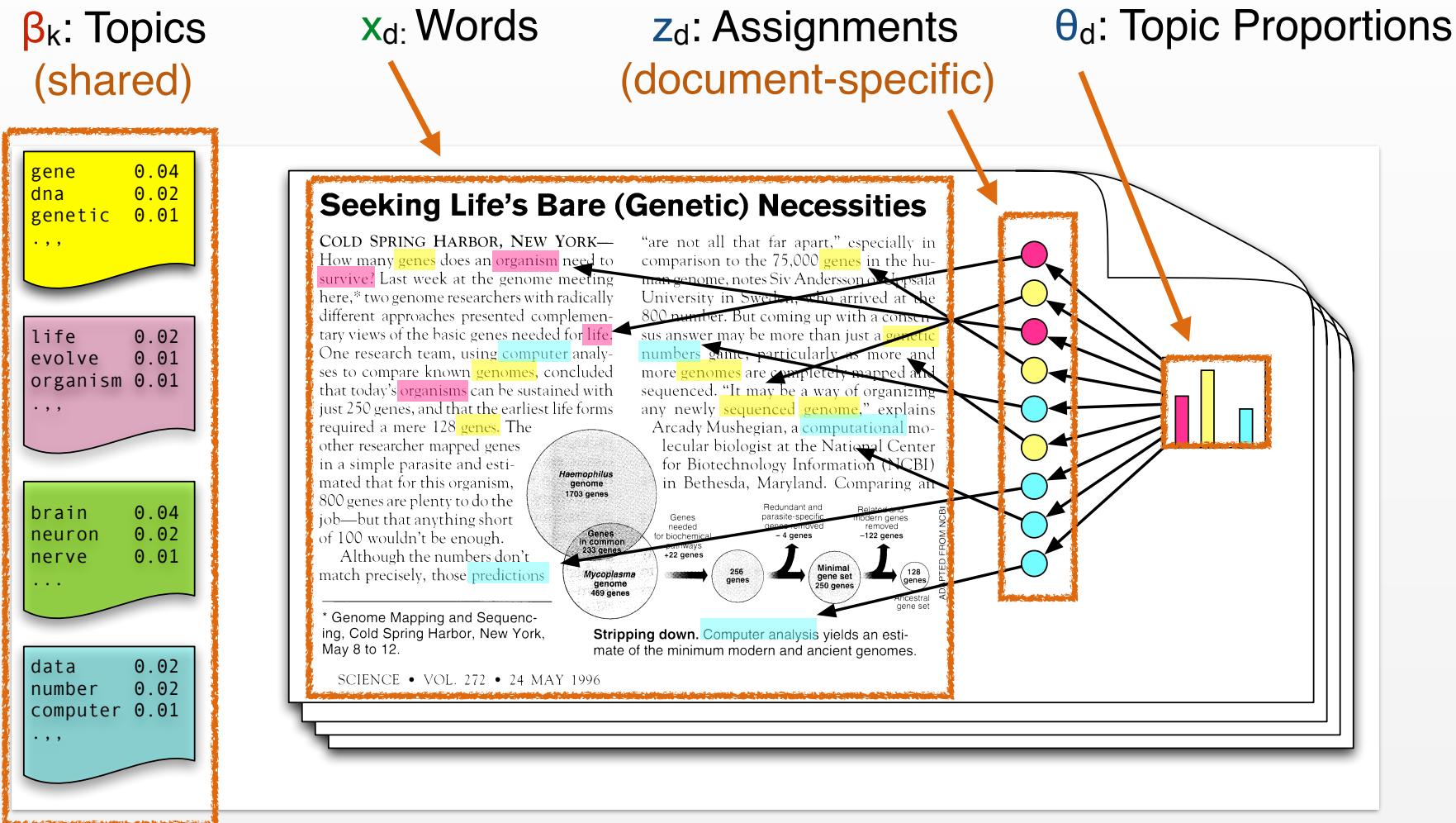




Latent Dirichlet Allocation

Topic Models with Dirichlet Priors

Review: Topic Modeling with PLSA/PLSI



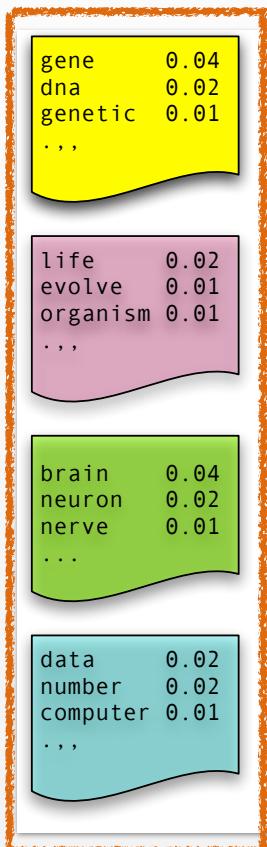
$$z_{d,n} \sim \text{Discrete}(\theta_d)$$

$$x_{d,n} | z_{d,n}=k \sim \text{Discrete}(\beta_k)$$

- Very flexible model.
- Overfits in practice.
- Doesn't give a sparse fit.

LDA: Add Dirichlet Priors

β_k : Topics
(shared)



x_d : Words

z_d : Assignments
(document-specific)

θ_d : Topic Proportions

Seeking Life's Bare (Genetic) Necessities

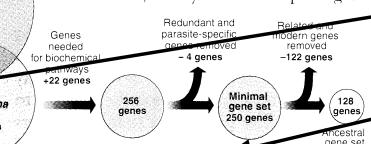
COLD SPRING HARBOR, NEW YORK—How many genes does an organism need to survive? Last week at the genome meeting here,* two genome researchers with radically different approaches presented complementary views of the basic genes needed for life. One research team, using computer analyses to compare known genomes, concluded that today's organisms can be sustained with just 250 genes, and that the earliest life forms required a mere 128 genes. The other researcher mapped genes in a simple parasite and estimated that for this organism, 800 genes are plenty to do the job—but that anything short of 100 wouldn't be enough.

Although the numbers don't match precisely, those predictions

* Genome Mapping and Sequencing, Cold Spring Harbor, New York, May 8 to 12.

SCIENCE • VOL. 272 • 24 MAY 1996

“are not all that far apart,” especially in comparison to the 75,000 genes in the human genome, notes Siv Andersson of Uppsala University in Sweden, who arrived at the 800 number. But coming up with a consensus answer may be more than just a genetic numbers game, particularly as more and more genomes are completely mapped and sequenced. “It may be a way of organizing any newly sequenced genome,” explains Arcady Mushegian, a computational molecular biologist at the National Center for Biotechnology Information (NCBI) in Bethesda, Maryland. Comparing all



Stripping down. Computer analysis yields an estimate of the minimum modern and ancient genomes.

$$z_{d,n} \sim \text{Discrete}(\theta_d)$$

$$x_{d,n} | z_{d,n}=k \sim \text{Discrete}(\beta_k)$$

$$\theta_d \sim \text{Dirichlet}(\alpha)$$

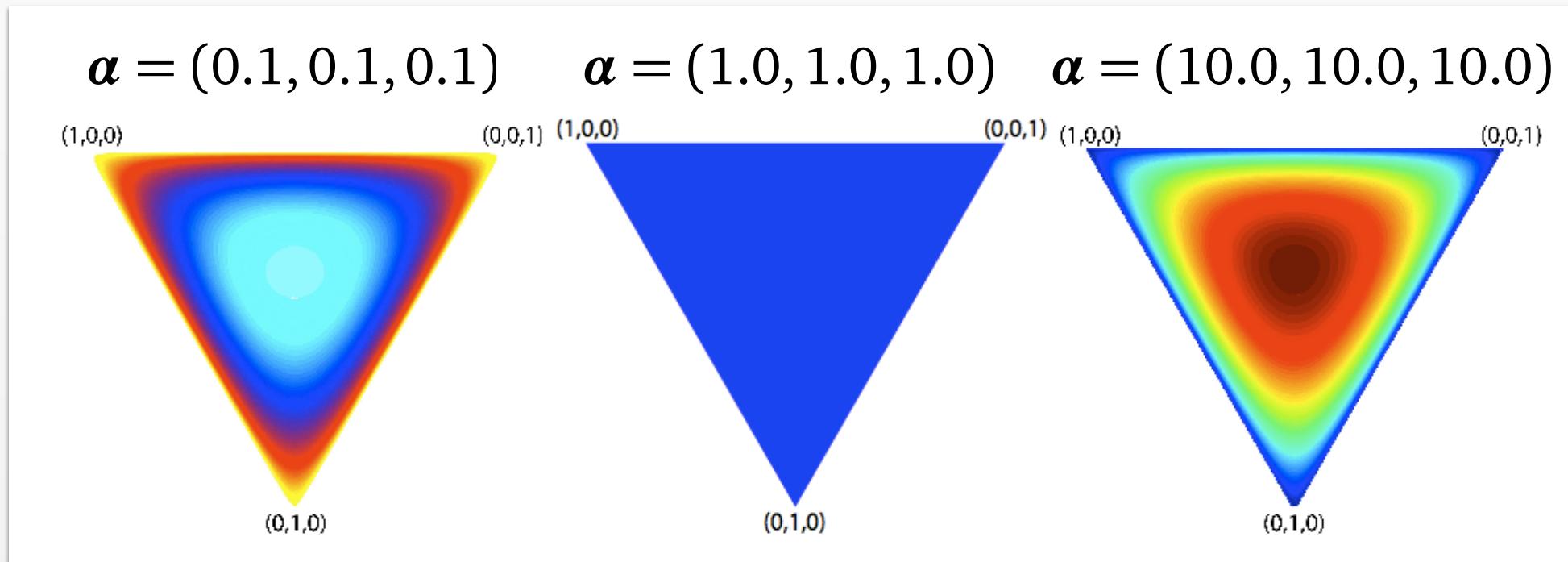
$$\beta_k \sim \text{Dirichlet}(\eta_k)$$

