmatrix}$$

Topic Assignment     Topics     Dataset

K Number of topics
M Number of documents
V Size of vocabulary

- If you use singular value decomposition (SVD), this technique is called latent semantic analysis.
- Popular in information retrieval.
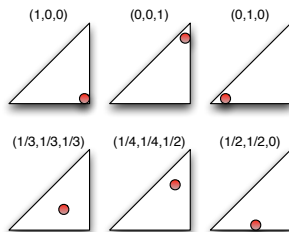
## Alternative: Generative Model

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

## Alternative: Generative Model

- How your data came to be
- Sequence of Probabilistic Steps
- Posterior Inference
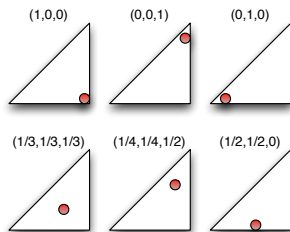- Blei, Ng, Jordan. Latent **Dirichlet** Allocation. JMLR, 2003.

## Multinomial Distribution

- Distribution over discrete outcomes
- Represented by non-negative vector that sums to one
- Picture representation

**Multinomial Distribution**

- Distribution over discrete outcomes
- Represented by non-negative vector that sums to one
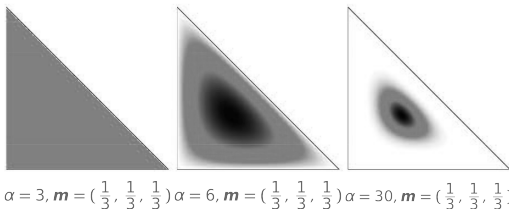- Picture representation



- Come from a Dirichlet distribution

## Dirichlet Distribution

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

## Dirichlet Distribution

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



$\alpha = 3, \boldsymbol{m} = (\frac{1}{3}, \frac{1}{3}, \frac{1}{3})$ $\alpha = 6, \boldsymbol{m} = (\frac{1}{3}, \frac{1}{3}, \frac{1}{3})$ $\alpha = 30, \boldsymbol{m} = (\frac{1}{3}, \frac{1}{3}, \frac{1}{3})$
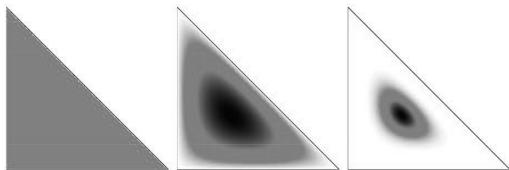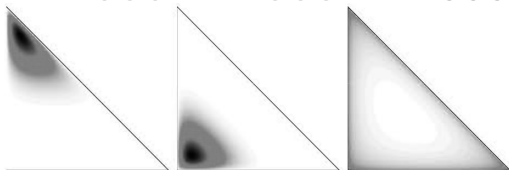
## Dirichlet Distribution

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



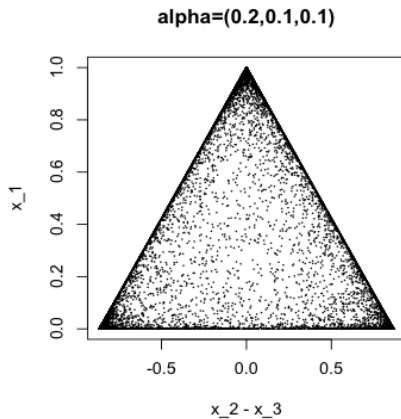$\alpha = 3, \boldsymbol{m} = (\frac{1}{3}, \frac{1}{3}, \frac{1}{3}) \quad \alpha = 6, \boldsymbol{m} = (\frac{1}{3}, \frac{1}{3}, \frac{1}{3}) \quad \alpha = 30, \boldsymbol{m} = (\frac{1}{3}, \frac{1}{3}, \frac{1}{3})$

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

## Dirichlet Distribution

## Dirichlet Distribution

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

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

$$\propto \prod_k \phi^{n_k} \prod_k \phi^{\alpha_k - 1} \tag{2}$$

$$\propto \prod_k \phi^{\alpha_k + n_k - 1} \tag{3}$$

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

## Dirichlet Distribution

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

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

$$\propto \prod_k \phi^{n_k} \prod_k \phi^{\alpha_k - 1} \tag{2}$$

$$\propto \prod_k \phi^{\alpha_k + n_k - 1} \tag{3}$$

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

## Generative Model

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

## Generative Model

## Generative Model



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 ...

