



Topic Models

Shantanu Jain



Topic Modeling Basics



Borrowing from:
David Blei
(Columbia)

Word Mixtures

Idea: Model text as a “bag” of words (ignore order)

Seeking Life's Bare (Genetic) Necessities

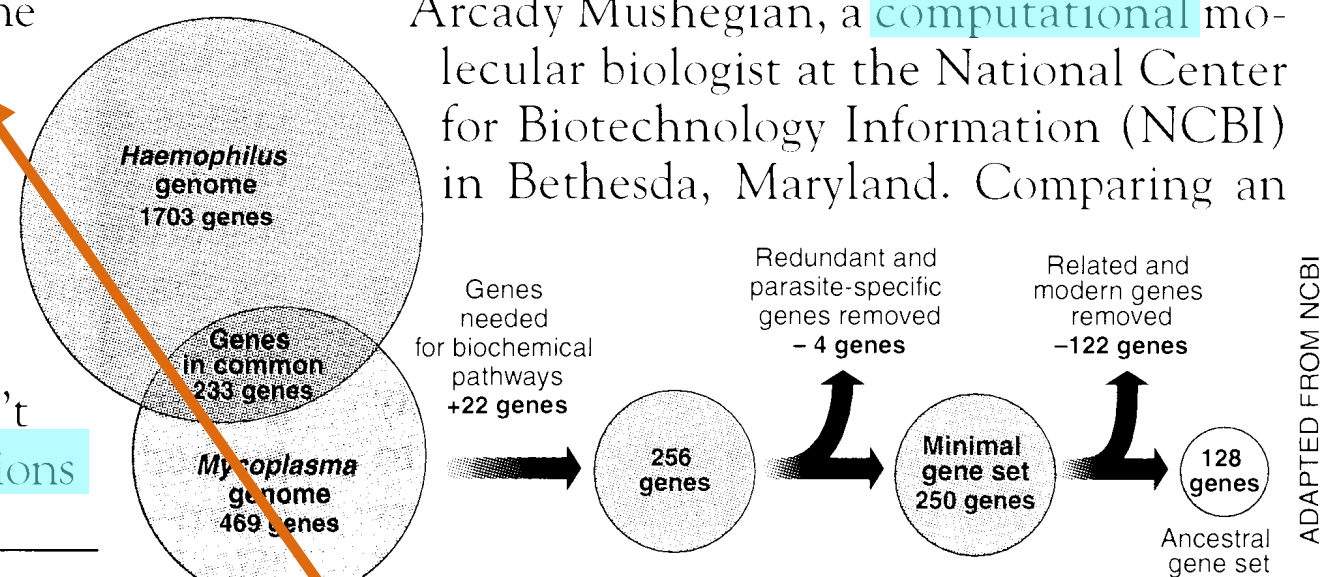
COLD SPRING HARBOR, NEW YORK—How many genes does an organism need to survive? Last week at the genome meeting here,* two genome researchers with radically different approaches presented complementary views of the basic genes needed for life. One research team, using computer analyses to compare known 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 an



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

ADAPTED FROM NCBI

gene	0.04
dna	0.02
genetic	0.01
...	

$$p(x | z = 1)$$

life	0.02
evolve	0.01
organism	0.01
...	

$$p(x | z = 2)$$

brain	0.04
neuron	0.02
nerve	0.01
...	

$$p(x | z = 3)$$

data	0.02
number	0.02
computer	0.01
...	

$$p(x | z = 4)$$

Word in vocabulary: $x_n \in \{1, \dots, V\}$

Topic assignment: $z_n \in \{1, \dots, K\}$

- Total N words in a document
- n denotes the index of the n^{th} word in the document.
- V is the number of words in the vocabulary.
- K is the number of topics.