Review: Dirichlet Distribution

$\theta \sim \text{Dirichlet}(\alpha)$

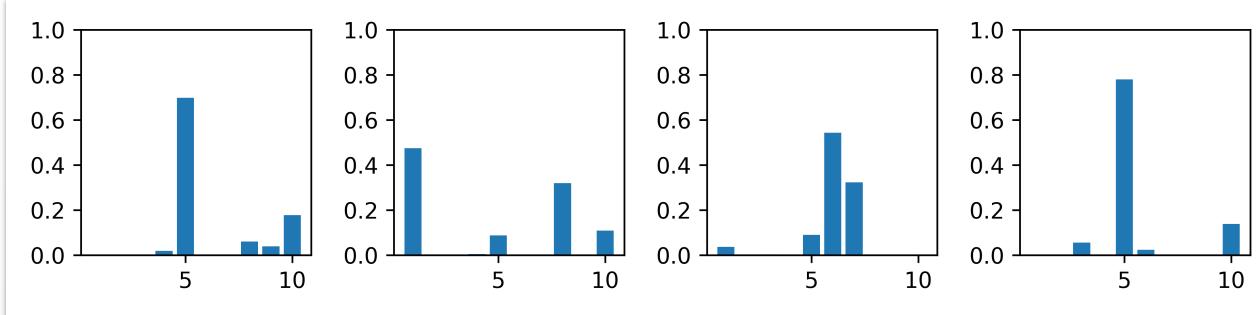
$$p(\theta) = \frac{1}{B(\alpha)} \prod_{k=1}^K \theta_k^{\alpha_k - 1}$$

$$B(\alpha) := \frac{\prod_{k=1}^K \Gamma(\alpha_k)}{\Gamma\left(\sum_{k=1}^K \alpha_k\right)}$$

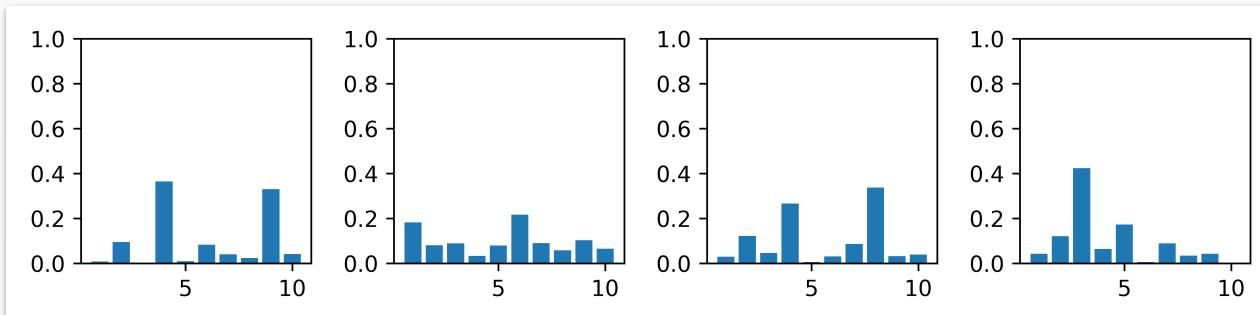


Review: Dirichlet Distribution

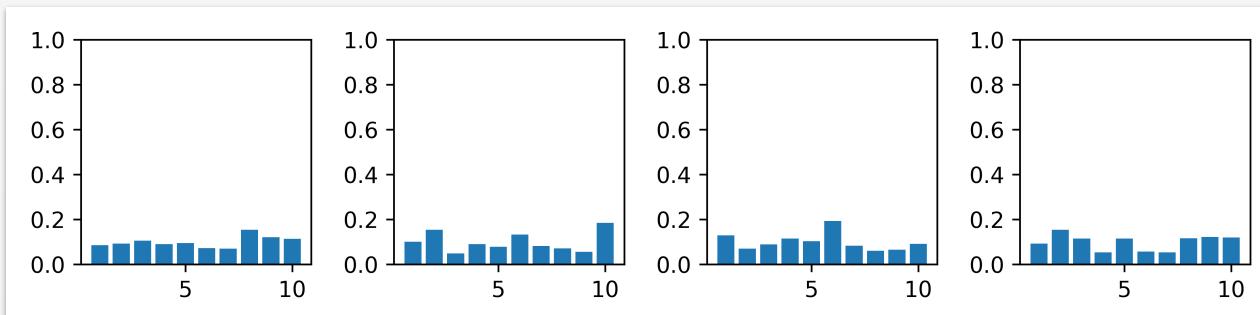
$\alpha_k = 0.1$



$\alpha_k = 1.0$



$\alpha_k = 10.0$



LDA: $\alpha_k = 0.001$ – Enforces Sparsity of Topic Weights θ_d

Maximum likelihood estimation

X : a dataset with N points:

$$X = \{x_n\}_{n=1}^N$$

$x_n \in \mathcal{X}$: the space where x_n takes values.
Typically $\mathcal{X} = \mathbb{R}^D$

Model Assumption:

$$x_n \sim P_\theta$$

A distribution parametrized by θ (unknown, but not random) defined on the sample space the space \mathcal{X}

Maximum likelihood estimate of θ .

$$\begin{aligned}\theta^{ML} &= \operatorname{argmax}_\theta p(X | \theta) \\ &= \operatorname{argmax}_\theta \log p(X | \theta)\end{aligned}$$

Does not depend on the prior distribution, $p(\theta)$.

Maximum likelihood is a frequentist approach. It assumes that the parameters are unknown, but fixed, in the sense that they do not come from a distribution. In the frequentist world the prior distribution does not exist.

Bayesian inference

X : a dataset with N points:

$$X = \{x_n\}_{n=1}^N$$

Posterior distribution

$$\begin{aligned} p(\theta | X, \alpha) &= \frac{p(X | \theta, \alpha)p(\theta | \alpha)}{P(X | \alpha)} \\ &= \frac{p(X | \theta)p(\theta | \alpha)}{P(X | \alpha)} \end{aligned}$$

Likelihood Prior
Intractable marginal

As the number of data points increase the contribution of the likelihood term increases and the prior term reduces

Model Assumption:

$$x_n \sim P_\theta$$

The data distribution, parametrized by θ

$$\theta \sim P_\alpha$$

The most general goal of Bayesian statistics is to learn the entire posterior distribution, not just an estimate of the parameter. However, this is often very difficult due to the intractable denominator. Have to rely on approximate inference.

$$P(X | \theta) = \int_{\theta} p(X | \theta)p(\theta | \alpha)d\theta$$

The prior distribution on the parameter θ . In the Bayesian setting, the parameter itself is a random variable having a “prior” distribution, P_α . The prior distribution parameter, α , is called the hyper-parameter. Which is typically assumed to be known or searched via grid-search.

Maximum a posteriori (MAP) estimation

X : a dataset with N points:

$$X = \{x_n\}_{n=1}^N$$

MAP Estimate

$$\begin{aligned}\theta^{MAP} &= \operatorname{argmax}_\theta p(\theta | X, \alpha) \\ &= \operatorname{argmax}_\theta \frac{p(X | \theta, \alpha)p(\theta | \alpha)}{P(X | \alpha)} \\ &= \operatorname{argmax}_\theta \frac{p(X | \theta)p(\theta | \alpha)}{P(X | \alpha)} \\ &= \operatorname{argmax}_\theta p(X | \theta) \times p(\theta | \alpha)\end{aligned}$$

↑
Posterior
Likelihood Prior

Not a function
of θ

As the number of data points increase the contribution of the likelihood term increases and the prior term reduces

Model Assumption:

$$x_n \sim P_\theta$$

The data distribution, parametrized by θ

$$\theta \sim P_\alpha$$

The prior distribution on the parameter θ . In the Bayesian setting, the parameter itself is a random variable having a “prior” distribution, P_α . The prior distribution parameter, α , is called the hyper-parameter. Which is typically assumed to be known or searched via grid-search.

$$\begin{aligned}\theta^{ML} &= \operatorname{argmax}_\theta p(X | \theta) \\ &= \operatorname{argmax}_\theta \log p(X | \theta)\end{aligned}$$

Estimating Model Parameters

Question: How can we estimate β_k and θ_d ?