## Generative Model
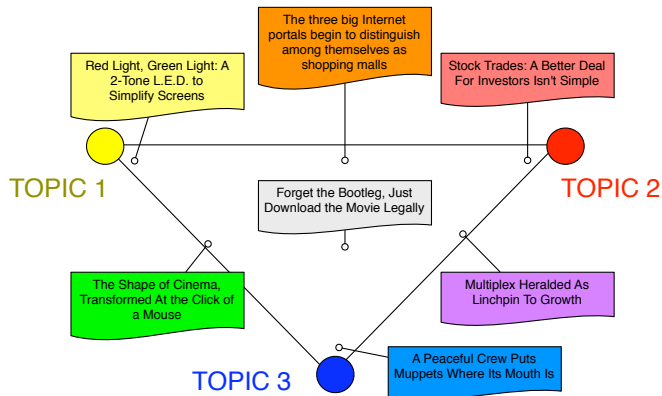


computer, technology, system, service, site, phone, internet, machine

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 ...
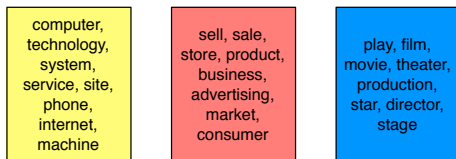
## Generative Model



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

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

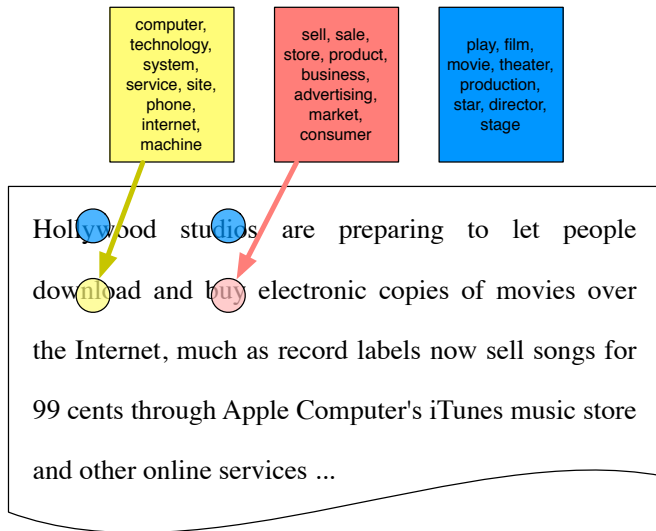play, film, movie, theater, production, star, director, stage

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 ...
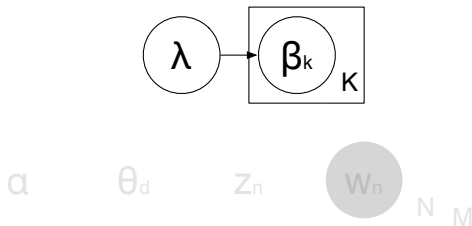
## Generative Model



| computer, technology, system, service, site, phone, internet, machine | sell, sale, store, product, business, advertising, market, consumer | play, film, movie, theater, production, star, director, stage |

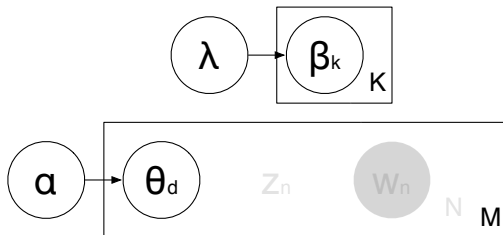Hollywood studios are preparing to let people download and buy electronic copies of movies over the Internet, much as record labels now sell songs for 99 cents through Apple Computer's iTunes music store and other online services ...

**Generative Model Approach**



- For each topic $k \in \{1, \ldots, K\}$, draw a multinomial distribution $\beta_k$ from a Dirichlet distribution with parameter $\lambda$

**Generative Model Approach**



- For each topic $k \in \{1, \dots, K\}$, draw a multinomial distribution $\beta_k$ from a Dirichlet distribution with parameter $\lambda$
- For each document $d \in \{1, \dots, M\}$, draw a multinomial distribution $\theta_d$ from a Dirichlet distribution with parameter $\alpha$

**Generative Model Approach**



- For each topic $k \in \{1, \ldots, K\}$, draw a multinomial distribution $\beta_k$ from a Dirichlet distribution with parameter $\lambda$
- For each document $d \in \{1, \ldots, M\}$, draw a multinomial distribution $\theta_d$ from a Dirichlet distribution with parameter $\alpha$
- For each word position $n \in \{1, \ldots, N\}$, select a hidden topic $z_n$ from the multinomial distribution parameterized by $\theta$.
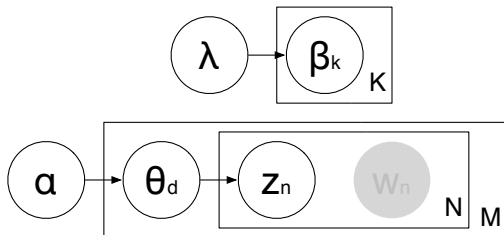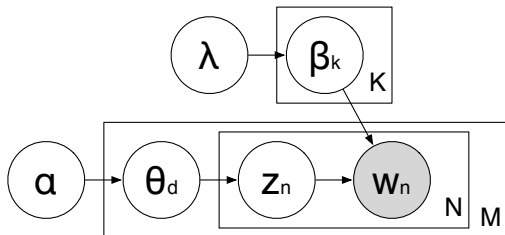
**Generative Model Approach**



