

10-708 Probabilistic Graphical Models

Case study with approximate inference: Topic Modeling

Readings:

Blei, Ng, & Jordan (2003) Griffiths & Steyvers (2004) Matt Gormley Lecture 17 March 16, 2016

Reminders

- Midway Project Report
 - Due March 23, 12:00 noon
- Mid-semester grades
- Feedback on HW1 and HW2

 Today: wrap up Slice Sampling and Hamiltonian Monte Carlo

Motivation:

Suppose you're given a massive corpora and asked to carry out the following tasks

- Organize the documents into thematic categories
- Describe the evolution of those categories over time
- Enable a domain expert to analyze and understand the content
- Find **relationships** between the categories
- Understand how **authorship** influences the content



Motivation:

Suppose you're given a massive corpora and asked to carry out the following tasks

- Organize the documents into thematic categories
- Describe the evolution of those categories over time
- Enable a domain expert to **analyze and understand** the content
- Find relationships between the categories
- Understand how authorship influences the content

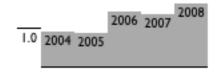
Topic Modeling:

A method of (usually unsupervised) discovery of latent or hidden structure in a corpus

- Applied primarily to text corpora, but techniques are more general
- Provides a modeling toolbox
- Has prompted the exploration of a variety of new inference methods to accommodate large-scale datasets

Dirichlet-multinomial regression (DMR) topic model on ICML (Mimno & McCallum, 2008)

Topic 0 [0.152]



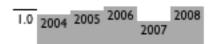
problem, optimization, problems, convex, convex optimization, linear, semidefinite programming, formulation, sets, constraints, proposed, margin, maximum margin, optimization problem, linear programming, programming, procedure, method, cutting plane, solutions

Topic 54 [0.051]



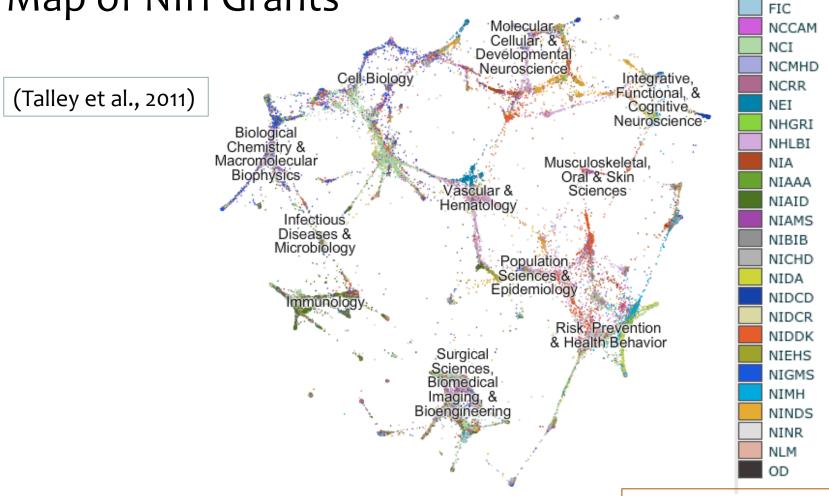
decision trees, trees, tree, decision tree, decision, tree ensemble, junction tree, decision tree learners, leaf nodes, arithmetic circuits, ensembles modts, skewing, ensembles, anytime induction decision trees, trees trees, random forests, objective decision trees, tree learners, trees grove, candidate split

Topic 99 [0.066]



inference, approximate inference, exact inference, markov chain, models, approximate, gibbs sampling, variational, bayesian, variational inference, variational bayesian, approximation, sampling, methods, exact, bayesian inference, dynamic bayesian, process, mcmc, efficient http://www.cs.umass.edu/~mimno/icml100.html

Map of NIH Grants

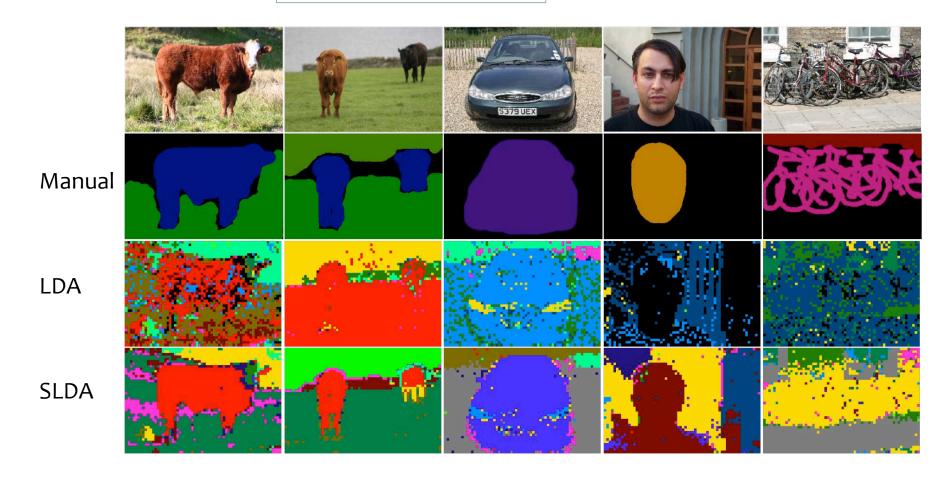


https://app.nihmaps.org/

Other Applications of Topic Models

Spacial LDA

(Wang & Grimson, 2007)



Other Applications of Topic Models

Word Sense Induction

Senses of drug (WSJ)

- 1. U.S., administration, federal, against, war, dealer
- 2. patient, people, problem, doctor, company, abuse
- 3. company, million, sale, maker, stock, inc.
- 4. administration, food, company, approval, FDA

(Brody & Lapata, 2009)

Senses of *drug* (BNC)

- 1. patient, treatment, effect, anti-inflammatory
- 2. alcohol, treatment, patient, therapy, addiction
- 3. patient, new, find, effect, choice, study
- 4. test, alcohol, patient, abuse, people, crime
- 5. trafficking, trafficker, charge, use, problem
- 6. abuse, against, problem, treatment, alcohol
- 7. people, wonder, find, prescription, drink, addict

• Selectional Preference 8. company, dealer, police, enforcement, patient

(Ritter et al., 2010)

Topic t	Arg1	Relations which assign highest probability to <i>t</i>	Arg2
18	The residue - The mixture - The reaction mixture - The solution - the mixture - the reaction mixture - the residue - The reaction - the solution - The filtrate - the reaction - The product - The crude product - The pellet - The organic layer - Thereto - This solution - The resulting solution - Next - The organic phase - The resulting mixture - C.)	was treated with, is treated with, was poured into, was extracted with, was purified by, was diluted with, was filtered through, is disolved in, is washed with	EtOAc - CH2Cl2 - H2O - CH.sub.2Cl.sub.2 - H.sub.2O - water - MeOH - NaHCO3 - Et2O - NHCl - CHCl.sub.3 - NHCl - drop- wise - CH2Cl.sub.2 - Celite - Et.sub.2O - Cl.sub.2 - NaOH - AcOEt - CH2Cl2 - the mixture - saturated NaHCO3 - SiO2 - H2O - N hydrochloric acid - NHCl - preparative HPLC - toO C

Outline

- Applications of Topic Modeling
- Review: Latent Dirichlet Allocation (LDA)
 - Beta-Bernoulli
 - 2. Dirichlet-Multinomial
 - 3. Dirichlet-Multinomial Mixture Model
 - 4. LDA
- Contrast of methods for Inference / Learning
 - Exact inference
 - EM
 - Monte Carlo EM
 - Gibbs sampler
 - Collapsed Gibbs sampler
- Extensions of LDA
 - Correlated topic models
 - Dynamic topic models
 - Polylingual topic models
 - Supervised LDA

Beta-Bernoulli Model

Beta Distribution

$$f(\phi|\alpha,\beta) = \frac{1}{B(\alpha,\beta)} x^{\alpha-1} (1-x)^{\beta-1}$$

$$\begin{array}{c} & & & \\$$

Beta-Bernoulli Model

Generative Process

```
\phi \sim \text{Beta}(\alpha, \beta) \qquad [draw \ distribution \ over \ words] For each word n \in \{1, \dots, N\} x_n \sim \text{Bernoulli}(\phi) \qquad [draw \ word]
```

Example corpus (heads/tails)

Н	Т	Т	Н	Н	Т	Т	Н	Н	Н
X ₁	X ₂	X ₃	X ₄	X ₅	X ₆	X ₇	X ₈	X ₉	X ₁₀

Dirichlet-Multinomial Model

Dirichlet Distribution

Dirichlet-Multinomial Model

Dirichlet Distribution

$$p(\vec{\phi}|\alpha) = \frac{1}{B(\alpha)} \prod_{k=1}^{K} \phi_k^{\alpha_k - 1} \quad \text{where } B(\alpha) = \frac{\prod_{k=1}^{K} \Gamma(\alpha_k)}{\Gamma(\sum_{k=1}^{K} \alpha_k)}$$