1. MAP with Expectation Maximization
2. Variational Inference
(high level)
3. Gibbs Sampling
(not in this module)

MAP generalization to Discrete distribution

Likelihood

$$p(Z|\theta) = \prod_{n=1}^N p(z_n|\theta) \\ = \prod_{n=1}^N \prod_{k=1}^K \theta_k^{I[x_n=k]}$$

$$p(\theta|Z, \alpha) \propto p(Z|\theta)p(\theta|\alpha)$$

Posterior

$$\propto \left(\prod_{n=1}^N \prod_{k=1}^K \theta_k^{I[z_n=k]} \right) \frac{\prod_{k=1}^K \theta_k^{\alpha_k-1}}{B(\alpha)} \quad \text{Not a function of } \mu \\ \propto \prod_{k=1}^K \theta_k^{\left[\sum_{n=1}^N I[z_n=k] \right]} \prod_{k=1}^K \theta_k^{\alpha_k-1} \\ \propto \prod_{k=1}^K \theta_k^{\text{TC}_k + \alpha_k - 1} \quad \text{TC}_k = \sum_{n=1}^N I[z_n = k] \\ = \text{Dirichlet}(\theta | \alpha_1, \alpha_2, \dots, \alpha_K)$$

$$\text{MAP Estimate} \quad \theta_k^{MAP} = \frac{\text{TC}_k + \alpha_k - 1}{N + \sum_{k=1}^K \alpha_k - K}$$

Prior

$$p(\theta|\alpha) = \frac{\prod_{k=1}^K \theta_k^{\alpha_k-1}}{B(\alpha)}$$

Simpler example: single document, observable topics

$$Z = \{z_n\}_{n=1}^N \quad z_n \in \{1, 2, \dots, K\}$$

$$z_n \sim \text{Discrete}(\theta) \quad \sum_{k=1}^K \theta_k = 1$$

$$p(z_n|\theta) = \begin{cases} \theta_1 & \text{if } z_n = 1 \\ \theta_2 & \text{if } z_n = 2 \\ \vdots & \\ \theta_K & \text{if } z_n = K \end{cases}$$

$$p(z_n|\theta) = \prod_{k=1}^K \theta_k^{I[z_n=k]}$$

$$\theta \sim \text{Dirichlet}(\alpha)$$

$$p(\theta|\alpha) = \frac{\prod_{k=1}^K \theta_k^{\alpha_k-1}}{B(\alpha)}$$

$$\alpha = [\alpha_1, \alpha_2, \dots, \alpha_K] \quad \alpha_k > 0$$

MAP generalization to Discrete distribution

MAP Estimate

$$\theta_{dk}^{MAP} = \frac{\text{TC}_{dk} + \alpha_k - 1}{N_d + \sum_{k=1}^K \alpha_k - K}$$

$$\text{TC}_{dk} = \sum_{n=1}^{N_d} I[z_{dn} = k]$$

Count of topic k in the d^{th} document

Extension to multiple document.

$$Z = \{z_{dn}\} \quad z_{dn} \in \{1, 2, \dots, K\}$$

$$z_{dn} \sim \text{Discrete}(\theta_d) \quad \sum_{k=1}^K \theta_{dk} = 1$$

$$p(z_{dn} | \theta) = \prod_{k=1}^K \theta_{dk}^{I[z_{dn}=k]}$$
$$\theta = \{\theta_d\}$$

$$\theta_d \sim \text{Dirichlet}(\alpha)$$

$$p(\theta_d | \alpha) = \frac{\prod_{k=1}^K \theta_{dk}^{\alpha_k - 1}}{B(\alpha)}$$

$$\alpha = [\alpha_1, \alpha_2, \dots, \alpha_K] \quad \alpha_k > 0$$

MAP generalization to Discrete distribution

MAP Estimate

$$\beta_{kv}^{MAP} = \frac{WC_{kv} + \eta_{kv} - 1}{\sum_{d=1}^D TC_{dk} + \sum_{v=1}^V \eta_{kv} - V}$$

$$WC_{kv} = \sum_{d=1}^D \sum_{n=1}^{N_d} I[z_{dn} = k] I[x_{dn} = v]$$

$$\theta_{dk}^{MAP} = \frac{TC_{dk} + \alpha_k - 1}{N_d + \sum_{k=1}^K \alpha_k - K}$$

$$TC_{dk} = \sum_{n=1}^{N_d} I[z_{dn} = k]$$

Count of topic k in the d^{th} document

Estimating word probabilities
multiple document.

Simplifying assumption: Topics
are observed

$$Z = \{z_{dn}\} \quad z_{dn} \in \{1, 2, \dots, K\}$$

$$X = \{x_{dn}\} \quad x_{dn} \in \{1, 2, \dots, V\}$$

$$z_{dn} \sim \text{Discrete}(\theta_d) \quad \sum_{k=1}^K \theta_{dk} = 1$$

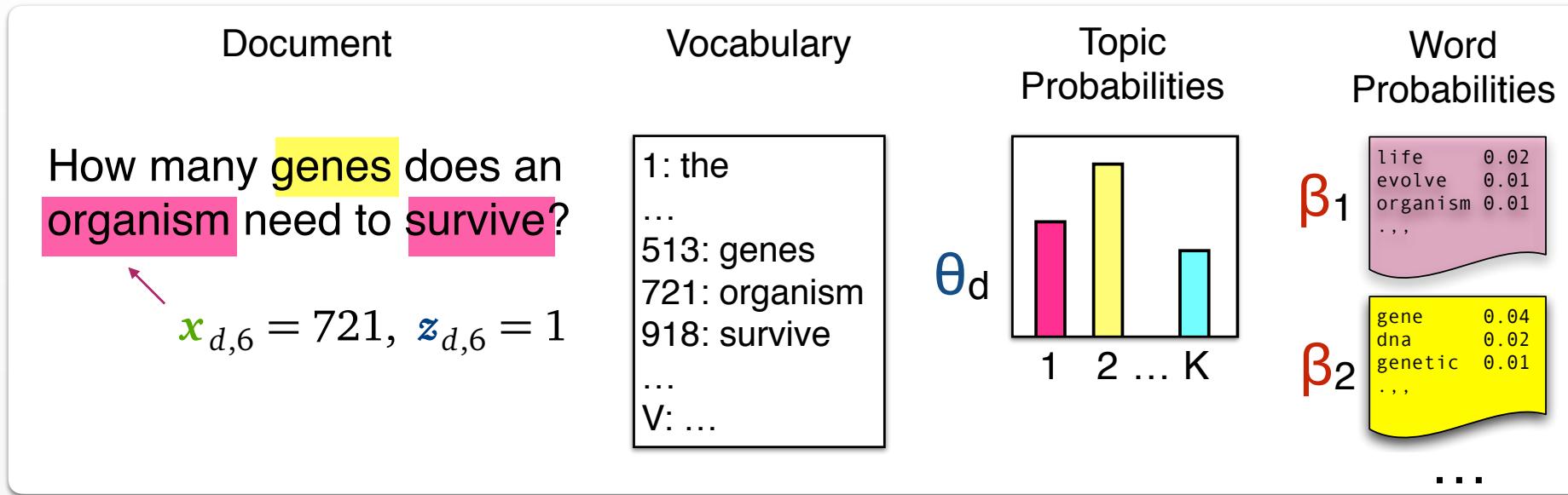
$$\theta_d \sim \text{Dirichlet}(\alpha) \quad \theta = \{\theta_d\}$$

$$x_{dn} | z_{dn} = k \sim \text{Discrete}(\beta_k)$$

$$\beta_k \sim \text{Dirichlet}(\eta_k) \quad \sum_{k=1}^K \beta_{kv} = 1$$

$$\beta = \{\beta_k\} \quad \eta = \{\eta_k\}$$

Estimating the Parameters



Maximum Likelihood: $\max_{\theta, \beta} \log p(\mathbf{x} | \theta, \beta)$

Maximum a Posteriori: $\max_{\theta, \beta} \log p(\theta, \beta | \mathbf{x}, \alpha, \eta) \quad \eta = \{\eta_k\}$

$$\mathbf{z}_{d,n} \sim \text{Discrete}(\boldsymbol{\theta}_d)$$

$$\mathbf{x}_{d,n} | \mathbf{z}_{d,n} = k \sim \text{Discrete}(\boldsymbol{\beta}_k)$$

$$\boldsymbol{\theta}_d \sim \text{Dirichlet}(\boldsymbol{\alpha})$$

$$\boldsymbol{\beta}_k \sim \text{Dirichlet}(\boldsymbol{\eta}_k)$$

MAP estimation for LDA with EM

Generative Model

$$\boldsymbol{\theta}_d \sim \text{Dirichlet}(\boldsymbol{\alpha})$$

$$\boldsymbol{\beta}_k \sim \text{Dirichlet}(\boldsymbol{\eta}_k)$$

$$\mathbf{z}_{d,n} \sim \text{Discrete}(\boldsymbol{\theta}_d)$$

$$\mathbf{x}_{d,n} \mid \mathbf{z}_{d,n} = k \sim \text{Discrete}(\boldsymbol{\beta}_k)$$

$$\text{TC}_{dk} = \sum_{n=1}^{N_d} I[z_{dn} = k]$$

$$\text{TC}_{dk} = \sum_{n=1}^{N_d} \phi_{dnk}$$

$$\text{WC}_{kv} = \sum_{d=1}^D \sum_{n=1}^{N_d} I[z_{dn} = k] I[x_{dn} = v]$$

$$\text{WC}_{kv} = \sum_{d=1}^D \sum_{n=1}^{N_d} \phi_{dnk} I[x_{dn} = v]$$

E-step: Update assignments

$$\begin{aligned} \phi_{d,n,k} &= p(\mathbf{z}_{d,n} = k \mid \mathbf{x}_{d,n} = v, \boldsymbol{\beta}, \boldsymbol{\theta}_d) \\ &= \frac{\boldsymbol{\theta}_{d,k} \left(\sum_{v=1}^V \boldsymbol{\beta}_{k,v} I[\mathbf{x}_{d,n} = v] \right)}{\sum_{l=1}^K \boldsymbol{\theta}_{d,l} \left(\sum_{v=1}^V \boldsymbol{\beta}_{l,v} I[\mathbf{x}_{d,n} = v] \right)} \end{aligned}$$