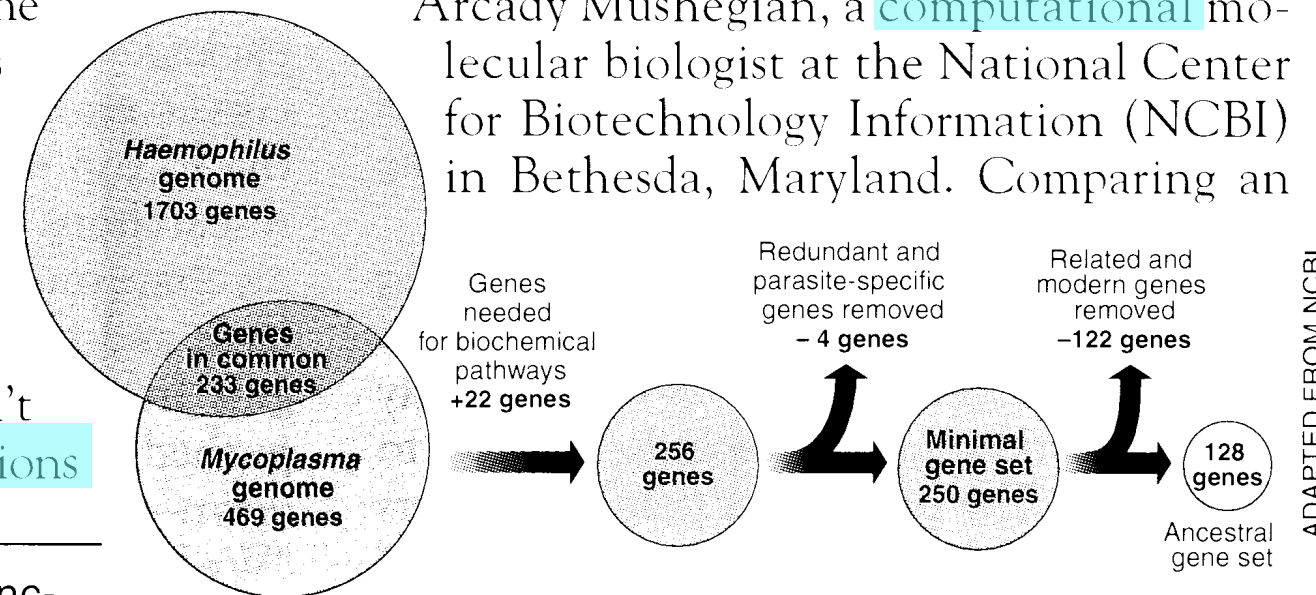
Word Mixtures

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.



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

ADAPTED FROM NCBI

“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 an

gene	0.04
dna	0.02
genetic	0.01
...	

$$p(x | z=1, \beta)$$

life	0.02
evolve	0.01
organism	0.01
...	

$$p(x | z=2, \beta)$$

brain	0.04
neuron	0.02
nerve	0.01
...	

$$p(x | z=3, \beta)$$

data	0.02
number	0.02
computer	0.01
...	

$$p(x | z=4, \beta)$$

θ : topic proportions/probabilities, probability over the K topics

$$\theta = [\theta_1, \theta_2, \dots, \theta_K] \quad \sum_k \theta_k = 1$$

$$p(z_n = k | \theta) = \theta_k$$

β_k : k^{th} topic's word probabilities over the vocabulary

$$\beta_k = [\beta_{k1}, \beta_{k2}, \dots, \beta_{kV}] \quad \sum_i \beta_{ki} = 1$$