- For each topic $k \in \{1, \ldots, K\}$, draw a multinomial distribution $\beta_k$ from a Dirichlet distribution with parameter $\lambda$
- For each document $d \in \{1, \ldots, M\}$, draw a multinomial distribution $\theta_d$ from a Dirichlet distribution with parameter $\alpha$
- For each word position $n \in \{1, \ldots, N\}$, select a hidden topic $z_n$ from the multinomial distribution parameterized by $\theta$.
- Choose the observed word $w_n$ from the distribution $\beta_{z_n}$.
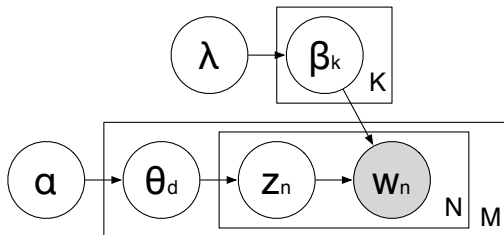
**Generative Model Approach**



- For each topic $k \in \{1, \ldots, K\}$, draw a multinomial distribution $\beta_k$ from a Dirichlet distribution with parameter $\lambda$
- For each document $d \in \{1, \ldots, M\}$, draw a multinomial distribution $\theta_d$ from a Dirichlet distribution with parameter $\alpha$
- For each word position $n \in \{1, \ldots, N\}$, select a hidden topic $z_n$ from the multinomial distribution parameterized by $\theta$.
- Choose the observed word $w_n$ from the distribution $\beta_{z_n}$.

**Topic Models: What's Important**

- Topic models
  - Topics to word types—multinomial distribution
  - Documents to topics—multinomial distribution
- Focus in this talk: statistical methods
  - Model: story of how your data came to be
  - Latent variables: missing pieces of your story
  - Statistical inference: filling in those missing pieces
- We use latent Dirichlet allocation (LDA), a fully Bayesian version of pLSI, probabilistic version of LSA
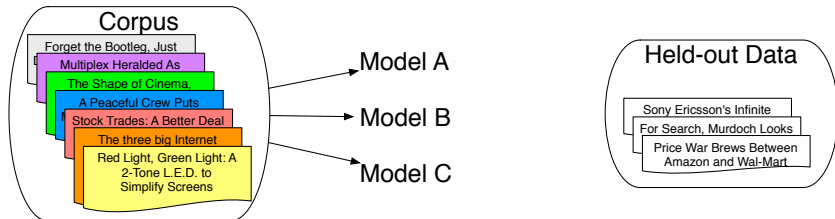
**Topic Models: What's Important**

- Topic models (latent variables)
  - Topics to word types—multinomial distribution
  - Documents to topics—multinomial distribution
- Focus in this talk: statistical methods
  - Model: story of how your data came to be
  - Latent variables: missing pieces of your story
  - Statistical inference: filling in those missing pieces
- We use latent Dirichlet allocation (LDA), a fully Bayesian version of pLSI, probabilistic version of LSA

## Evaluation



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

How you compute it is important too (Wallach et al. 2009)

## Evaluation



Measures predictive power, not what the topics are

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

How you compute it is important too (Wallach et al. 2009)

## Word Intrusion

TOPIC 1

TOPIC 2

TOPIC 3

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

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

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

**Word Intrusion**

1. Take the highest probability words from a topic

Original Topic

dog, cat, horse, pig, cow

**Word Intrusion**

1. Take the highest probability words from a topic

Original Topic

dog, cat, horse, pig, cow

2. Take a high-probability word from another topic and add it

Topic with Intruder

dog, cat, apple, horse, pig, cow

**Word Intrusion**

1. Take the highest probability words from a topic

Original Topic

dog, cat, horse, pig, cow

2. Take a high-probability word from another topic and add it

Topic with Intruder

dog, cat, apple, horse, pig, cow

3. We ask users to find the word that doesn't belong

Hypothesis

If the topics are interpretable, users will consistently choose true intruder

## Word Intrusion

| 1 / 10 | | | | | |
|---|---|---|---|---|---|
| crash | accident | board | agency | tibetan | safety |

| 2 / 10 | | | | | |
|---|---|---|---|---|---|
| commercial | network | television | advertising | viewer | layoff |

| 3 / 10 | | | | | |
|---|---|---|---|---|---|
| arrest | crime | inmate | pitcher | prison | death |

| 4 / 10 | | | | | |
|---|---|---|---|---|---|
| hospital | doctor | health | care | medical | tradition |

**Word Intrusion**



| 1 / 10 | | | Reveal additional response | | |
| --- | --- | --- | --- | --- | --- |
| crash | accident | board | agency | tibetan | safety |

| 2 / 10 | | | | | |
| --- | --- | --- | --- | --- | --- |
| commercial | network | television | advertising | viewer | layoff |

| 3 / 10 | | | | | |
| --- | --- | --- | --- | --- | --- |
| arrest | crime | inmate | pitcher | prison | death |