M-step: Update parameters

$$\beta_{kv}^{MAP} = \frac{\text{WC}_{kv} + \eta_{kv} - 1}{\sum_{d=1}^D \text{TC}_{dk} + \sum_{v=1}^V \eta_{kv} - V}$$

$$\theta_{dk}^{MAP} = \frac{\text{TC}_{dk} + \alpha_k - 1}{N_d + \sum_{k=1}^K \alpha_k - K}$$

(not used in practice; requires $\alpha_k > 1$ and $\eta_{kv} > 1$)

Estimating the posterior

$$\begin{aligned} p(Z, \theta, \beta | X, \alpha, \eta) &= \frac{p(Z, \theta, \beta, X | \alpha, \eta)}{p(X | \alpha, \eta)} \\ &= \frac{p(X | Z, \beta)p(Z | \theta)p(\theta | \alpha)p(\beta | \eta)}{p(X | \alpha, \eta)} \\ &= \frac{p(\theta | \alpha)p(\beta | \eta)\prod_d \prod_n p(x_{dn} | z_{dn}, \beta)p(z_{dn} | \theta_d)}{\int_{\theta} \int_{\beta} \sum_Z p(Z, \theta, \beta, X | \alpha, \eta) d\theta d\beta} \end{aligned}$$

Intractable due to complex
dependencies

Observed variables:
 X

Unobserved variables:
 Z, θ, β
 $\theta = \{\theta_d\}$
 $\beta = \{\beta_k\}$

Hyper-parameters:
 α, η
 $\eta = \{\eta_k\}$



Variational Inference approach

Approximate with a simple probability distribution,
removing dependencies

Mean-field assumption

$$q(Z, \theta, \beta | X, \alpha, \eta) = q(Z | X, \alpha, \beta) \times q(\theta | X, \alpha, \beta) \times q(\beta | X, \alpha, \eta)$$

$$q(Z, \theta, \beta | \phi, \gamma, \lambda) = q(Z | \phi(X)) \times q(\theta | \gamma(X)) \times q(\beta | \lambda(X))$$

Variational Parameters to be
learnt from the observed data

$$z_{dn} \sim \text{Discrete}(\phi_{dn1}, \phi_{dn2} \dots \phi_{dnK})$$

$$\theta_d \sim \text{Dirichlet}(\gamma_{d1}, \gamma_{d2} \dots \gamma_{dK})$$

$$\beta_k \sim \text{Dirichlet}(\lambda_{k1}, \lambda_{k2} \dots \lambda_{kV})$$

Observed variables:

$$X$$

Unobserved variables:

$$Z, \theta, \beta$$

$$\theta = \{\theta_d\}$$

$$\beta = \{\beta_k\}$$

Hyper-parameters:

$$\alpha, \eta$$

$$\eta = \{\eta_k\}$$

$$\phi^*, \gamma^*, \lambda^* = \underset{\phi, \gamma, \lambda}{\operatorname{argmin}} \text{KL}(q(Z, \theta, \beta | \phi, \gamma, \lambda) \| p(Z, \theta, \beta | X, \alpha, \beta))$$

Variational Expectation Maximization (high-level)

Variational E-step: Update ϕ

$$\phi_{d,n,k} \propto \exp \left(\mathbb{E}_q \left[\log \theta_{d,k} + \sum_{v=1}^V I[\mathbf{x}_{d,n} = v] \log \beta_{k,v} \right] \right)$$

(won't derive this – but can be computed in closed form)

$$\theta_d \sim \text{Dirichlet}(\gamma_{d1}, \gamma_{d2} \dots \gamma_{dK}) \quad \beta_k \sim \text{Dirichlet}(\lambda_{k1}, \lambda_{k2} \dots \lambda_{kV})$$

EM (non-variational)

$$\begin{aligned} \phi_{d,n,k} &= p(z_{d,n}=k \mid \mathbf{x}_{d,n} = v, \beta, \theta_d) \\ &= \frac{\theta_{d,k} \left(\sum_{v=1}^V \beta_{k,v} I[\mathbf{x}_{d,n} = v] \right)}{\sum_{l=1}^K \theta_{d,l} \left(\sum_{v=1}^V \beta_{l,v} I[\mathbf{x}_{d,n} = v] \right)} \end{aligned}$$

Variational M-step: Update γ and λ

$$\gamma_{d,k} = \alpha_k + \text{TC}_{dk}$$

$$\lambda_{k,v} = \eta_{k,v} + \text{WC}_{kv}$$

(analogous to MAP estimation – need to know this)

$$\phi_{d,n,k} \propto \theta_{d,k} \left(\sum_{v=1}^V \beta_{k,v} I[\mathbf{x}_{d,n} = v] \right)$$

Take log on both sides and compare to the variational EM

Variational Expectation Maximization (high-level)

The output of Variational EM is

$$\phi^*, \gamma^*, \lambda^* = \operatorname{argmin}_{\phi, \gamma, \lambda} \text{KL}(q(Z, \theta, \beta | \phi, \gamma, \lambda) \| p(Z, \theta, \beta | X, \alpha, \beta))$$

This gives distribution over
 $z_{dn}, \theta_d, \beta_k$

$$z_{dn} \sim \text{Discrete}(\phi_{dn1}^*, \phi_{dn2}^* \dots \phi_{dnK}^*)$$

$$\theta_d \sim \text{Dirichlet}(\gamma_{d1}^*, \gamma_{d2}^* \dots \gamma_{dK}^*)$$

$$\beta_k \sim \text{Dirichlet}(\lambda_{k1}^*, \lambda_{k2}^* \dots \lambda_{kV}^*)$$

Point estimates can be obtained as

$$\theta_{dk}^* = \mathbf{E}[\theta_{dk} | \gamma_d^*] = \frac{\gamma_{dk}}{\sum_{k=1}^K \gamma_{dk}}$$

$$\beta_{kv}^* = \mathbf{E}[\beta_{kv} | \lambda_k^*] = \frac{\lambda_{kv}}{\sum_{v=1}^V \lambda_{kv}}$$

$$z_{dn}^* = \operatorname{argmax}_k \phi_{dnk}^*$$

The point estimated of the topic (if needed) is obtained as the topic having the highest probability

Point estimates of θ_d, β_k are needed to obtain a representation of a document and a topic, respectively.

EM vs. Variational EM

EM $\theta, \beta = \underset{\theta, \beta}{\operatorname{argmax}} \log p(\mathbf{x} | \theta, \beta)$

E-step: $\phi_{d,n,k} \propto \theta_{d,k} \left(\sum_{v=1}^V \beta_{k,v} I[\mathbf{x}_{d,n} = v] \right)$

M-step: $\theta_{d,k} = \frac{\text{TC}_{dk}}{N_d}$ $\beta_{k,v} = \frac{\text{WC}_{kv}}{\sum_{d=1}^D \text{TC}_{dk}}$

Variational EM $\phi^*, \gamma^*, \lambda^* = \underset{\phi, \gamma, \lambda}{\operatorname{argmin}} \text{KL}\left(q(Z, \theta, \beta | \phi, \gamma, \lambda) \| p(Z, \theta, \beta | X, \alpha, \beta)\right)$

E-step: $\phi_{d,n,k} \propto \exp \left(\mathbb{E}_q \left[\log \theta_{d,k} + \sum_{v=1}^V I[\mathbf{x}_{d,n} = v] \log \beta_{k,v} \right] \right)$

M-step: $\gamma_{d,k} = \alpha_k + \text{TC}_{dk}$ $\lambda_{k,v} = \eta_{k,v} + \text{WC}_{kv}$

EM vs. Variational EM

Commonalities: Both compute sufficient statistics

$$\phi_{d,n,k}$$

Probability that word n in document d belongs to topic k

$$TC_{dk}$$

Number of words in document d that belong to topic k

$$WC_{kv}$$

Number of times word v appears in topic k
(across all documents in corpus)