$$p(x_n = i | z_n = k, \beta) = \beta_{ki}$$

$$z_n \sim \text{Discrete}(\theta)$$

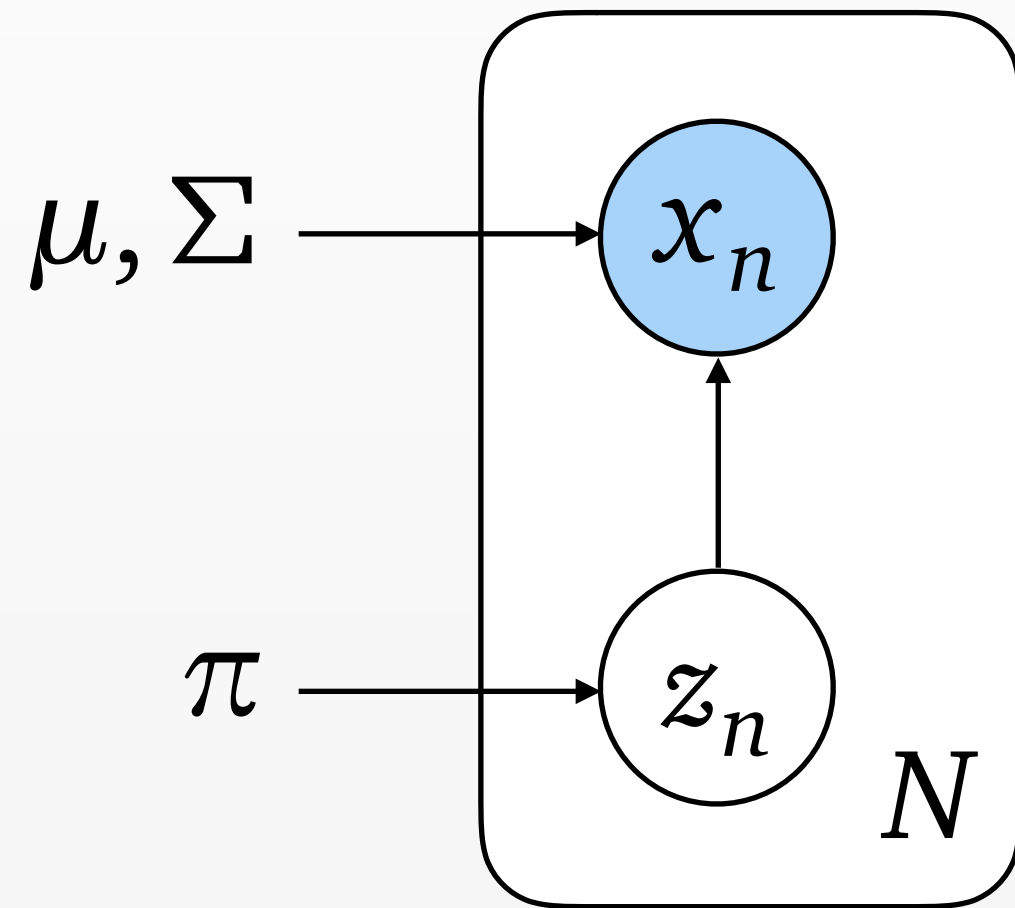
Pick a topic

$$x_n | z_n = k \sim \text{Discrete}(\beta_k)$$

Pick a word given topic

Gaussian Mixtures vs Word Mixtures

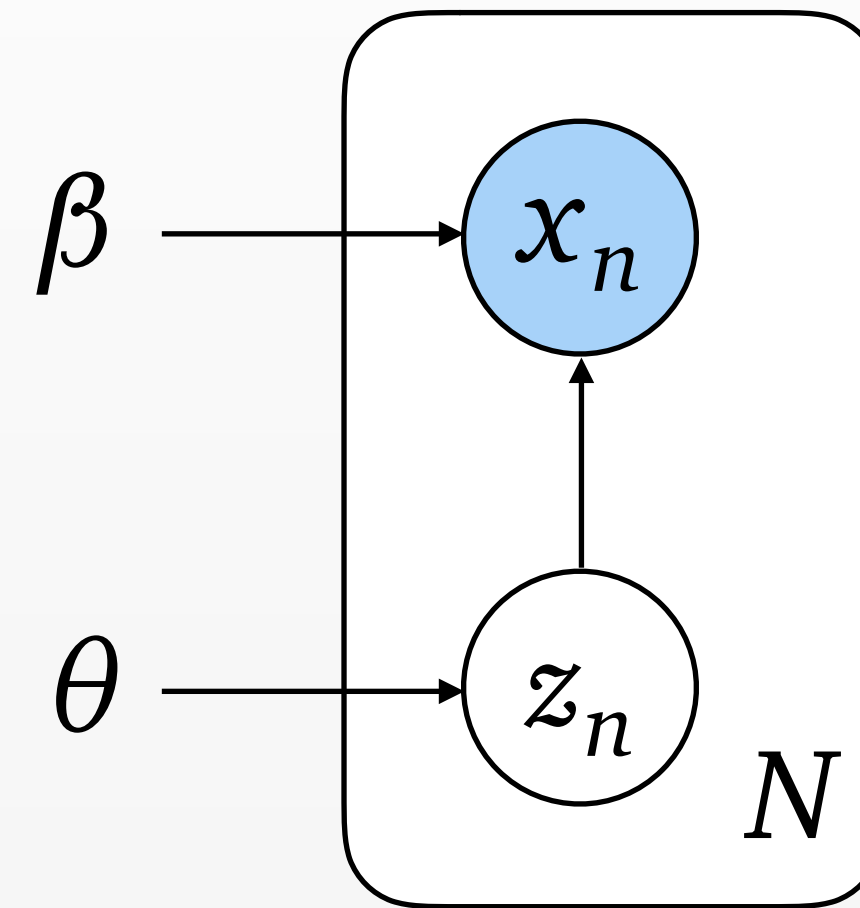
Gaussian Mixture Model



$$z_n \sim \text{Discrete}(\pi_1, \dots, \pi_K)$$

$$x_n \mid z_n = k \sim \text{Normal}(\mu_k, \Sigma_k)$$

Discrete Mixture Model



$$z_n \sim \text{Discrete}(\theta_1, \dots, \theta_K)$$

$$x_n \mid z_n = k \sim \text{Discrete}(\beta_{k,1}, \dots, \beta_{k,V})$$

Difference: Replace Gaussian with Discrete

Topic Modeling

Topics
(shared)

Words in Document
(mixture over topics)

Topic Proportions
(document-specific)

gene	0.04
dna	0.02
genetic	0.01
...	

life	0.02
evolve	0.01
organism	0.01
...	

brain	0.04
neuron	0.02
nerve	0.01
...	

data	0.02
number	0.02
computer	0.01
...	

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** "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 an

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

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

SCIENCE • VOL. 272 • 24 MAY 1996

Idea: Model *corpus* of documents with *shared* topics

Topic Modeling

β_k : Topics

gene	0.04
dna	0.02
genetic	0.01
...	
life	0.02
evolve	0.01
organism	0.01
...	
brain	0.04
neuron	0.02
nerve	0.01
...	
data	0.02
number	0.02
computer	0.01
...	

(shared across documents)

x_d : Words

z_d : Assignments

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

“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 an

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

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

SCIENCE • VOL. 272 • 24 MAY 1996

(unique to each document)

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

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

Distribution over Topic Assignments

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 researcher estimates that for a simple organism just 250 genes are required a mere 128 genes. The other researcher estimated that for a more complex organism 800 genes are needed—but that a mere 100 would do the job. Although the match precise.