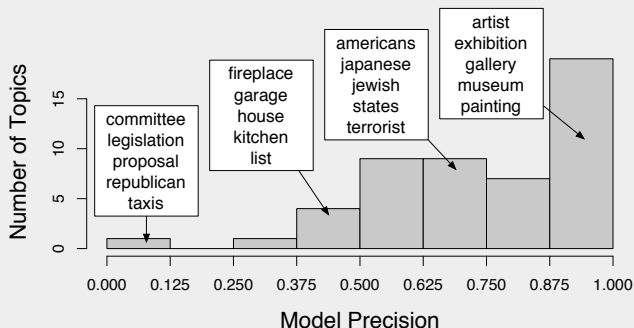| 4 / 10 | | | | | |
| --- | --- | --- | --- | --- | --- |
| hospital | doctor | health | care | medical | tradition |

- Order of words was shuffled
- Which intruder was selected varied
- Model precision: percentage of users who clicked on intruder

**Word Intrusion: Which Topics are Interpretable?**
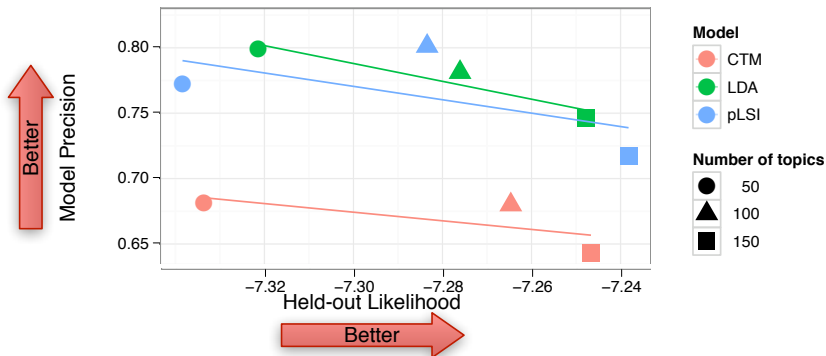
New York Times, 50 LDA Topics



Model Precision: percentage of correct intruders found
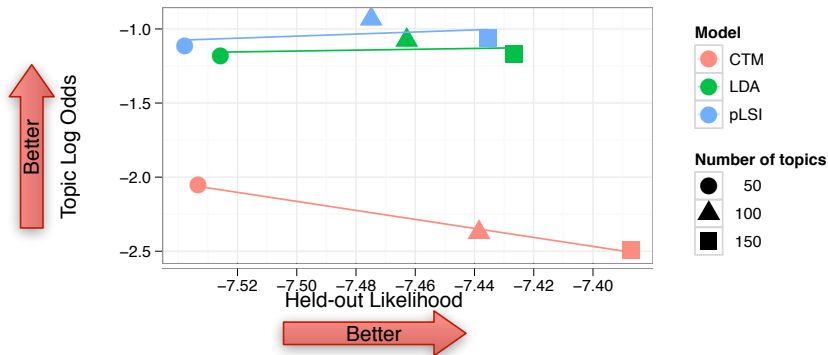
## Interpretability and Likelihood



Model Precision on New York Times

within a model, higher likelihood $\neq$ higher interpretability

## Interpretability and Likelihood



Topic Log Odds on Wikipedia

across models, higher likelihood $\neq$ higher interpretability

## Evaluation Takeaway

- Measure what you care about
- If you care about prediction, likelihood is good
- If you care about a particular task, measure that

## Inference

- We are interested in posterior distribution

$$p(Z|X, \Theta) \qquad (4)$$

**Inference**

- We are interested in posterior distribution

$$p(Z|X, \Theta) \qquad (4)$$

- Here, latent variables are topic assignments $z$ and topics $\theta$. $X$ is the words (divided into documents), and $\Theta$ are hyperparameters to Dirichlet distributions: $\alpha$ for topic proportion, $\lambda$ for topics.

$$p(\mathbf{z}, \boldsymbol{\beta}, \boldsymbol{\theta}|\mathbf{w}, \alpha, \lambda) \qquad (5)$$

**Inference**

- We are interested in posterior distribution

$$p(Z|X, \Theta) \qquad (4)$$

- Here, latent variables are topic assignments $z$ and topics $\theta$. $X$ is the words (divided into documents), and $\Theta$ are hyperparameters to Dirichlet distributions: $\alpha$ for topic proportion, $\lambda$ for topics.

$$p(\mathbf{z}, \boldsymbol{\beta}, \boldsymbol{\theta} | \mathbf{w}, \alpha, \lambda) \qquad (5)$$

$$p(\mathbf{w}, \mathbf{z}, \boldsymbol{\theta}, \boldsymbol{\beta} | \alpha, \lambda) = \\ \prod_k p(\beta_k | \lambda) \prod_d p(\theta_d | \alpha) \prod_n p(z_{d,n} | \theta_d) p(w_{d,n} | \beta_{z_{d,n}})$$