Differences: Point estimates vs Distributions

EM: Computes most likely values for parameters

$$\theta_{d,k}$$

Fraction of words in document d for topic k

$$\beta_{k,v}$$

Fraction of words in topic k for vocabulary entry v

Variational EM: Estimate *Posterior* over Parameters

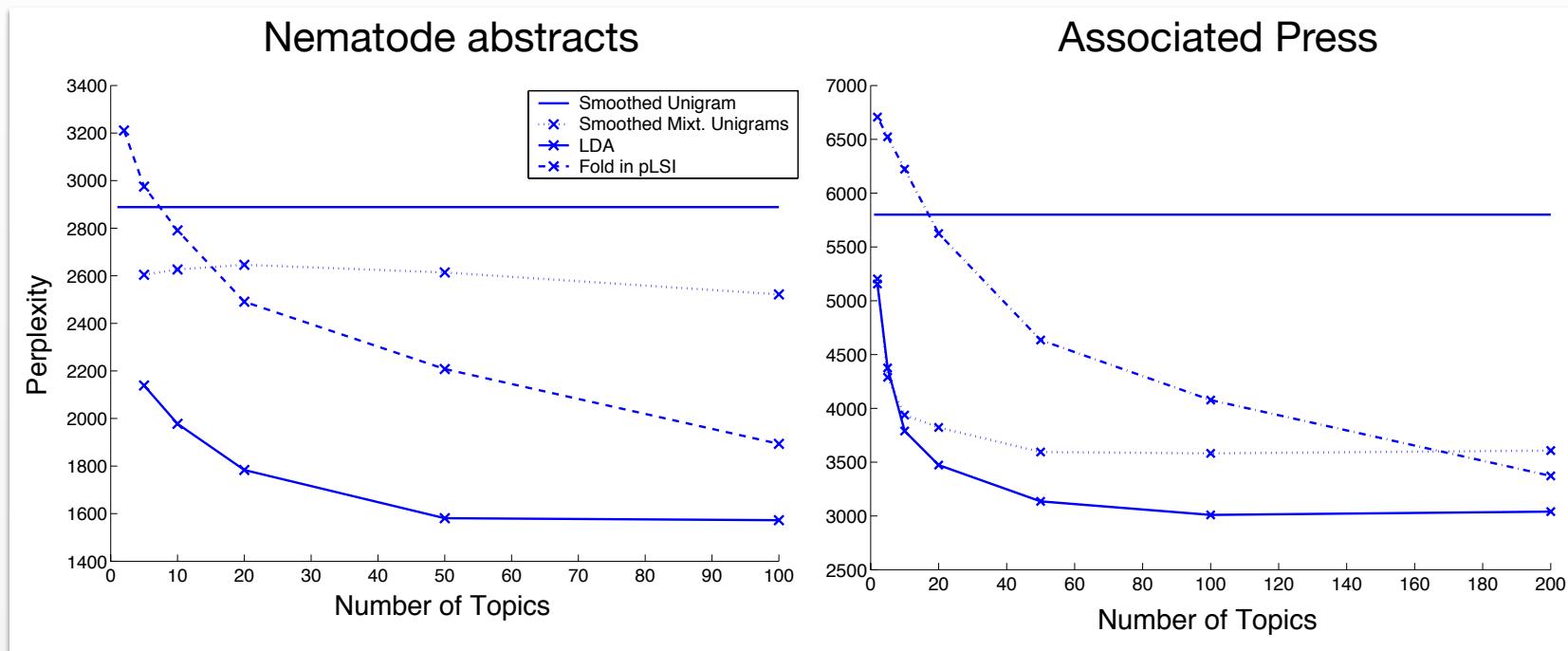
$$q(\theta_d; \gamma_d)$$

Approximation of topic distribution for document d

$$q(\beta_k; \lambda_k)$$

Approximation of word distribution for topic k

Performance Metric: Perplexity



$$\text{Perplexity} = \exp \left[-\frac{1}{D'} \sum_{d=1}^{D'} \frac{1}{N_d} \log(\mathbf{x}'_d | \mathbf{x}, \boldsymbol{\alpha}, \boldsymbol{\eta}) \right] \quad \{\mathbf{x}'_1, \dots, \mathbf{x}'_{D'}\}$$

Exponent of per-word log predictive probability



Topic Models

Shantanu Jain



Extensions of LDA



Borrowing from:
David Blei
(Columbia)

Extensions of LDA

Latent dirichlet allocation

[DM Blei, AY Ng, MI Jordan - Journal of machine Learning research, 2003 - jmlr.org](#)

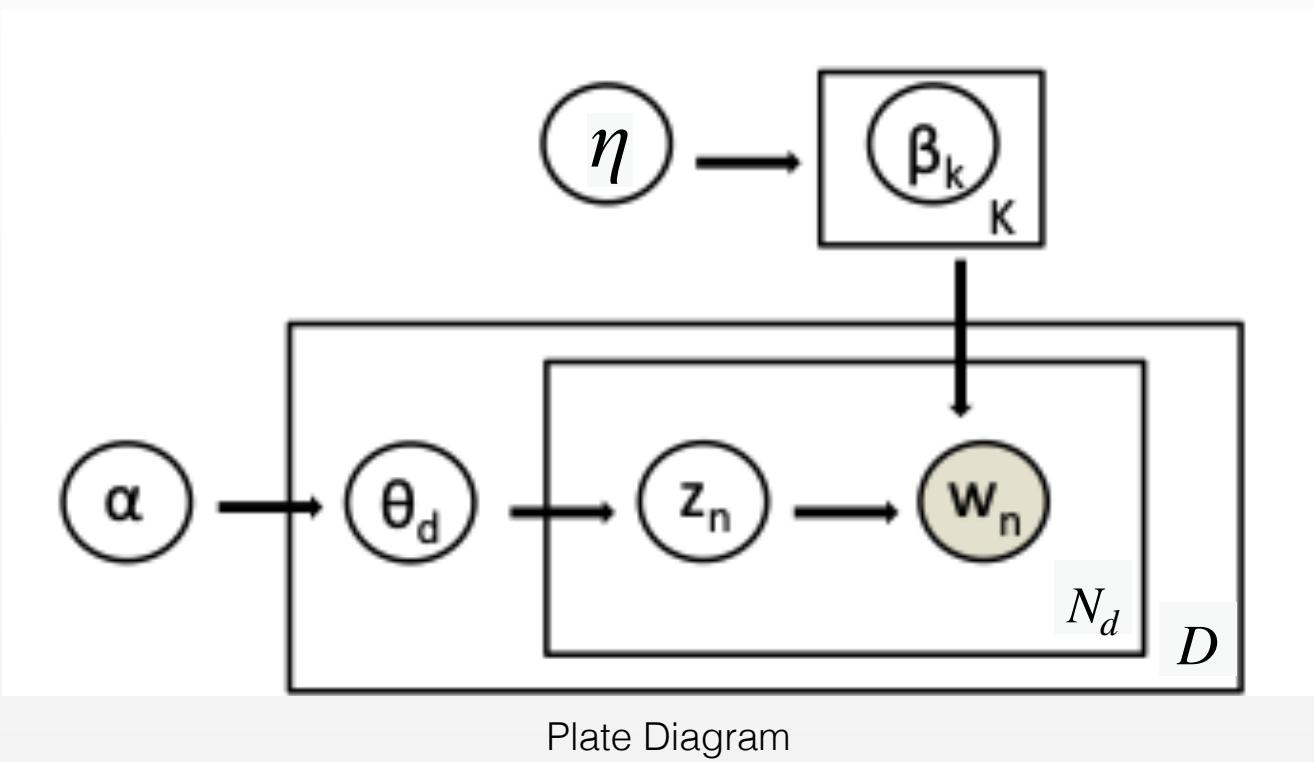
Abstract We describe latent Dirichlet allocation (LDA), a generative probabilistic model for collections of discrete data such as text corpora. LDA is a three-level hierarchical Bayesian model, in which each item of a collection is modeled as a finite mixture over an underlying ...

Cited by 15971 Related articles All 124 versions Cite Save

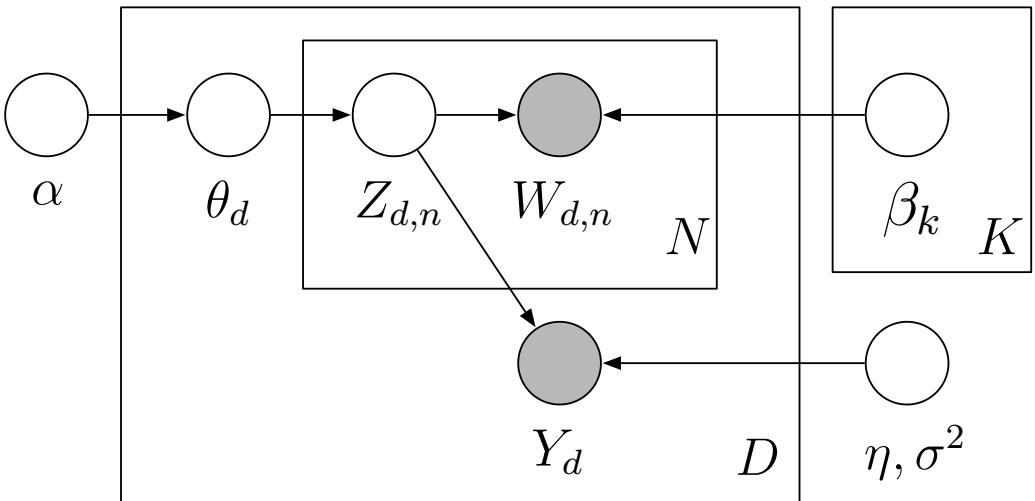
Reasons for popularity of LDA:

- Dirichlet prior gives sparser vectors θ_d
- LDA can be extended to more sophisticated models