Dirichlet-Multinomial Model

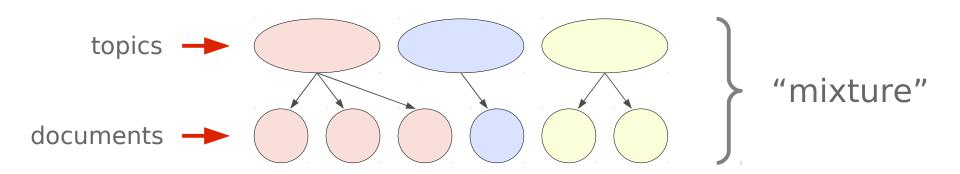
Generative Process

Example corpus

the	he	is	the	and	the	she	she	is	is
X ₁	X ₂	X ₃	X ₄	X ₅	X ₆	X ₇	x ₈	x ₉	X ₁₀

Dirichlet-Multinomial Mixture Model

Generative Process



Example corpus

the	he	is
X ₁₁	X ₁₂	X ₁₃

Document 1

the	and	the
X ₂₁	X ₂₂	X ₂₃

Document 2

she	she	is	is
X ₃₁	X ₃₂	X ₃₃	X ₃₄

Document 3

Dirichlet-Multinomial Mixture Model

Generative Process

```
\begin{array}{ll} \text{For each topic } k \in \{1, \dots, K\}: \\ \boldsymbol{\phi}_k \sim \text{Dir}(\boldsymbol{\beta}) & [\textit{draw distribution over words}] \\ \boldsymbol{\theta} \sim \text{Dir}(\boldsymbol{\alpha}) & [\textit{draw distribution over topics}] \\ \text{For each document } m \in \{1, \dots, M\} \\ z_m \sim \text{Mult}(1, \boldsymbol{\theta}) & [\textit{draw topic assignment}] \\ \text{For each word } n \in \{1, \dots, N_m\} \\ x_{mn} \sim \text{Mult}(1, \boldsymbol{\phi}_{z_m}) & [\textit{draw word}] \end{array}
```

Example corpus

the	he	is
X ₁₁	X ₁₂	X ₁₃

the	and	the
X ₂₁	X ₂₂	X ₂₃

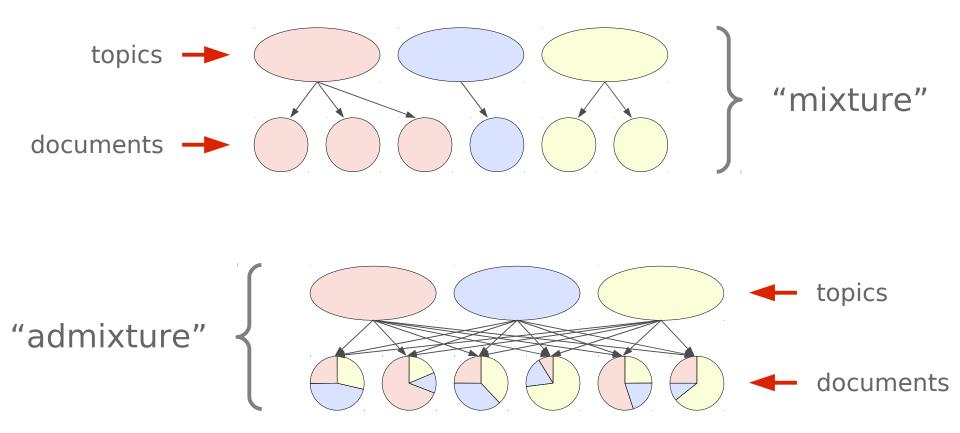
she	she	is	is
X ₃₁	X ₃₂	X ₃₃	X ₃₄

Document 1

Document 2

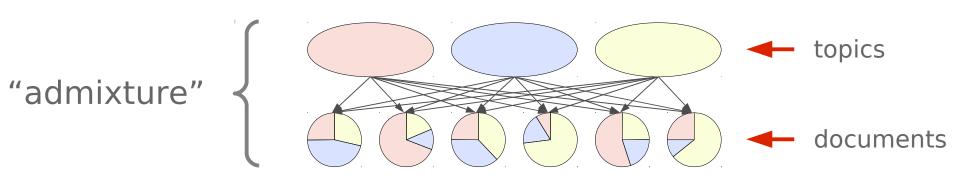
Document 3

Mixture vs. Admixture (LDA)



Diagrams from Wallach, JHU 2011, slides

Generative Process



Example corpus

the	he	is
X ₁₁	X ₁₂	X ₁₃

D^{\sim}	~ ·	m	Or	\+	4
Do	CL	1111	וכו	ΙL	ı

the	and	the	
X ₂₁	X ₂₂	X ₂₃	

Document 2

she	she	is	is
X ₃₁	X ₃₂	X ₃₃	X ₃₄

Document 3

Generative Process

```
For each topic k \in \{1, \dots, K\}:  \phi_k \sim \operatorname{Dir}(\boldsymbol{\beta}) \qquad [draw\ distribution\ over\ words]  For each document m \in \{1, \dots, M\}  \boldsymbol{\theta}_m \sim \operatorname{Dir}(\boldsymbol{\alpha}) \qquad [draw\ distribution\ over\ topics]  For each word n \in \{1, \dots, N_m\}  z_{mn} \sim \operatorname{Mult}(1, \boldsymbol{\theta}_m) \qquad [draw\ topic\ assignment]   x_{mn} \sim \boldsymbol{\phi}_{z_{mi}} \qquad [draw\ word]
```

Example corpus

the	he	is
X ₁₁	X ₁₂	X ₁₃

the	and	the
X ₂₁	X ₂₂	X ₂₃

she	she	is	is
X ₃₁	X ₃₂	X ₃₃	X ₃₄

Document 1

Document 2

Document 3

Plate Diagram

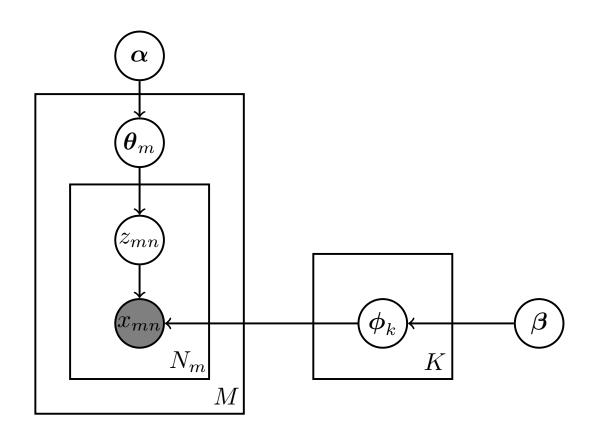
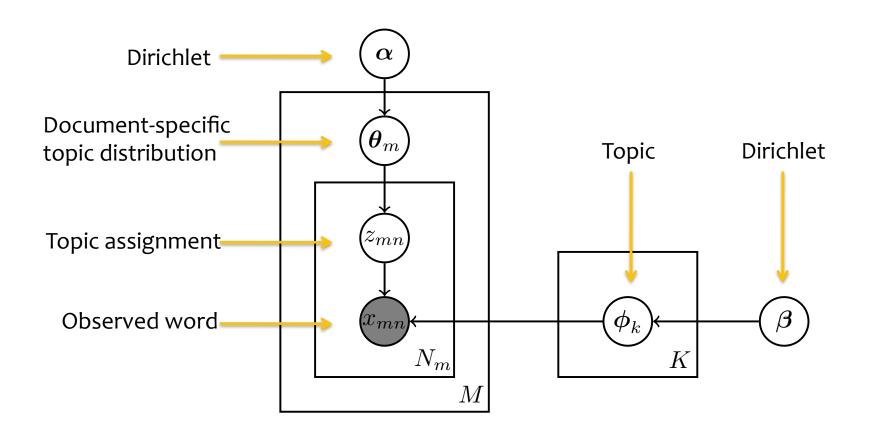
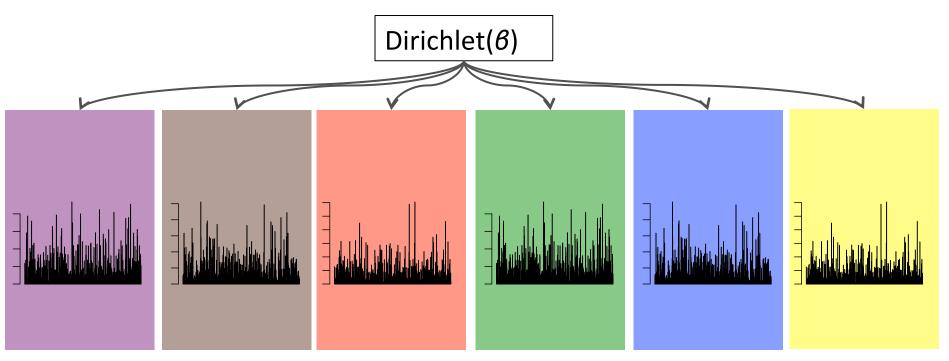
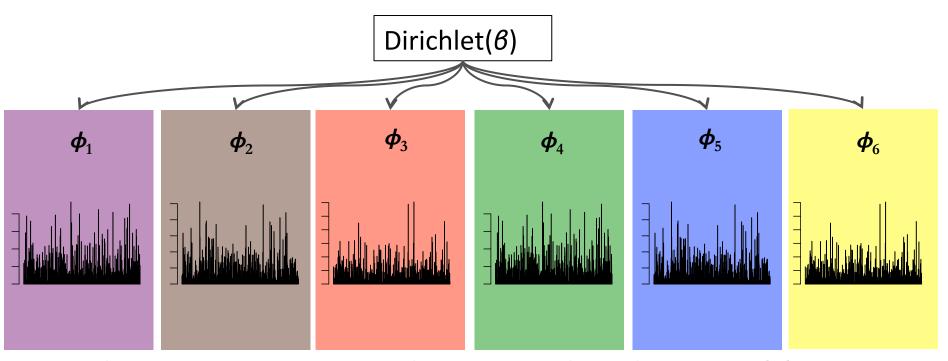