## Gibbs Sampling

- A form of Markov Chain Monte Carlo
- Chain is a sequence of random variable states
- Given a state $\{z_1, \ldots z_N\}$ given certain technical conditions, drawing $z_k \sim p(z_1, \ldots z_{k-1}, z_{k+1}, \ldots z_N | X, \Theta)$ for all $k$ (repeatedly) results in a Markov Chain whose stationary distribution *is* the posterior.
- For notational convenience, call **z** with $z_{d,n}$ removed $\mathbf{z}_{-d,n}$

## Inference

## 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 ...

## Inference

## Inference

## Inference



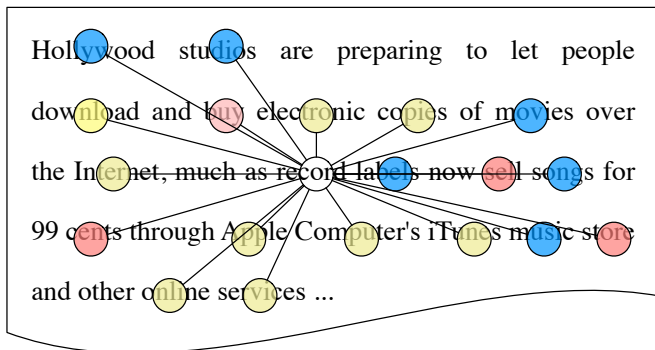Machine Learning: Jordan Boyd-Graber | Boulder

## Inference

## 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 ...

**Gibbs Sampling**
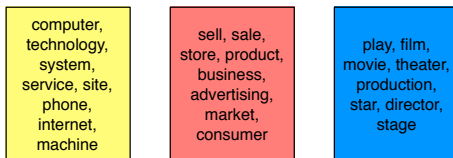
- For LDA, we will sample the topic assignments
- Thus, we want:

$$p(z_{d,n} = k | \mathbf{z}_{-d,n}, \mathbf{w}, \alpha, \lambda) = \frac{p(z_{d,n} = k, \mathbf{z}_{-d,n} | \mathbf{w}, \alpha, \lambda)}{p(\mathbf{z}_{-d,n} | \mathbf{w}, \alpha, \lambda)}$$

**Gibbs Sampling**

- For LDA, we will sample the topic assignments
- Thus, we want:

$$p(z_{d,n} = k | \mathbf{z}_{-d,n}, \mathbf{w}, \alpha, \lambda) = \frac{p(z_{d,n} = k, \mathbf{z}_{-d,n} | \mathbf{w}, \alpha, \lambda)}{p(\mathbf{z}_{-d,n} | \mathbf{w}, \alpha, \lambda)}$$

- The topics and per-document topic proportions are integrated out / marginalized
- Let $n_{d,i}$ be the number of words taking topic $i$ in document $d$. Let $v_{k,w}$ be the number of times word $w$ is used in topic $k$.

$$= \frac{\int_{\theta_d} \left( \prod_{i \neq k} \theta_d^{\alpha_i + n_{d,i} - 1} \right) \theta_d^{\alpha_k + n_{d,i}} d\theta_d \int_{\beta_k} \left( \prod_{i \neq w_{d,n}} \beta_{k,i}^{\lambda_i + v_{k,i} - 1} \right) \beta_{k,w_{d,n}}^{\lambda_i + v_{k,i}} d\beta_k}{\int_{\theta_d} \left( \prod_i \theta_d^{\alpha_i + n_{d,i} - 1} \right) d\theta_d \int_{\beta_k} \left( \prod_i \beta_{k,i}^{\lambda_i + v_{k,i} - 1} \right) d\beta_k}$$

**Gibbs Sampling**

- Integral is normalizer of Dirichlet distribution

$$\int_{\beta_k} \left( \prod_i \beta_{k,i}^{\lambda_i + v_{k,i} - 1} \right) d\beta_k = \frac{\prod_i^V \Gamma \left( \beta_i + v_{k,i} \right)}{\Gamma \left( \sum_i^V \beta_i + v_{k,i} \right)}$$

**Gibbs Sampling**

- Integral is normalizer of Dirichlet distribution

$$\int_{\beta_k} \left( \prod_i \beta_{k,i}^{\lambda_i + v_{k,i} - 1} \right) d\beta_k = \frac{\prod_i^V \Gamma\left(\beta_i + v_{k,i}\right)}{\Gamma\left(\sum_i^V \beta_i + v_{k,i}\right)}$$