Extensions of LDA



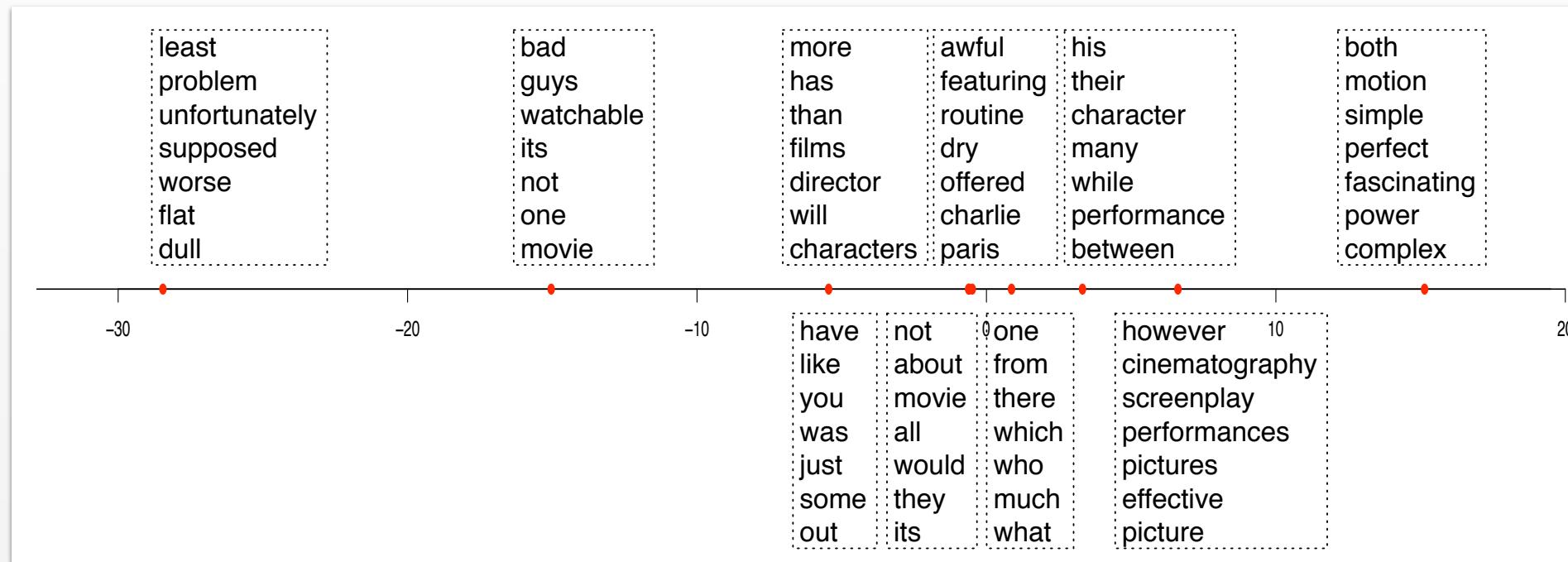
Extensions: Supervised LDA



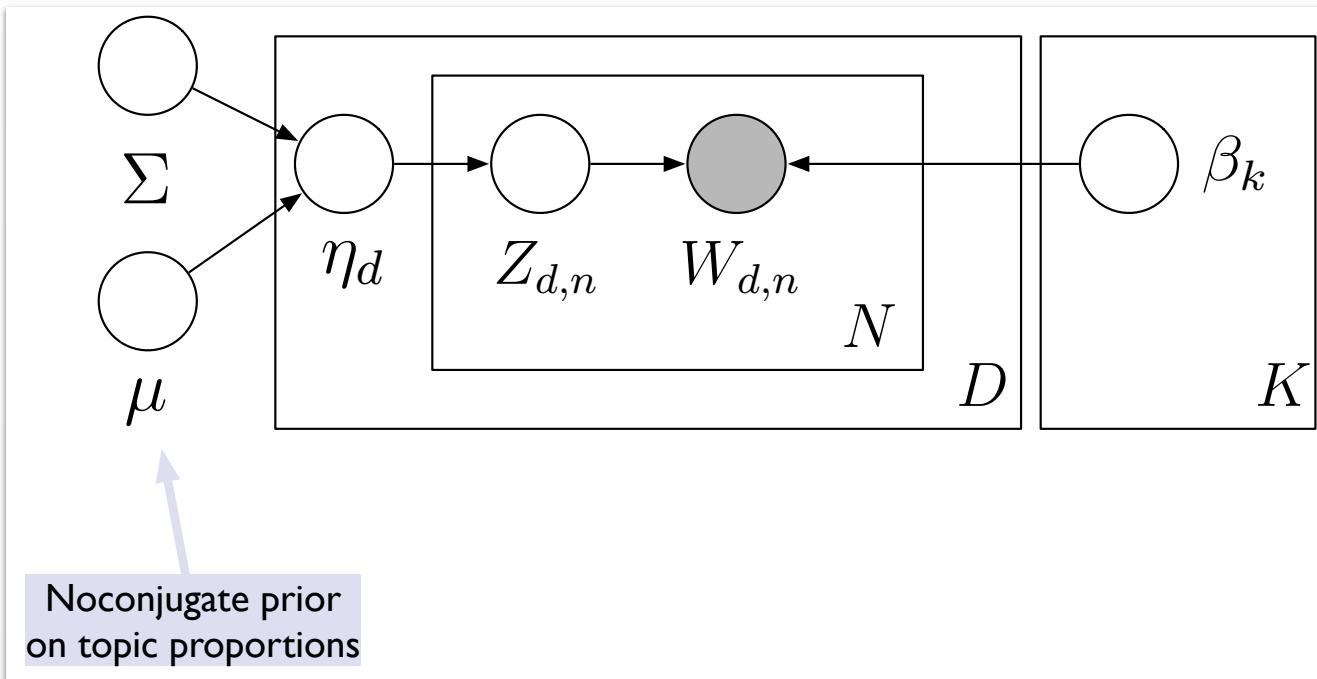
- ① Draw topic proportions $\theta | \alpha \sim \text{Dir}(\alpha)$.
- ② For each word
 - Draw topic assignment $z_n | \theta \sim \text{Mult}(\theta)$.
 - Draw word $w_n | z_n, \beta_{1:K} \sim \text{Mult}(\beta_{z_n})$.
- ③ Draw response variable $y | z_{1:N}, \eta, \sigma^2 \sim \mathcal{N}(\eta^\top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