Plate Diagram

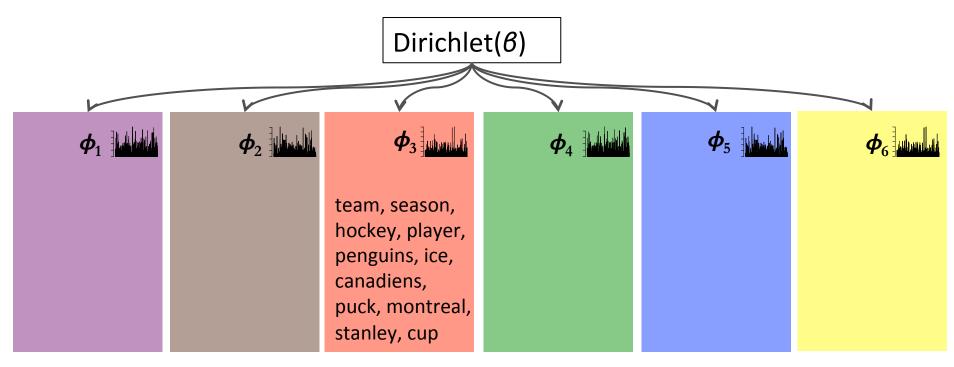




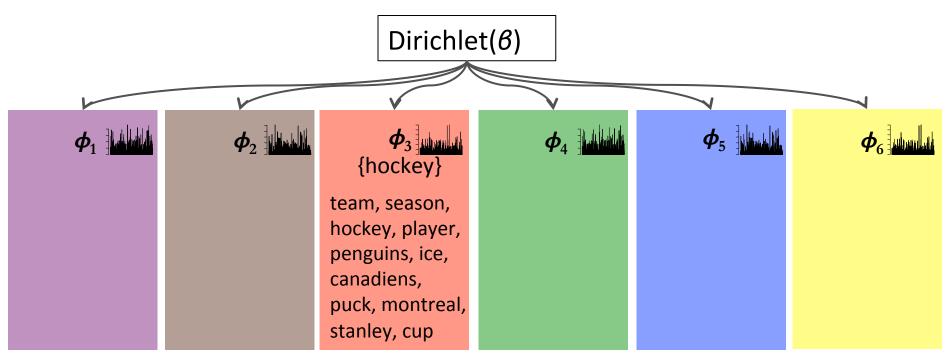
- The generative story begins with only a Dirichlet prior over the topics.
- Each **topic** is defined as a **Multinomial distribution** over the vocabulary, parameterized by $\phi_{\mathbf{k}}$



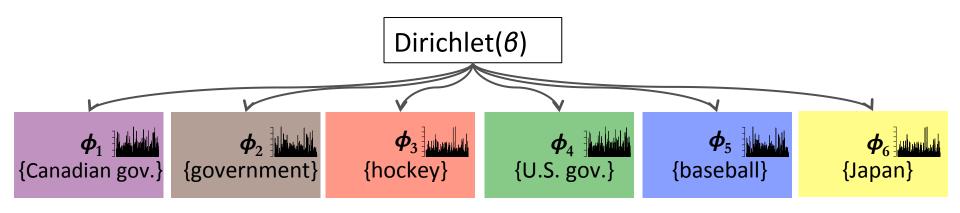
- The generative story begins with only a Dirichlet prior over the topics.
- Each **topic** is defined as a **Multinomial distribution** over the vocabulary, parameterized by ϕ_k



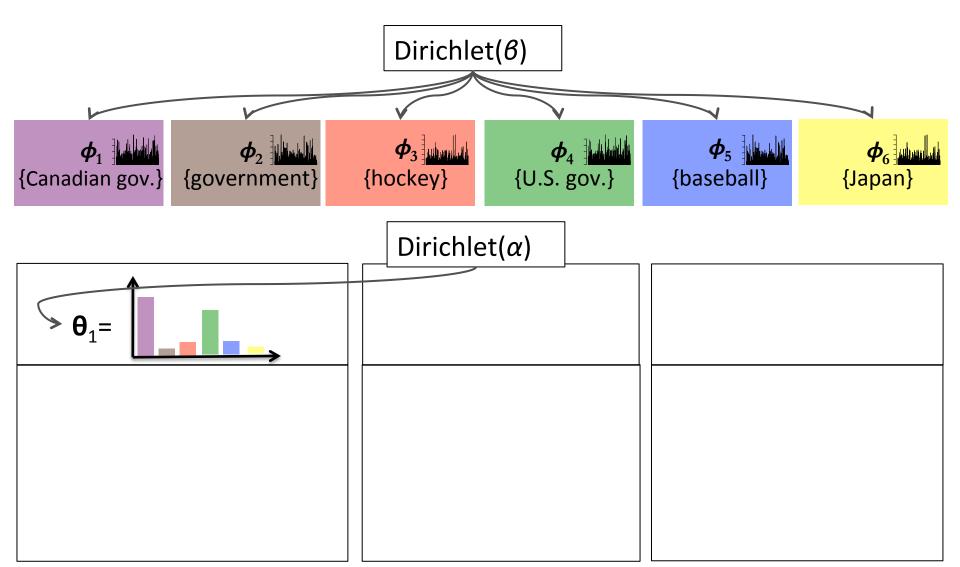
 A topic is visualized as its high probability words.

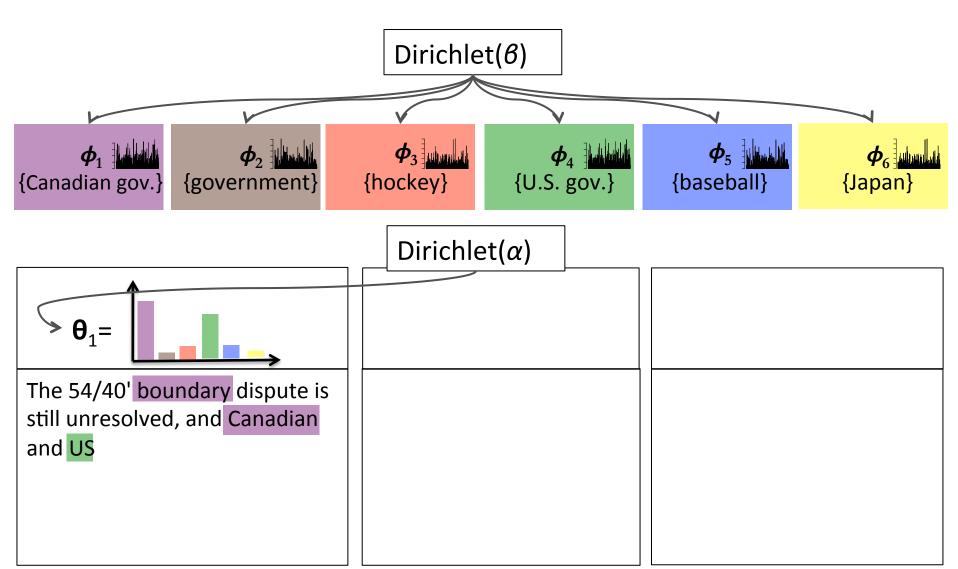


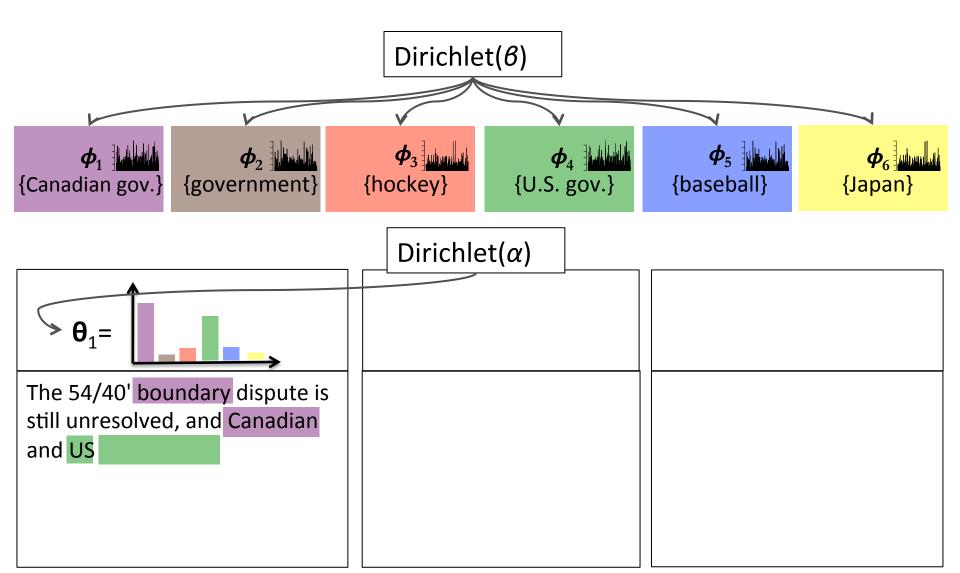
- A topic is visualized as its high probability words.
- A pedagogical label is used to identify the topic.