- So we can simplify

$$\frac{\int_{\theta_d} \left( \prod_{i \neq k} \theta_d^{\alpha_i + n_{d,i} - 1} \right) \theta_d^{\alpha_k + n_{d,i}} d\theta_d \int_{\beta_k} \left( \prod_{i \neq w_{d,n}} \beta_{k,i}^{\lambda_i + v_{k,i} - 1} \right) \beta_{k,w_{d,n}}^{\lambda_i + v_{k,i}} d\beta_k}{\int_{\theta_d} \left( \prod_i \theta_d^{\alpha_i + n_{d,i} - 1} \right) d\theta_d \int_{\beta_k} \left( \prod_i \beta_{k,i}^{\lambda_i + v_{k,i} - 1} \right) d\beta_k} =$$

$$\frac{\frac{\Gamma\left(\alpha_k + n_{d,k} + 1\right)}{\Gamma\left(\sum_i^K \alpha_i + n_{d,i} + 1\right)} \prod_{i \neq k}^K \Gamma\left(\alpha_k + n_{d,k}\right)}{\frac{\prod_i^K \Gamma\left(\alpha_i + n_{d,i}\right)}{\Gamma\left(\sum_i^K \alpha_i + n_{d,i}\right)}} \frac{\frac{\Gamma\left(\lambda_{w_{d,n}} + v_{k,w_{d,n}} + 1\right)}{\Gamma\left(\sum_i^V \lambda_i + v_{k,i} + 1\right)} \prod_{i \neq w_{d,n}}^V \Gamma\left(\lambda_k + v_{k,w_{d,n}}\right)}{\frac{\prod_i^V \Gamma\left(\lambda_i + v_{k,i}\right)}{\Gamma\left(\sum_i^V \lambda_i + v_{k,i}\right)}}$$

Gamma Function Identity

$$z = \frac{\Gamma(z+1)}{\Gamma(z)} \tag{6}$$

$$\frac{\frac{\Gamma\left(\alpha_k+n_{d,k}+1\right)}{\Gamma\left(\sum_i^K \alpha_i+n_{d,i}+1\right)} \prod_{i \neq k}^K \Gamma\left(\alpha_k+n_{d,k}\right)}{\frac{\prod_i^K \Gamma\left(\alpha_i+n_{d,i}\right)}{\Gamma\left(\sum_i^K \alpha_i+n_{d,i}\right)}} \frac{\frac{\Gamma\left(\lambda_{w_{d,n}}+v_{k,w_{d,n}}+1\right)}{\Gamma\left(\sum_i^V \lambda_i+v_{k,i}+1\right)} \prod_{i \neq w_{d,n}}^V \Gamma\left(\lambda_k+v_{k,w_{d,n}}\right)}{\frac{\prod_i^V \Gamma\left(\lambda_i+v_{k,i}\right)}{\Gamma\left(\sum_i^V \lambda_i+v_{k,i}\right)}}$$

$$= \frac{n_{d,k}+\alpha_k}{\sum_i^K n_{d,i}+\alpha_i} \frac{v_{k,w_{d,n}}+\lambda_{w_{d,n}}}{\sum_i v_{k,i}+\lambda_i}$$

**Gibbs Sampling Equation**

$$\frac{n_{d,k} + \alpha_k}{\sum_i^K n_{d,i} + \alpha_i} \frac{v_{k,w_{d,n}} + \lambda_{w_{d,n}}}{\sum_i v_{k,i} + \lambda_i} \tag{7}$$

- Number of times document $d$ uses topic $k$
- Number of times topic $k$ uses word type $w_{d,n}$
- Dirichlet parameter for document to topic distribution
- Dirichlet parameter for topic to word distribution
- How much this document likes topic $k$
- How much this topic likes word $w_{d,n}$

**Gibbs Sampling Equation**

$$\frac{n_{d,k} + \alpha_k}{\sum_i^K n_{d,i} + \alpha_i} \frac{v_{k,w_{d,n}} + \lambda_{w_{d,n}}}{\sum_i v_{k,i} + \lambda_i} \tag{7}$$

- Number of times document $d$ uses topic $k$
- Number of times topic $k$ uses word type $w_{d,n}$
- Dirichlet parameter for document to topic distribution
- Dirichlet parameter for topic to word distribution
- How much this document likes topic $k$
- How much this topic likes word $w_{d,n}$

**Gibbs Sampling Equation**

$$\frac{n_{d,k} + \alpha_k}{\sum_i^K n_{d,i} + \alpha_i} \frac{v_{k,w_{d,n}} + \lambda_{w_{d,n}}}{\sum_i v_{k,i} + \lambda_i} \tag{7}$$

- Number of times document $d$ uses topic $k$
- Number of times topic $k$ uses word type $w_{d,n}$
- Dirichlet parameter for document to topic distribution
- Dirichlet parameter for topic to word distribution
- How much this document likes topic $k$
- How much this topic likes word $w_{d,n}$

**Gibbs Sampling Equation**

$$\frac{n_{d,k} + \alpha_k}{\sum_i^K n_{d,i} + \alpha_i} \frac{v_{k,w_{d,n}} + \lambda_{w_{d,n}}}{\sum_i v_{k,i} + \lambda_i} \tag{7}$$

- Number of times document $d$ uses topic $k$
- Number of times topic $k$ uses word type $w_{d,n}$
- Dirichlet parameter for document to topic distribution
- Dirichlet parameter for topic to word distribution
- How much this document likes topic $k$
- How much this topic likes word $w_{d,n}$