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

Extensions: Supervised LDA

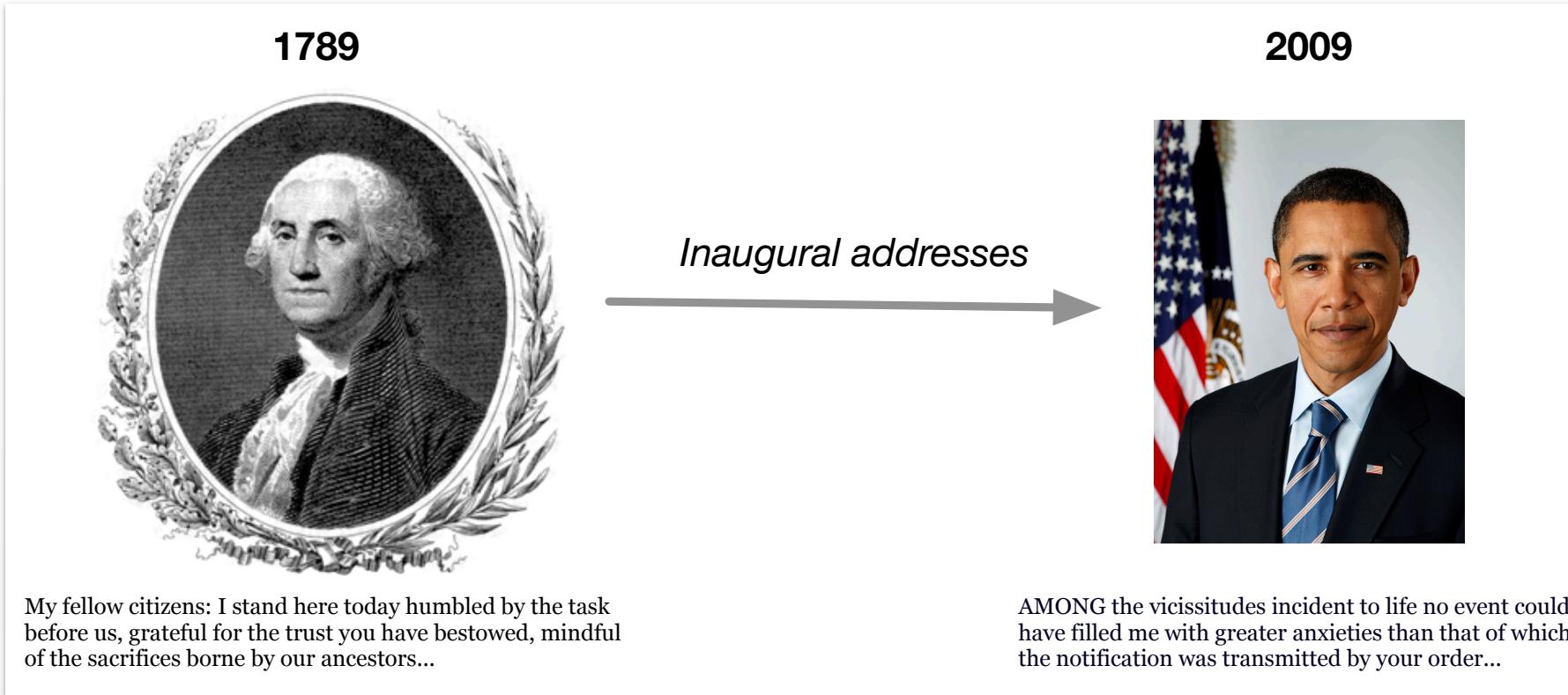


Extensions: Correlated Topic Model



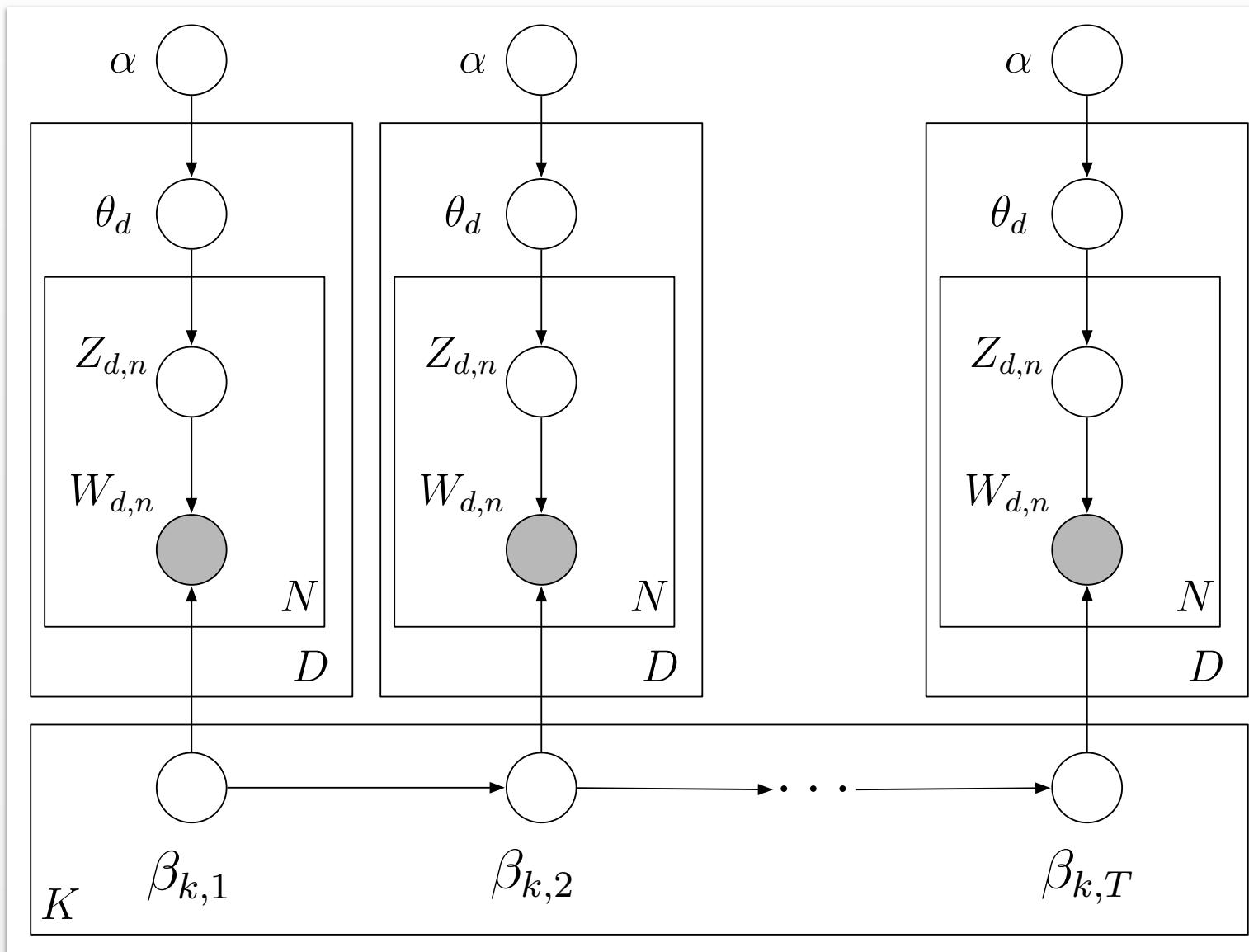
Estimate a covariance matrix Σ that parameterizes correlations between topics in a document

Extensions: Dynamic Topic Models

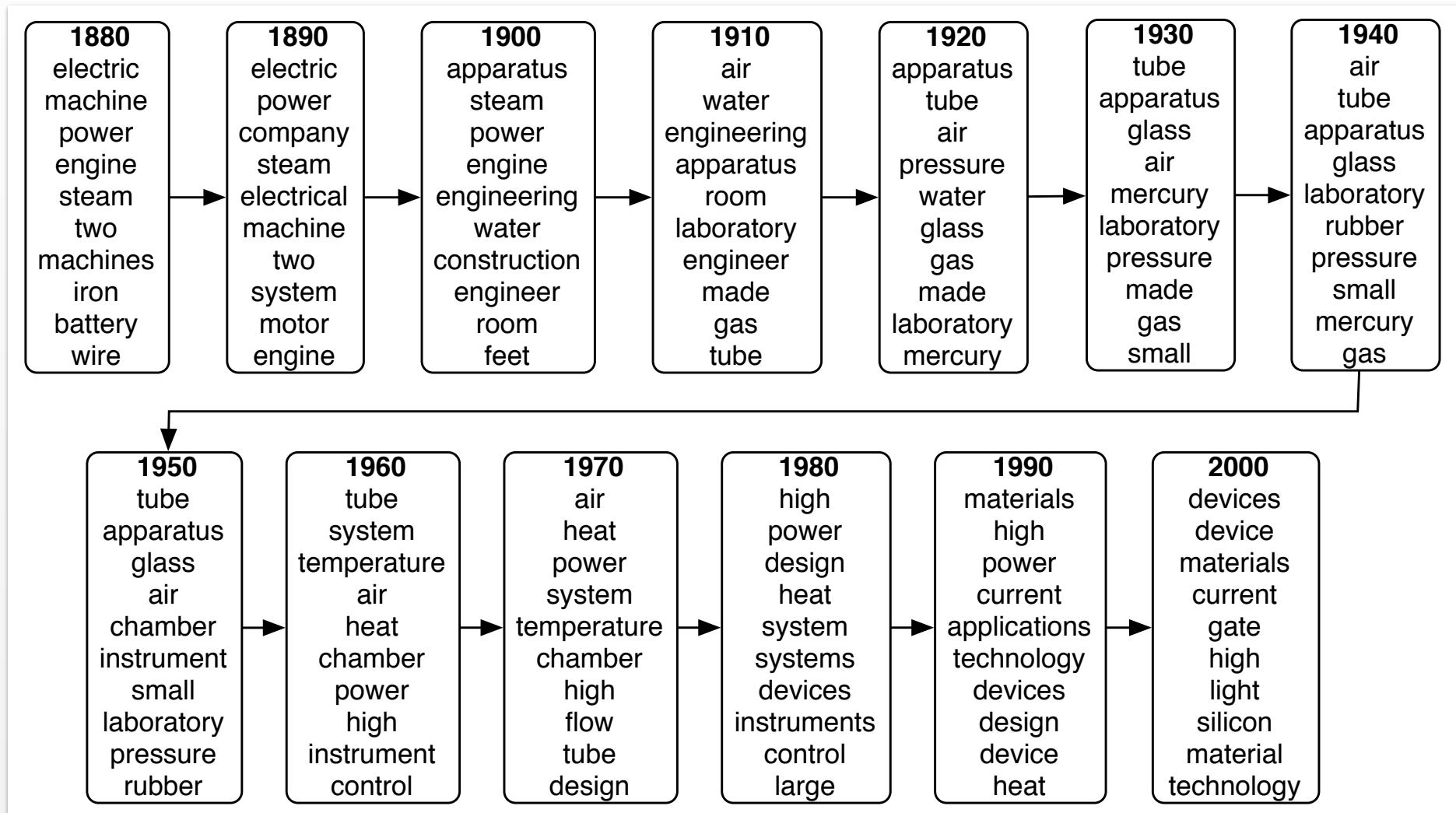


Track changes in word distributions
associated with a topic over time.