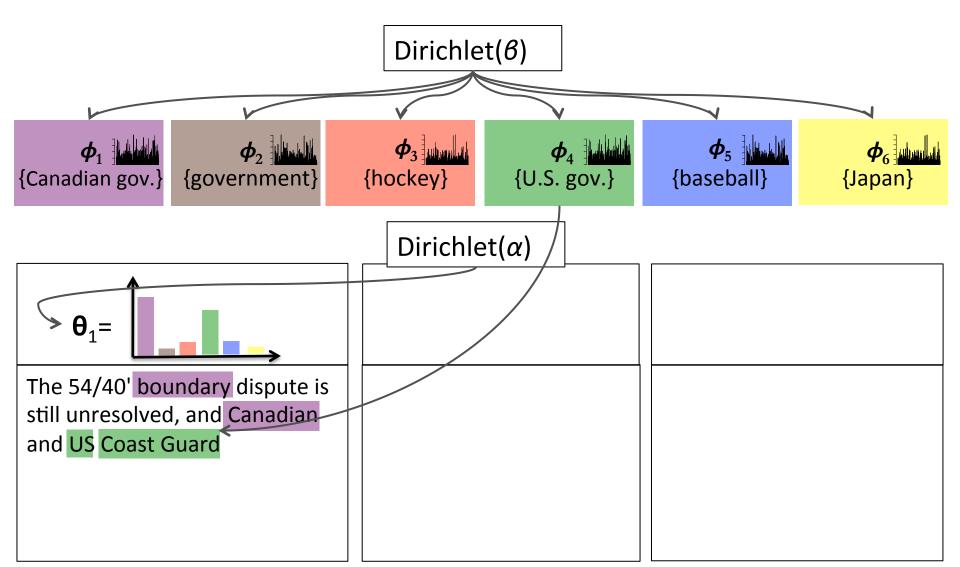


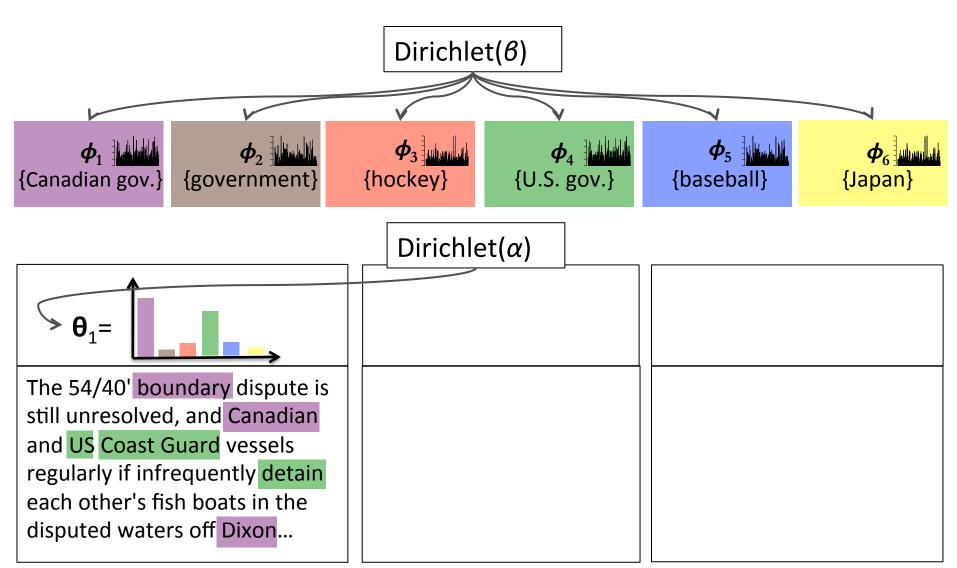
- A topic is visualized as its high probability words.
- A pedagogical label is used to identify the topic.

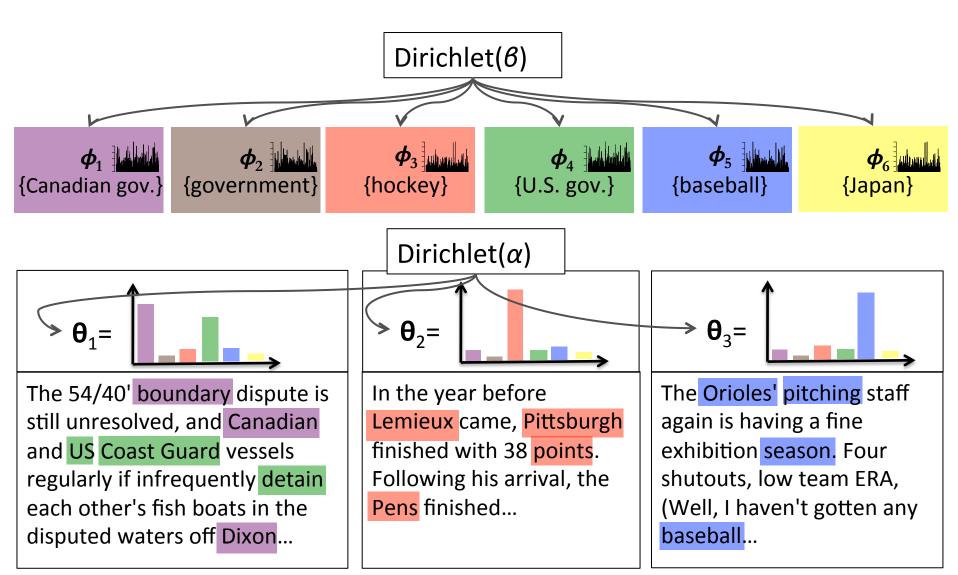


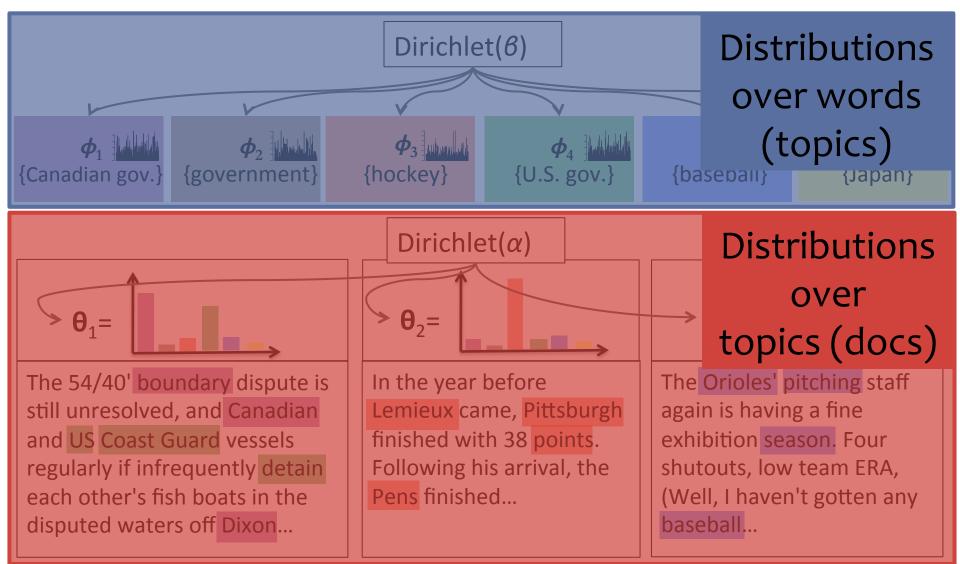


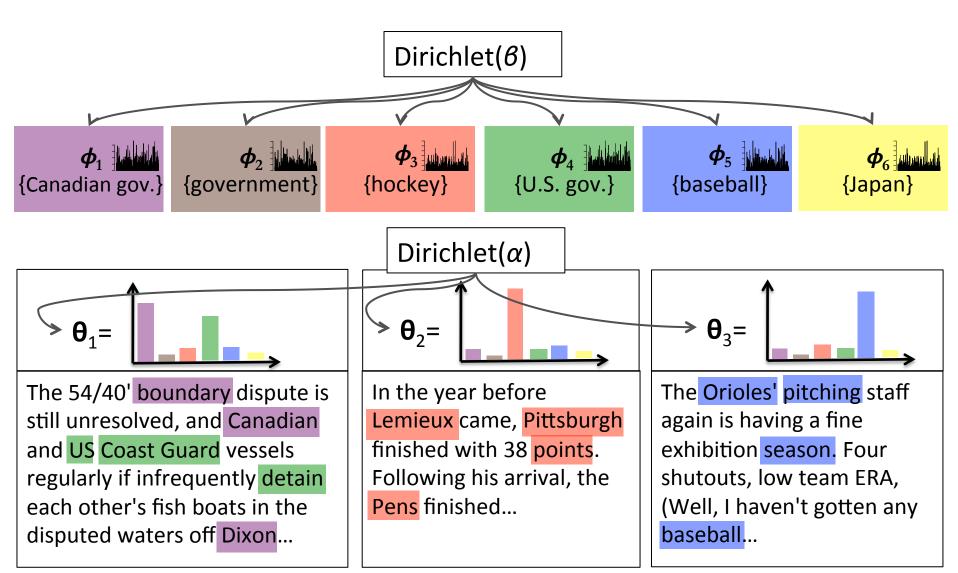












Inference and learning start with only the data

Dirichlet()

ф₁ =

 $\phi_2 =$

 $\phi_3 =$

 $\phi_4 =$

 $oldsymbol{\phi}_5$ =

φ₆ =

Dirichlet()

The 54/40' boundary dispute is still unresolved, and Canadian and US Coast Guard vessels regularly if infrequently detain each other's fish boats in the disputed waters off Dixon...

• **θ**₂=

In the year before Lemieux came, Pittsburgh finished with 38 points. Following his arrival, the Pens finished... **> θ**₃=

The Orioles' itching staff again is having a fine exhibition season. Four shutouts, low team ERA, (Well, I haven't gotten any baseball...

Questions:

 Is this a believable story for the generation of a corpus of documents?

Why might it work well anyway?

Latent Dirichlet Allocation

Why does LDA "work"?

- LDA trades off two goals.
 - Tor 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 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.

Latent Dirichlet Allocation

How does this relate to my other favorite model for capturing low-dimensional representations of a corpus?