## Gibbs Sampling Equation

$$\frac{n_{d,k} + \alpha_k}{\sum_i^K n_{d,i} + \alpha_i} \frac{v_{k,w_{d,n}} + \lambda_{w_{d,n}}}{\sum_i v_{k,i} + \lambda_i} \tag{7}$$

- Number of times document $d$ uses topic $k$
- Number of times topic $k$ uses word type $w_{d,n}$
- Dirichlet parameter for document to topic distribution
- Dirichlet parameter for topic to word distribution
- How much this document likes topic $k$
- How much this topic likes word $w_{d,n}$

**Gibbs Sampling Equation**

$$\frac{n_{d,k} + \alpha_k}{\sum_i^K n_{d,i} + \alpha_i} \frac{v_{k,w_{d,n}} + \lambda_{w_{d,n}}}{\sum_i v_{k,i} + \lambda_i} \tag{7}$$

- Number of times document $d$ uses topic $k$
- Number of times topic $k$ uses word type $w_{d,n}$
- Dirichlet parameter for document to topic distribution
- Dirichlet parameter for topic to word distribution
- How much this document likes topic $k$
- How much this topic likes word $w_{d,n}$

## Sample Document

| | | | | |
|---|---|---|---|---|
| Etruscan | trade | price | temple | market |

## Sample Document

| | | | | |
|---|---|---|---|---|
| Etruscan | trade | price | temple | market |

## Randomly Assign Topics



| z | 3 | 2 | 1 | 3 | 1 |
|---|---|---|---|---|---|
| w | Etruscan | trade | price | temple | market |

## Randomly Assign Topics

## Total Topic Counts

| 3 | 2 | 1 | 3 | 1 |
|---|---|---|---|---|
| Etruscan | trade | price | temple | market |

|  | 1 | 2 | 3 |
|---|---|---|---|
| Etruscan | 1 | 0 | 35 |
| market | 50 | 0 | 1 |
| price | 42 | 1 | 0 |
| temple | 0 | 0 | 20 |
| trade | 10 | 8 | 1 |
| ... |  |  |  |

Total counts from **all** docs

## Total Topic Counts

| 3 | 2 | 1 | 3 | 1 |
|---|---|---|---|---|
| Etruscan | trade | price | temple | market |

Total

|  | 1 | 2 | 3 |
|---|---|---|---|
| Etruscan | 1 | 0 | 35 |
| market | 50 | 0 | 1 |

**Sampling Equation**

$$\frac{n_{d,k} + \alpha_k}{\sum_i^K n_{d,i} + \alpha_i} \frac{v_{k,w_{d,n}} + \lambda_{w_{d,n}}}{\sum_i v_{k,i} + \lambda_i}$$

...

## Total Topic Counts

| 3 | 2 | 1 | 3 | 1 |
|---|---|---|---|---|
| Etruscan | trade | price | temple | market |

Total

|  | 1 | 2 | 3 |
|---|---|---|---|
| Etruscan | 1 | 0 | 35 |
| market | 50 | 0 | 1 |

### Sampling Equation

$$\frac{n_{d,k} + \alpha_k}{\sum_i^K n_{d,i} + \alpha_i} \frac{v_{k,w_{d,n}} + \lambda_{w_{d,n}}}{\sum_i v_{k,i} + \lambda_i}$$

...

## We want to sample this word . . .

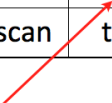| 3 | **2** | 1 | 3 | 1 |
|---|---|---|---|---|
| Etruscan | trade | price | temple | market |

| | 1 | 2 | 3 |
|---|---|---|---|
| Etruscan | 1 | 0 | 35 |
| market | 50 | 0 | 1 |
| price | 42 | 1 | 0 |
| temple | 0 | 0 | 20 |
| trade | 10 | 8 | 1 |
| ... | | | |

## We want to sample this word . . .

| 3 | **2** | 1 | 3 | 1 |
|---|---|---|---|---|
| Etruscan | trade | price | temple | market |

|  | 1 | 2 | 3 |
|---|---|---|---|
| Etruscan | 1 | 0 | 35 |
| market | 50 | 0 | 1 |
| price | 42 | 1 | 0 |
| temple | 0 | 0 | 20 |
| trade | 10 | **8** | 1 |
| ... |  |  |  |

## Decrement its count

| 3 | ? | 1 | 3 | 1 |
|---|---|---|---|---|
| Etruscan | trade | price | temple | market |

|          | 1  | 2 | 3  |
|----------|----|---|----|
| Etruscan | 1  | 0 | 35 |
| market   | 50 | 0 | 1  |
| price    | 42 | 1 | 0  |
| temple   | 0  | 0 | 20 |
| trade    | 10 | **7** | 1  |
| ...      |    |   |    |

**What is the conditional distribution for this topic?**

| 3 | ? | 1 | 3 | 1 |
|---|---|---|---|---|
| Etruscan | trade | price | temple | market |

**Part 1: How much does this document like each topic?**

| 3 | ? | 1 | 3 | 1 |
|---|---|---|---|---|
| Etruscan | trade | price | temple | market |