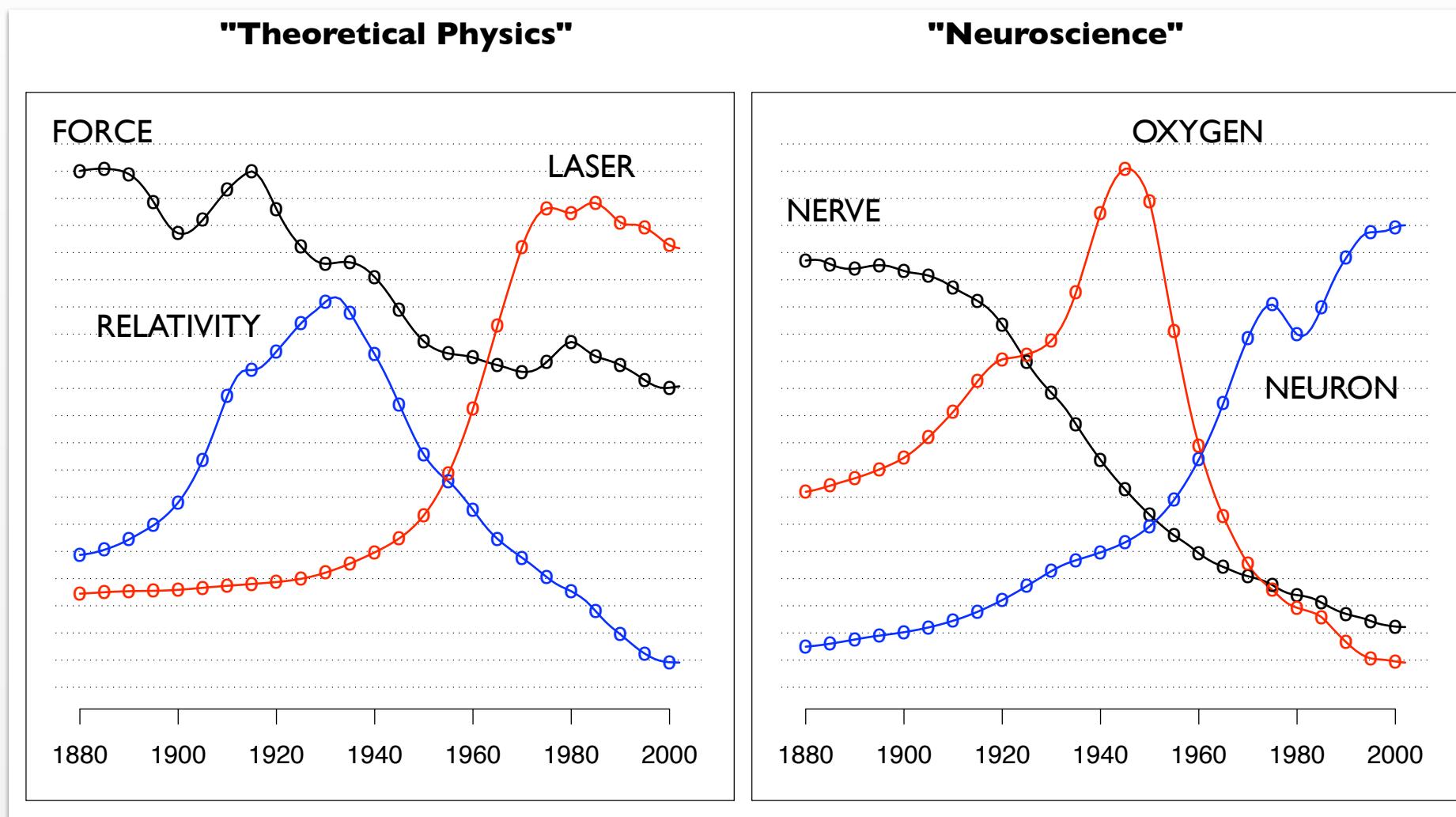
Extensions: Dynamic Topic Models



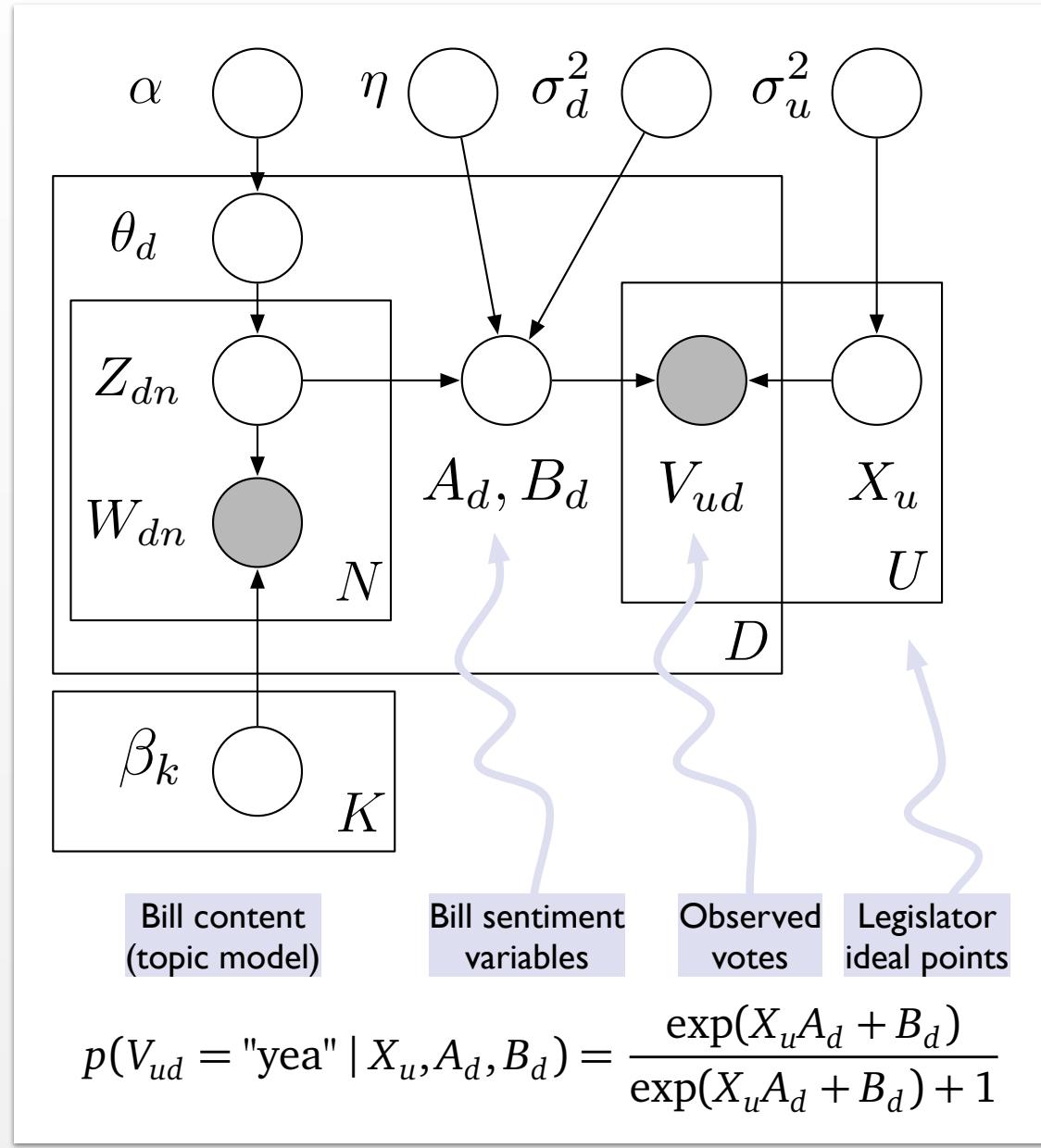
Extensions: Dynamic Topic Models



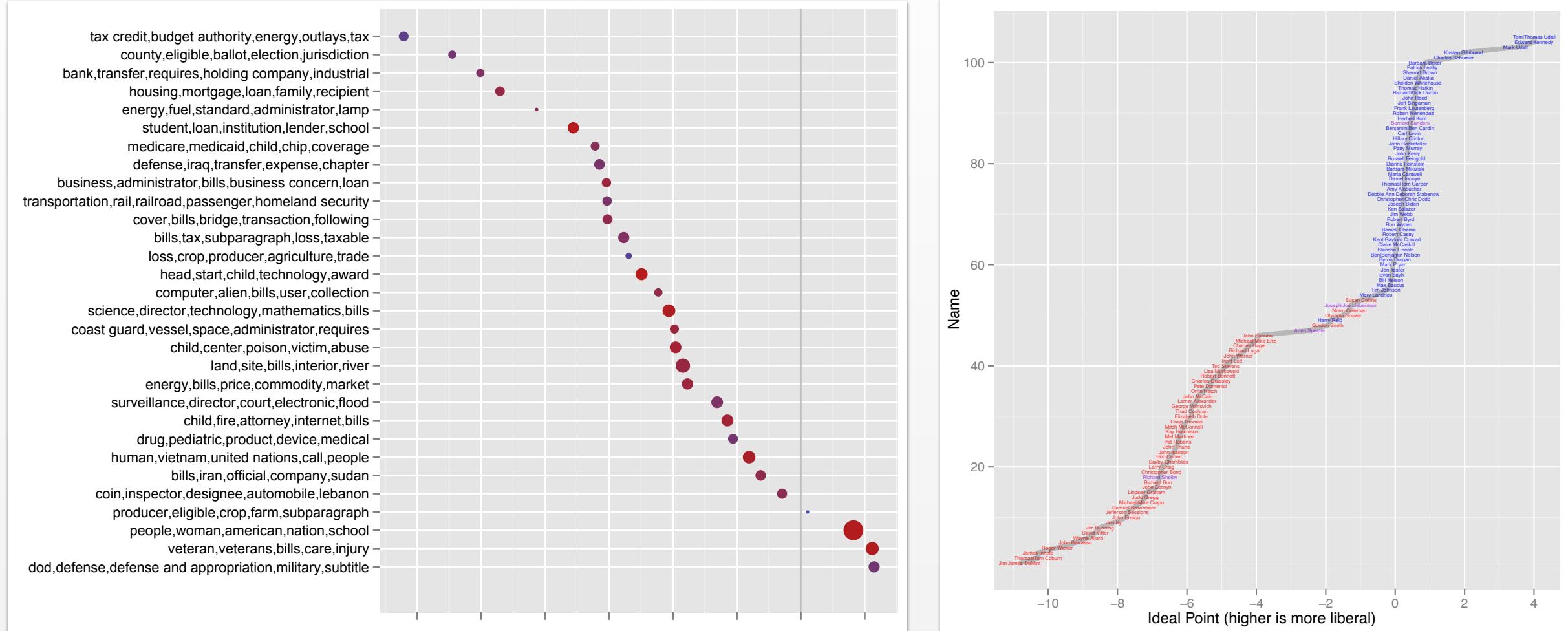
Extensions: Dynamic Topic Models



Extensions: Ideal Point Topic Models



Extensions: Ideal Point Topic Models



$$p(V_{ud} = \text{"yea"} | X_u, A_d, B_d) = \frac{\exp(X_u A_d + B_d)}{\exp(X_u A_d + B_d) + 1}$$

LDA: *Summary*

- Idea: Model documents as *mixtures* over topics
- Model parameters (estimate with VBEM)
 - θ_d Topic probabilities for each document
(K-dimensional vector for each document)
 - β_k Word probabilities for each topic
(V-dimensional vector for each topic)
- Dirichlet Priors: Enforce sparsity, associate a small number of topics which each document
- Extensions: Can design graphical models that build on LDA for a variety of modeling tasks