- Builds on latent semantic analysis (Deerwester et al., 1990; Hofmann, 1999)
- It is a mixed-membership model (Erosheva, 2004).
- It relates to PCA and matrix factorization (Jakulin and Buntine, 2002)
- Was independently invented for genetics (Pritchard et al., 2000)

Outline

- Applications of Topic Modeling
- Review: Latent Dirichlet Allocation (LDA)
 - Beta-Bernoulli
 - 2. Dirichlet-Multinomial
 - 3. Dirichlet-Multinomial Mixture Model
 - 4. LDA
- Contrast of methods for Inference / Learning
 - Exact inference
 - EM
 - Monte Carlo EM
 - Gibbs sampler
 - Collapsed Gibbs sampler
- Extensions of LDA
 - Correlated topic models
 - Dynamic topic models
 - Polylingual topic models
 - Supervised LDA

Unsupervised Learning

Three learning paradigms:

1. Maximum likelihood

$$\arg \max_{\theta} p(X|\theta)$$

2. Maximum a posteriori (MAP)

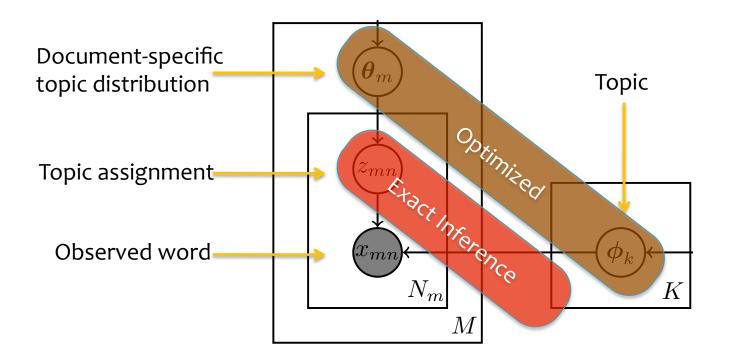
$$\arg \max_{\theta} p(\theta|X) \propto p(X|\theta)p(\theta)$$

3. Bayesian approach

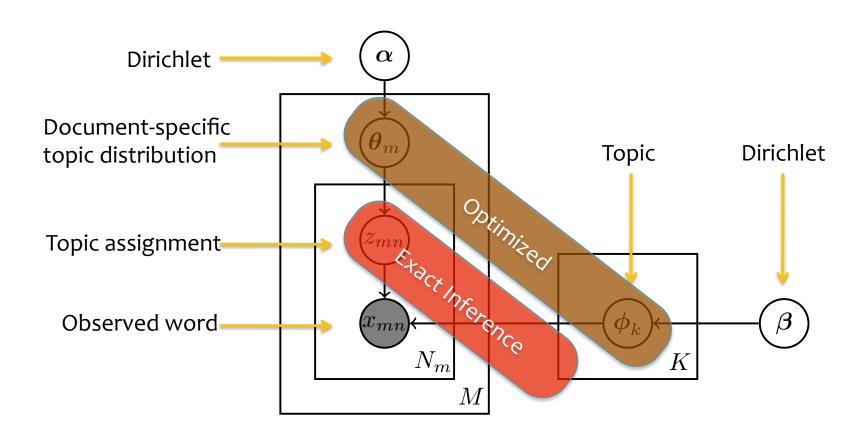
Estimate the posterior:

$$p(\theta|X) = \dots$$

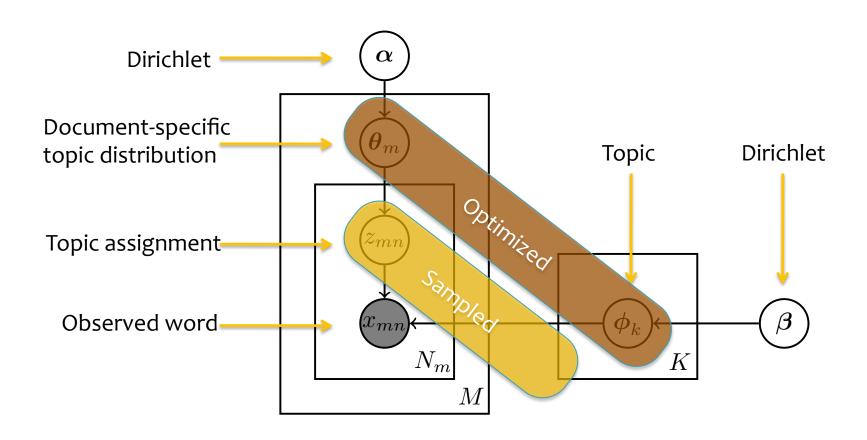
Standard EM (Maximum Likelihood)



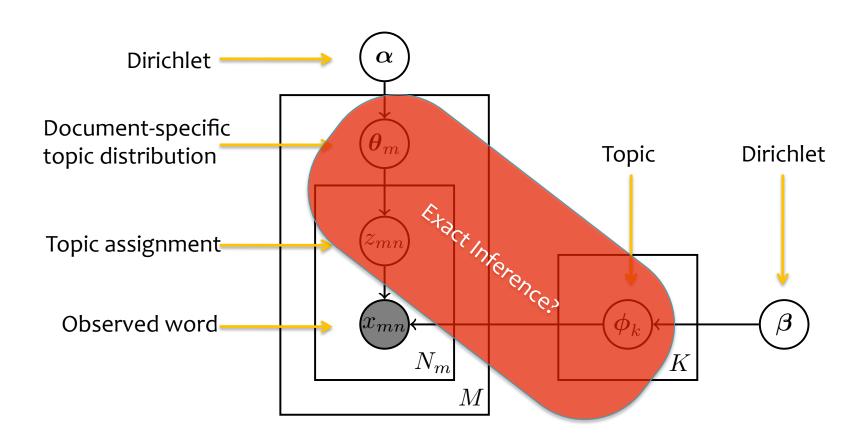
Standard EM (MAP)



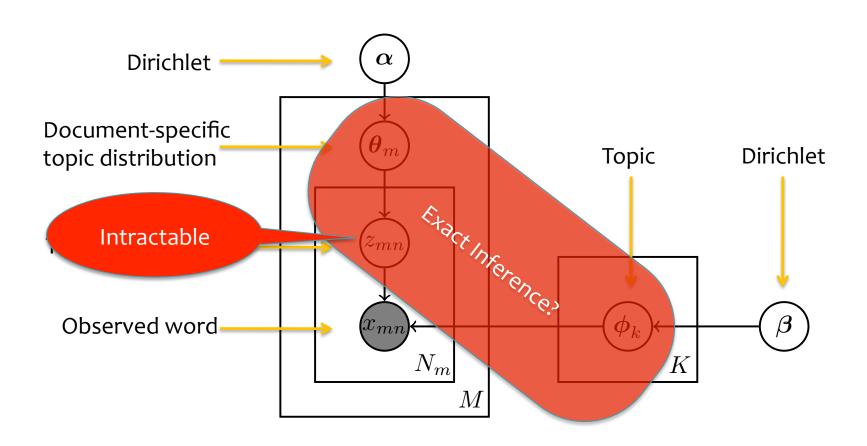
Monte Carlo EM



Bayesian Approach



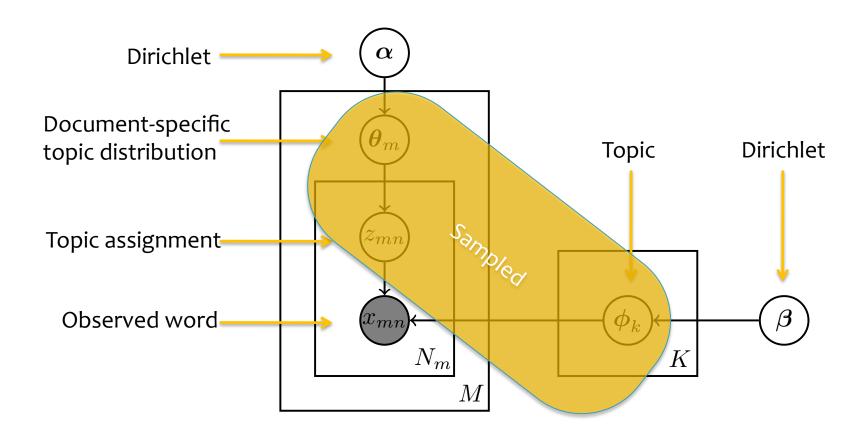
Bayesian Approach



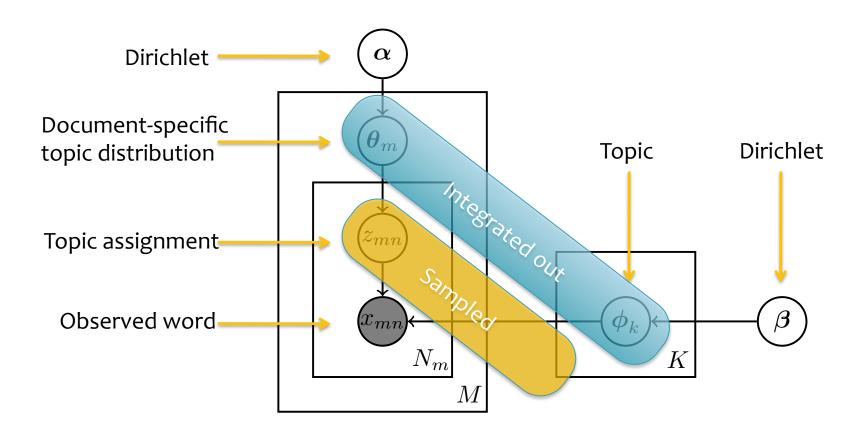
Exact Inference in LDA

- Exactly computing the posterior is intractable in LDA
 - Junction tree algorithm: exact inference in general graphical models
 - 1. "moralization" converts directed to undirected
 - 2. "triangulation" breaks 4-cycles by adding edges
 - 3. Cliques arranged into a junction tree
 - Time complexity is exponential in size of cliques
 - LDA cliques will be large (at least O(# topics)), so complexity is O(2^{# topics})
- Exact MAP inference in LDA is NP-hard for a large number of topics (Sontag & Roy, 2011)

Explicit Gibbs Sampler



Collapsed Gibbs Sampler



Sampling

Goal:

- Draw samples from the posterior $p(Z|X,\alpha,\beta)$
- Integrate out topics ϕ and document-specific distribution over topics θ

Algorithm:

- While not done...
 - For each document, *m*:
 - For each word, n:
 - » Resample a single topic assignment using the full conditionals for z_{mn}

Sampling

- What can we do with samples of z_{mn} ?
 - Mean of z_{mn}
 - Mode of z_{mn}
 - Estimate posterior over z_{mn}
 - Estimate of topics ϕ and document-specific distribution over topics θ