## Part 1: How much does this document like each topic?

| 3 | ? | 1 | 3 | 1 |
|---|---|---|---|---|
| Etruscan | trade | price | temple | market |

Topic 1      Topic 2      Topic 3

**Part 1: How much does this document like each topic?**

| 3 | ? | 1 | 3 | 1 |
|---|---|---|---|---|
| Etruscan | trade | price | temple | market |

Topic 1        Topic 2        Topic 3

Sampling Equation

$$\frac{n_{d,k} + \alpha_k}{\sum_i^K n_{d,i} + \alpha_i} \frac{v_{k,w_{d,n}} + \lambda_{w_{d,n}}}{\sum_i v_{k,i} + \lambda_i}$$

## Part 1: How much does this document like each topic?

| 3 | ? | 1 | 3 | 1 |
|---|---|---|---|---|
| Etruscan | trade | price | temple | market |

Topic 1          Topic 2          Topic 3

### Sampling Equation

$$\frac{n_{d,k} + \alpha_k}{\sum_i^K n_{d,i} + \alpha_i} \frac{v_{k,w_{d,n}} + \lambda_{w_{d,n}}}{\sum_i v_{k,i} + \lambda_i}$$

## Part 2: How much does each topic like the word?

| 3 | ? | 1 | 3 | 1 |
|---|---|---|---|---|
| Etruscan | trade | price | temple | market |

Topic 1            Topic 2            Topic 3

|       | 1  | 2 | 3 |
|-------|----|---|---|
| trade | 10 | 7 | 1 |

## Part 2: How much does each topic like the word?

| 3 | ? | 1 | 3 | 1 |
|---|---|---|---|---|
| Etruscan | trade | price | temple | market |

**Topic 1**  **Topic 2**  **Topic 3**

**Sampling Equation**

$$\frac{n_{d,k} + \alpha_k}{\sum_i^K n_{d,i} + \alpha_i} \frac{v_{k,w_{d,n}} + \lambda_{w_{d,n}}}{\sum_i v_{k,i} + \lambda_i}$$

| trade | 10 | 7 | 1 |

## Part 2: How much does each topic like the word?

| 3 | ? | 1 | 3 | 1 |
|---|---|---|---|---|
| Etruscan | trade | price | temple | market |

Topic 1          Topic 2          Topic 3

$$\frac{n_{d,k} + \alpha_k}{\sum_i^K n_{d,i} + \alpha_i} \frac{v_{k,w_{d,n}} + \lambda_{w_{d,n}}}{\sum_i v_{k,i} + \lambda_i}$$

| trade | 10 | 7 | 1 |

## Geometric interpretation

| 3 | ? | 1 | 3 | 1 |
|---|---|---|---|---|
| Etruscan | trade | price | temple | market |

Topic 1          Topic 2          Topic 3

## Geometric interpretation

| 3 | ? | 1 | 3 | 1 |
|---|---|---|---|---|
| Etruscan | trade | price | temple | market |



Topic 1          Topic 2          Topic 3

## Geometric interpretation

| 3 | ? | 1 | 3 | 1 |
|---|---|---|---|---|
| Etruscan | trade | price | temple | market |



Topic 1          Topic 2          Topic 3

## Update counts

| 3 | ? | 1 | 3 | 1 |
|---|---|---|---|---|
| Etruscan | trade | price | temple | market |

|  | 1 | 2 | 3 |
|---|---|---|---|
| Etruscan | 1 | 0 | 35 |
| market | 50 | 0 | 1 |
| price | 42 | 1 | 0 |
| temple | 0 | 0 | 20 |
| trade | **10** | 7 | 1 |
| ... | | | |

## Update counts

| 3 | **1** | 1 | 3 | 1 |
|---|---|---|---|---|
| Etruscan | trade | price | temple | market |

| | 1 | 2 | 3 |
|---|---|---|---|
| Etruscan | 1 | 0 | 35 |
| market | 50 | 0 | 1 |
| price | 42 | 1 | 0 |
| temple | 0 | 0 | 20 |
| trade | **11** | 7 | 1 |
| ... | | | |

## Update counts

| 3 | 1 | 1 | 3 | 1 |
|---|---|---|---|---|
| Etruscan | trade | price | temple | market |



Topic 1    Topic 2    Topic 3

## Details: how to sample from a distribution



$$\frac{n_{d,k} + \alpha_k}{\sum_i^K n_{d,i} + \alpha_i} \frac{v_{k,w_{d,n}} + \lambda_{w_{d,n}}}{\sum_i v_{k,i} + \lambda_i}$$

0.0

Topic 1

Topic 2

0.112

Topic 3

Topic 4

Normalize

Topic 5

1.0

## Algorithm

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

**Implementation**

## Algorithm

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

**Desiderata**

- 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)

**Wrapup**

- Topic Models: Tools to uncover themes in large document collections
- Another example of Gibbs Sampling
- In class: Gibbs sampling example