Gibbs Sampling for LDA

Full conditionals

$$p(z_i = k | Z^{-i}, X, \boldsymbol{\alpha}, \boldsymbol{\beta}) = \frac{n_{kt}^{-i} + \beta_t}{\sum_{v=1}^T n_{kv}^{-i} + \beta_v} \cdot \frac{n_{mk}^{-i} + \alpha_k}{\sum_{j=1}^K n_{mj}^{-i} + \alpha_j}$$
where t, m are given by i

 n_{kt} = # times topic k appears with type t n_{mk} = # times topic k appears in document m

Gibbs Sampling for LDA

Sketch of the derivation of the full conditionals

$$\begin{split} p(z_i = k|Z^{-i}, X, \pmb{\alpha}, \pmb{\beta}) &= \frac{p(X, Z|\pmb{\alpha}, \pmb{\beta})}{p(X, Z^{-i}|\pmb{\alpha}, \pmb{\beta})} \\ &\propto p(X, Z|\pmb{\alpha}, \pmb{\beta}) \\ &= p(X|Z, \pmb{\beta}) p(Z|\pmb{\alpha}) \\ &= \int_{\Phi} p(X|Z, \Phi) p(\Phi|\pmb{\beta}) \, d\Phi \, \int_{\Theta} p(Z|\Theta) p(\Theta|\pmb{\alpha}) \, d\Theta \\ &= \left(\prod_{k=1}^K \frac{B(\vec{n}_k + \pmb{\beta})}{B(\pmb{\beta})}\right) \left(\prod_{m=1}^M \frac{B(\vec{n}_m + \pmb{\alpha})}{B(\pmb{\alpha})}\right) \\ &= \frac{n_{kt}^{-i} + \beta_t}{\sum_{v=1}^T n_{kv}^{-i} + \beta_v} \cdot \frac{n_{mk}^{-i} + \alpha_k}{\sum_{j=1}^K n_{mj}^{-i} + \alpha_j} \\ &\qquad \text{where } t, m \text{ are given by } i \end{split}$$

Dirichlet-Multinomial Model

The Dirichlet is conjugate to the Multinomial

```
\phi \sim \operatorname{Dir}(\boldsymbol{\beta})   [draw distribution over words] For each word n \in \{1, \dots, N\}   [draw word]
```

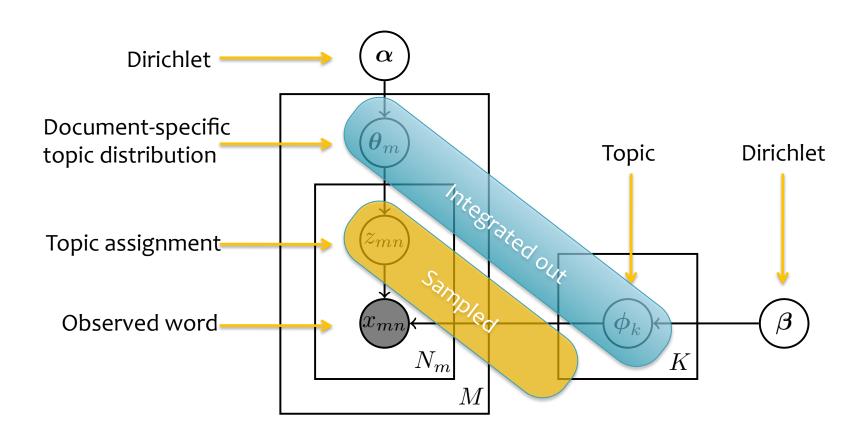
- The posterior of ϕ is $p(\phi|X) = \frac{p(X|\phi)p(\phi)}{P(X)}$
- Define the count vector n such that n_t denotes the number of times word t appeared
- Then the posterior is also a Dirichlet distribution: $p(\phi|X) \sim \text{Dir}(\beta + n)$

Dirichlet-Multinomial Model

Why conjugacy is so useful

$$\begin{split} p(X|\alpha) &= \int_{\phi} p(X|\vec{\phi}) p(\vec{\phi}|\alpha) \; d\phi \\ &= \int_{\phi} \left(\prod_{v=1}^{V} \phi_{v}^{n_{v}} \right) \left(\frac{1}{B(\alpha)} \prod_{v=1}^{V} \phi_{v}^{\alpha_{v}-1} \right) d\phi \\ &= \frac{1}{B(\alpha)} \int_{\phi} \prod_{v=1}^{V} \phi_{v}^{n_{v}+\alpha_{v}-1} \; d\phi \\ &= \frac{1}{B(\alpha)} \int_{\phi} \frac{B(\vec{n}+\alpha)}{B(\vec{n}+\alpha)} \prod_{v=1}^{V} \phi_{v}^{n_{v}+\alpha_{v}-1} \; d\phi \\ &= \frac{B(\vec{n}+\alpha)}{B(\alpha)} \int_{\phi} \underbrace{\frac{1}{B(\vec{n}+\alpha)} \prod_{v=1}^{V} \phi_{v}^{n_{v}+\alpha_{v}-1}}_{Dir(\vec{n}+\alpha)} \; d\phi \\ &= \frac{B(\vec{n}+\alpha)}{B(\alpha)} \end{split}$$

Collapsed Gibbs Sampler



Gibbs Sampling for LDA

Algorithm

```
zero all count variables, n_m^{(k)}, n_m, n_k^{(t)}, n_k

for all documents m \in [1, M] do

for all words n \in [1, N_m] in document m do

sample topic index z_{m,n} = k \sim \text{Mult}(1/K)

increment document—topic count: n_m^{(k)} += 1

increment topic—term count: n_k^{(t)} += 1

increment topic—term sum: n_k += 1
```

Gibbs Sampling for LDA

Algorithm

```
// Gibbs sampling over burn-in period and sampling period
 while not finished do
         for all documents m \in [1, M] do
                for all words n \in [1, N_m] in document m do
                        // for the current assignment of k to a term t for word w_{m,n}:
decrement counts and ...

// multinomial sampling acc. to Eq. ...

sample topic index \tilde{k} \sim p(z_i | \vec{z}_{\neg i}, \vec{w})

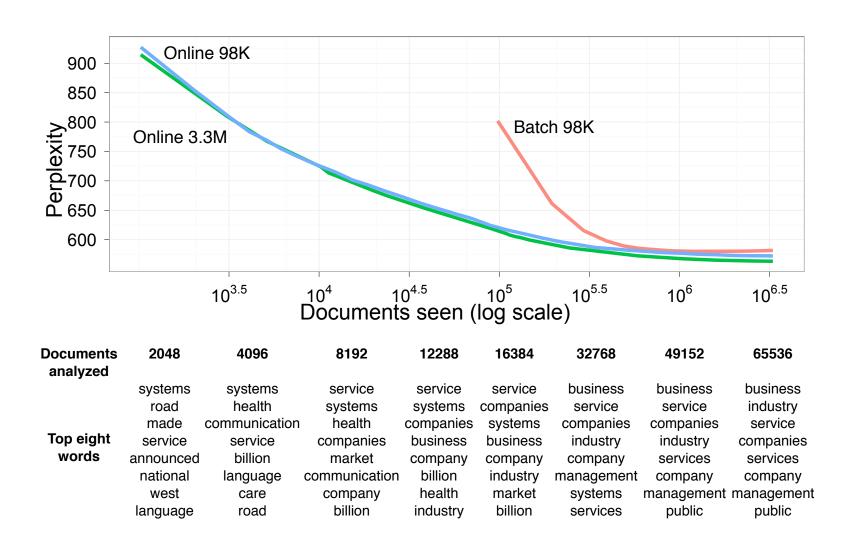
// for the new assignment of z_{m,n} to the term t for word w_{m,n}:

increment counts and sums: n_m^{(\tilde{k})} += 1; n_m += 1; n_{\tilde{k}}^{(t)} += 1; n_{\tilde{k}} += 1
                       // multinomial sampling acc. to Eq. 78 (decrements from previous step):
```

Why does Gibbs sampling work?

- Metropolis-Hastings
 - Markov chains
 - Stationary distribution
 - MH Algorithm
 - Constructs a Markov chain whose stationary distribution is the desired distribution
 - Proof that samples will be from desired distribution:
 - Sufficient conditions for constructing a markov chain with desired stationary distribution:
 - ergodicity
 - detailed balance (stronger, than what we need, but easier for the proof)
- Gibbs Sampling is a special case of Metropolis-Hastings
 - a special proposal distribution, which ensures the hastings ratio is always 1.0

Online Variational Inference for LDA



Outline

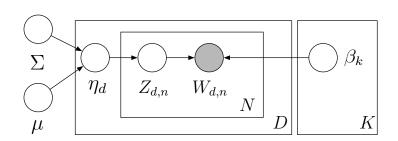
- Applications of Topic Modeling
- Review: Latent Dirichlet Allocation (LDA)
 - Beta-Bernoulli
 - 2. Dirichlet-Multinomial
 - 3. Dirichlet-Multinomial Mixture Model
 - 4. LDA
- Contrast of methods for Inference / Learning
 - Exact inference
 - EM
 - Monte Carlo EM
 - Gibbs sampler
 - Collapsed Gibbs sampler

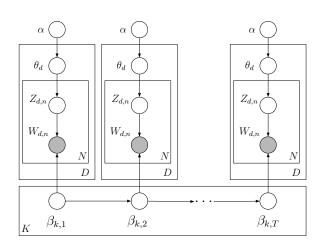
Extensions of LDA

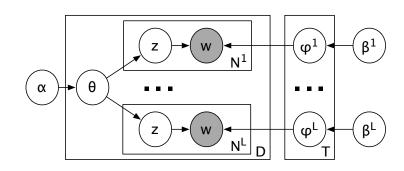
- Correlated topic models
- Dynamic topic models
- Polylingual topic models
- Supervised LDA

Extensions to the LDA Model

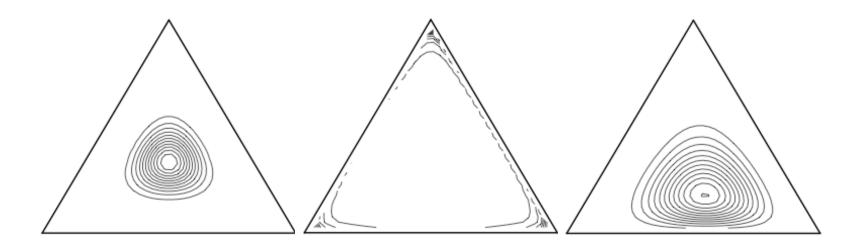
- Correlated topic models
 - Logistic normal prior over topic assignments
- Dynamic topic models
 - Learns topic changes over time
- Polylingual topic models
 - Learns topics aligned across multiple languages



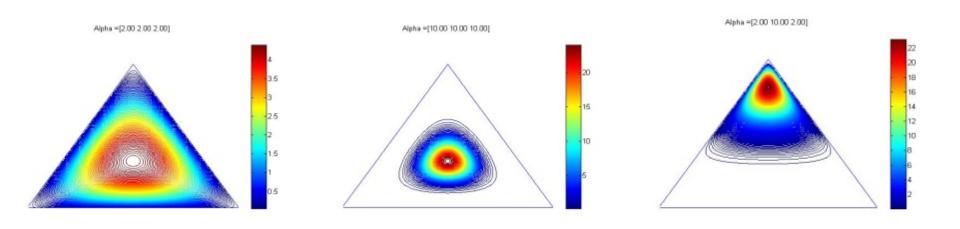




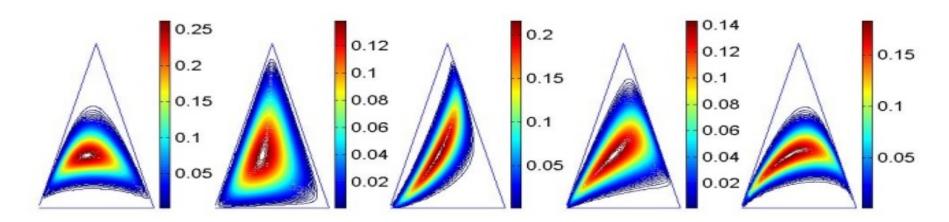
• • •



- The Dirichlet is a distribution on the simplex, positive vectors that sum to 1.
- It assumes that components are nearly independent.
- In real data, an article about *fossil fuels* is more likely to also be about *geology* than about *genetics*.



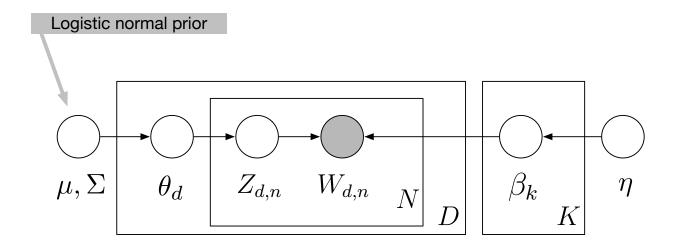
- The Dirichlet is a distribution on the simplex, positive vectors that sum to 1.
- It assumes that components are nearly independent.
- In real data, an article about *fossil fuels* is more likely to also be about *geology* than about *genetics*.



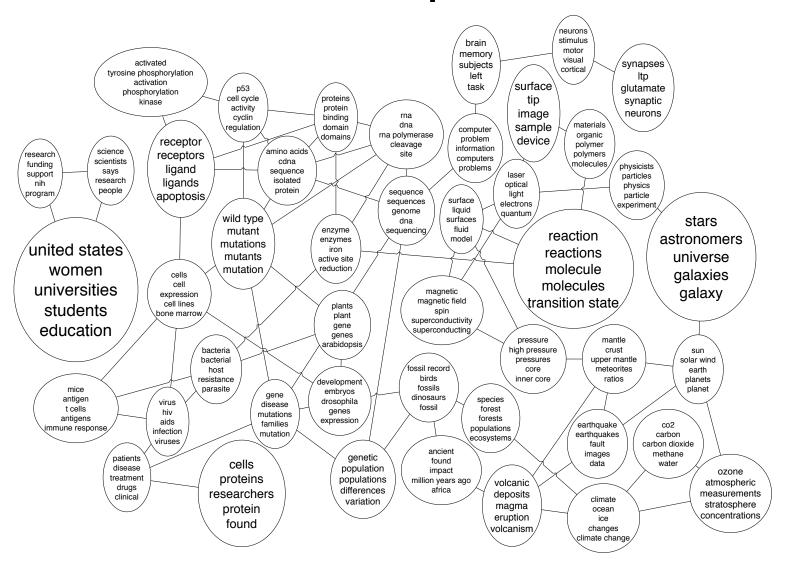
- The logistic normal is a distribution on the simplex that can model dependence between components (Aitchison, 1980).
- The log of the parameters of the multinomial are drawn from a multivariate Gaussian distribution,

$$X \sim \mathcal{N}_K(\mu, \Sigma)$$

$$\theta_i \propto \exp\{x_i\}.$$

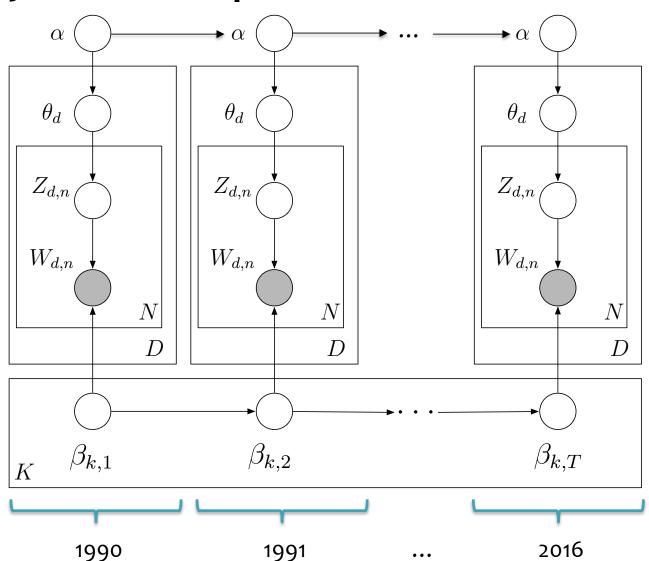


- Draw topic proportions from a logistic normal
- This allows topic occurrences to exhibit correlation.
- Provides a "map" of topics and how they are related
- Provides a better fit to text data, but computation is more complex



High-level idea:

- Divide the documents up by year
- Start with a separate topic model for each year
- Then add a dependence of each year on the previous one



1789 2009



Inaugural addresses



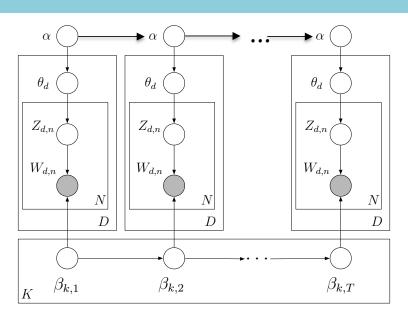
My fellow citizens: I stand here today humbled by the task before us, grateful for the trust you have bestowed, mindful of the sacrifices borne by our ancestors...

AMONG the vicissitudes incident to life no event could have filled me with greater anxieties than that of which the notification was transmitted by your order...

- LDA assumes that the order of documents does not matter.
- Not appropriate for sequential corpora (e.g., that span hundreds of years)
- Further, we may want to track how language changes over time.
- Dynamic topic models let the topics drift in a sequence.

Generative Story

- 1. Draw topics $\beta_t \mid \beta_{t-1} \sim \mathcal{N}(\beta_{t-1}, \sigma^2 I)$.
- 2. Draw $\alpha_t \mid \alpha_{t-1} \sim \mathcal{N}(\alpha_{t-1}, \delta^2 I)$.
- 3. For each document:
 - (a) Draw $\eta \sim \mathcal{N}(\alpha_t, a^2 I)$
 - (b) For each word:
 - i. Draw $Z \sim Mult(\pi(\eta))$.
 - ii. Draw $W_{t,d,n} \sim Mult(\pi(\beta_{t,z}))$.

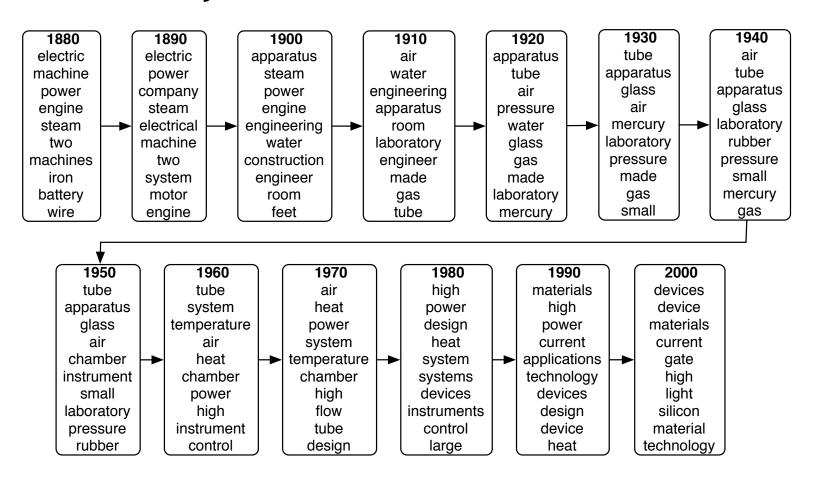


Logistic-normal priors

The pi function maps from the natural parameters to the mean parameters:

$$\pi(\beta_{k,t})_w = \frac{\exp(\beta_{k,t,w})}{\sum_w \exp(\beta_{k,t,w})}.$$

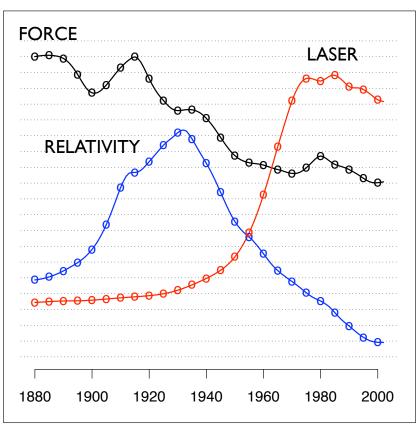
Top ten most likely words in a "drifting" topic shown at 10-year increments

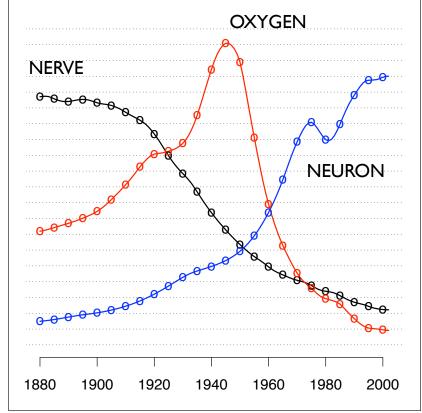


Posterior estimate of word frequency as a function of year for three words each in two separate topics:

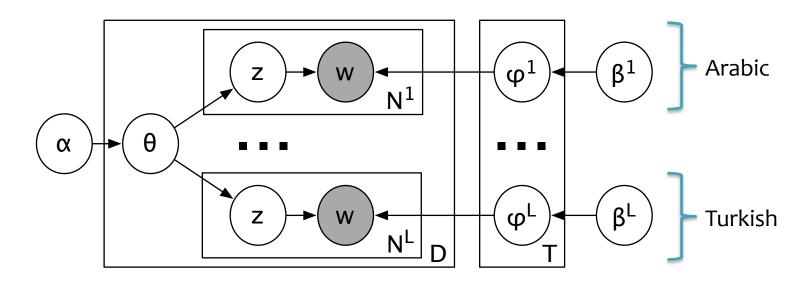
"Theoretical Physics"

"Neuroscience"





- Data Setting: Comparable versions of each document exist in multiple languages (e.g. the Wikipedia article for "Barak Obama" in twelve languages)
- **Model:** Very similar to LDA, except that the topic assignments, z, and words, w, are sampled separately for each language.



Topic 1 (twelve languages)

CY sadwrn blaned gallair at lloeren mytholeg

DE space nasa sojus flug mission

EL διαστημικό sts nasa αγγλ small

EN space mission launch satellite nasa spacecraft

فضایی ماموریت ناسا مدار فضانورد ماهواره FA

FI sojuz nasa apollo ensimmäinen space lento

FR spatiale mission orbite mars satellite spatial

החלל הארץ חלל כדור א תוכנית HE

IT spaziale missione programma space sojuz stazione

PL misja kosmicznej stacji misji space nasa

RU космический союз космического спутник станции

TR uzay soyuz ay uzaya salyut sovyetler

Topic 2 (twelve languages)

CY sbaen madrid el la josé sbaeneg

DE de spanischer spanischen spanien madrid la

EL ισπανίας ισπανία de ισπανός ντε μαδρίτη

EN de spanish spain la madrid y

ترین de اسپانیا اسپانیا یکوبا مادرید

FI espanja de espanjan madrid la real

FR espagnol espagne madrid espagnole juan y

ספרד ספרדית דה מדריד הספרדית קובה HE

IT de spagna spagnolo spagnola madrid el

PL de hiszpański hiszpanii la juan y

RU де мадрид испании испания испанский de

TR ispanya ispanyol madrid la küba real

Topic 3 (twelve languages)

```
CY bardd gerddi iaith beirdd fardd gymraeg

DE dichter schriftsteller literatur gedichte gedicht werk

EL ποιητής ποίηση ποιητή έργο ποιητές ποιήματα
```

EN poet poetry literature literary poems poem

شاعر شعر ادبیات فارسی ادبی آثار FA

FI runoilija kirjailija kirjallisuuden kirjoitti runo julkaisi

FR poète écrivain littérature poésie littéraire ses

משורר ספרות שירה סופר שירים המשורר

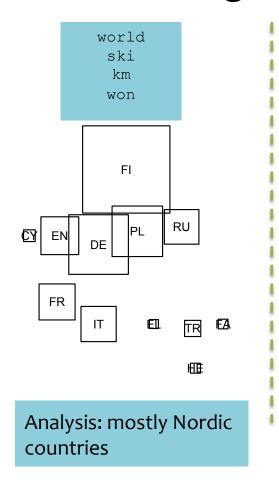
IT poeta letteratura poesia opere versi poema

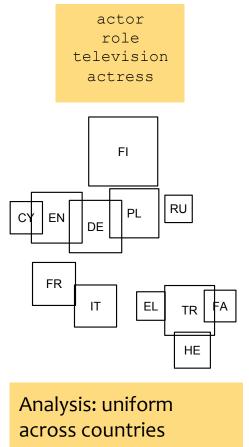
PL poeta literatury poezji pisarz in jego

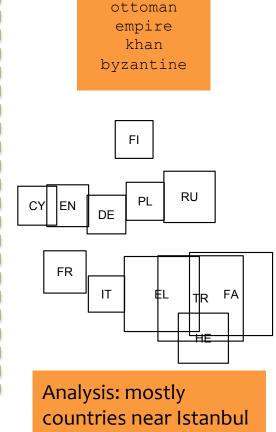
RU поэт его писатель литературы поэзии драматург

TR şair edebiyat şiir yazar edebiyatı adlı

Size of each square represents proportion of tokens assigned to the specified topic.



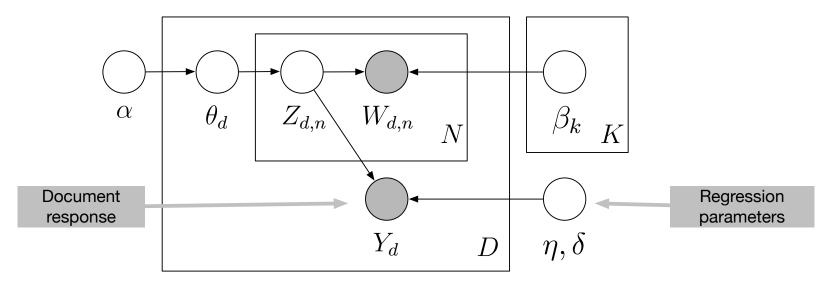




Supervised LDA

- LDA is an unsupervised model. How can we build a topic model that is good at the task we care about?
- Many data are paired with response variables.
 - User reviews paired with a number of stars
 - Web pages paired with a number of "likes"
 - Documents paired with links to other documents
 - Images paired with a category
- Supervised LDA are topic models of documents and responses.
 They are fit to find topics predictive of the response.

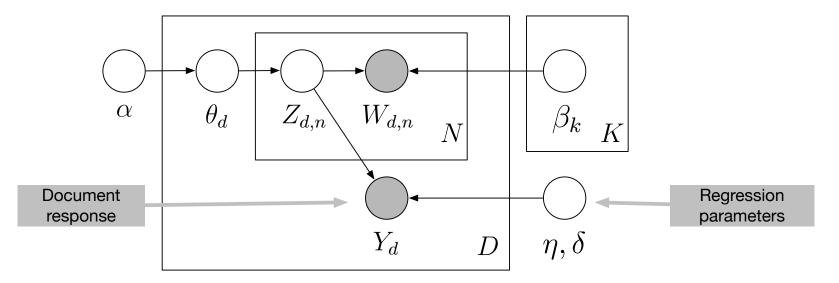
Supervised LDA



- ① Draw topic proportions $\theta \mid \alpha \sim Dir(\alpha)$.
- 2 For each word
 - Draw topic assignment $z_n | \theta \sim \text{Mult}(\theta)$.
 - Draw word $w_n | z_n, \beta_{1:K} \sim \text{Mult}(\beta_{z_n})$.
- 3 Draw response variable $y | z_{1:N}, \eta, \sigma^2 \sim N(\eta^T \overline{z}, \sigma^2)$, where

$$\overline{z} = (1/N) \sum_{n=1}^{N} z_n.$$

Supervised LDA



- Fit sLDA parameters to documents and responses. This gives: topics $\beta_{1:K}$ and coefficients $\eta_{1:K}$.
- Given a new document, predict its response using the expected value:

$$\mathrm{E}\left[Y|w_{1:N},\alpha,\beta_{1:K},\eta,\sigma^{2}\right]=\eta^{\top}\mathrm{E}\left[\bar{Z}|w_{1:N}\right]$$

This blends generative and discriminative modeling.

Summary: Approximate Inference

- Markov Chain Monte Carlo (MCMC)
 - Metropolis-Hastings, Gibbs sampling, Hamiltonion MCMC, slice sampling, etc.
- Variational inference
 - Minimizes KL(q||p) where q is a simpler graphical model than the original p
- Loopy Belief Propagation
 - Belief propagation applied to general (loopy) graphs
- Expectation propagation
 - Approximates belief states with moments of simpler distributions
- Spectral methods
 - Uses tensor decompositions (e.g. SVD)

Summary: Topic Modeling

The Task of Topic Modeling

- Topic modeling enables the analysis of large (possibly unannotated) corpora
- Applicable to more than just bags of words
- Extrinsic evaluations are often appropriate for these unsupervised methods

Constructing Models

- LDA is comprised of simple building blocks (Dirichlet, Multinomial)
- LDA itself can act as a building block for other models

Approximate Inference

 Many different approaches to inference (and learning) can be applied to the same model

What if we don't know the number of topics, K, ahead of time?

Next week: Bayesian Nonparametrics

- New modeling constructs:
 - Chinese Restaurant Process (Dirichlet Process)
 - Indian Buffet Process
- e.g. an infinite number of topics in a finite amount of space