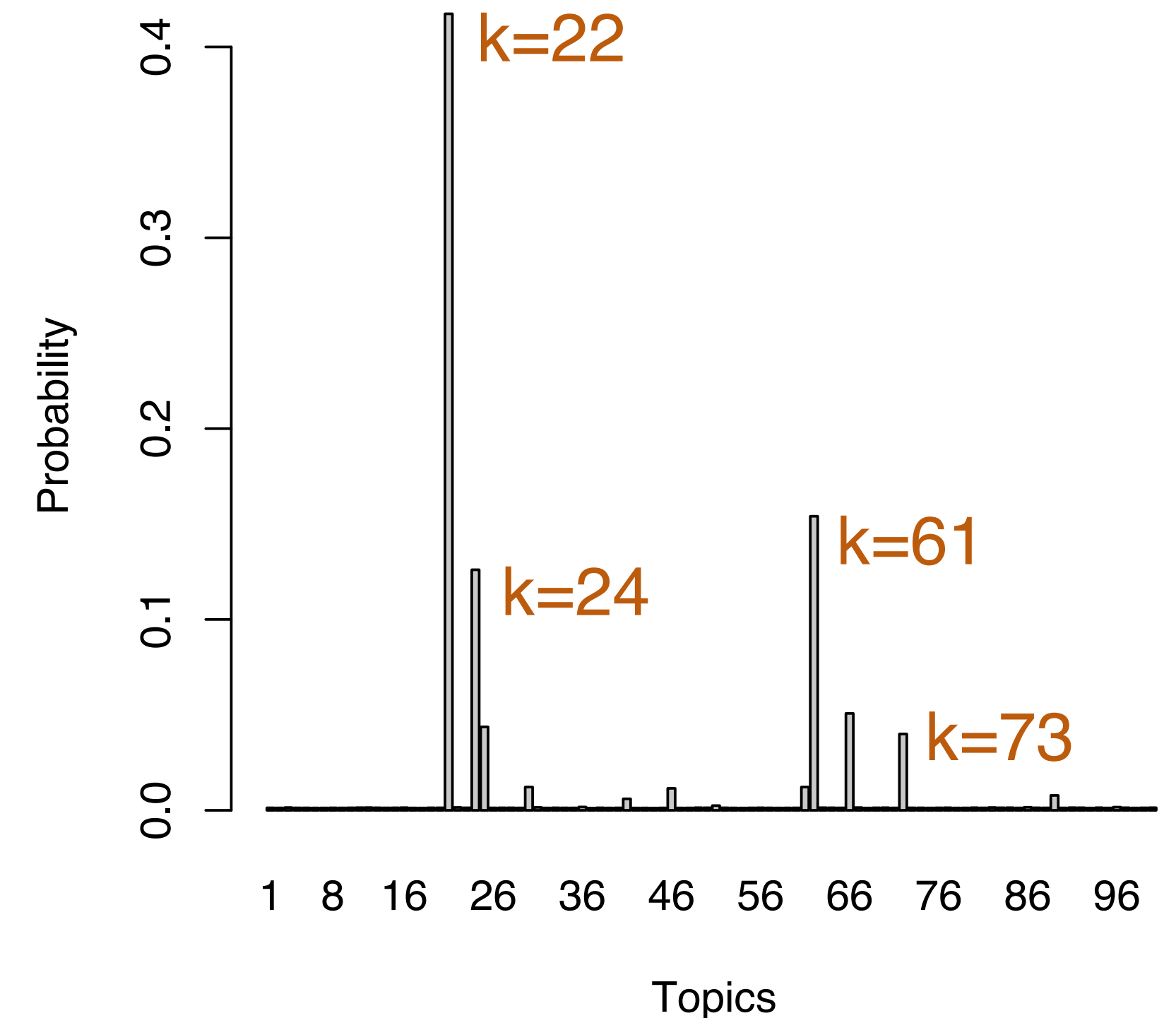
* Genome Mapping, Cold Spring Harbor, May 8 to 12

SCIENCE

“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

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

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




Next Slide: Frequent words in these topics

Most Probable Words in Topics

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

Most frequent
(*within topic*)



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
k=22	k=24	k=61	k=73

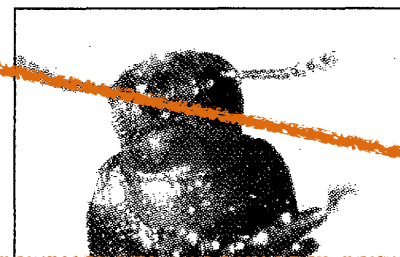
Each Document has Different Topics

Chaotic Beetles

Charles Godfray and Michael Hassell

Ecologists have known since the pioneering work of May in the mid-1970s (1) that the population dynamics of animals and plants can be exceedingly complex. This complexity arises from two sources: The tangled web of interactions that constitute any natural community provide a myriad of different pathways for species to interact, both directly and indirectly. And even in isolated populations the nonlinear feedback processes present in all natural populations result in complex dynamic behavior. Natural populations can show persistent oscillations and chaos, the latter characterized by extreme sensitivity to initial conditions. If such chaotic dynamics were common in nature, then this would have important implications for the management and conservation of natural resources. On page 389 of this issue, Costantino *et al.* (2) provide the

convincing evidence to date of complex dynamics and chaos in a biological population—of the flour beetle, *Tribolium castaneum* (see figure). It has proven extremely dif-



move over the surface of the attractor, sets of adjacent trajectories are pulled apart, then stretched and folded, so that it becomes impossible to predict exact population densities into the future. The strength of the mixing that gives rise to the extreme sensitivity to initial conditions can be measured mathematically estimating the Liapunov exponent, which is positive for chaotic dynamics and nonpositive otherwise. There have been many attempts to estimate attractor dimension and Liapunov exponents from time series

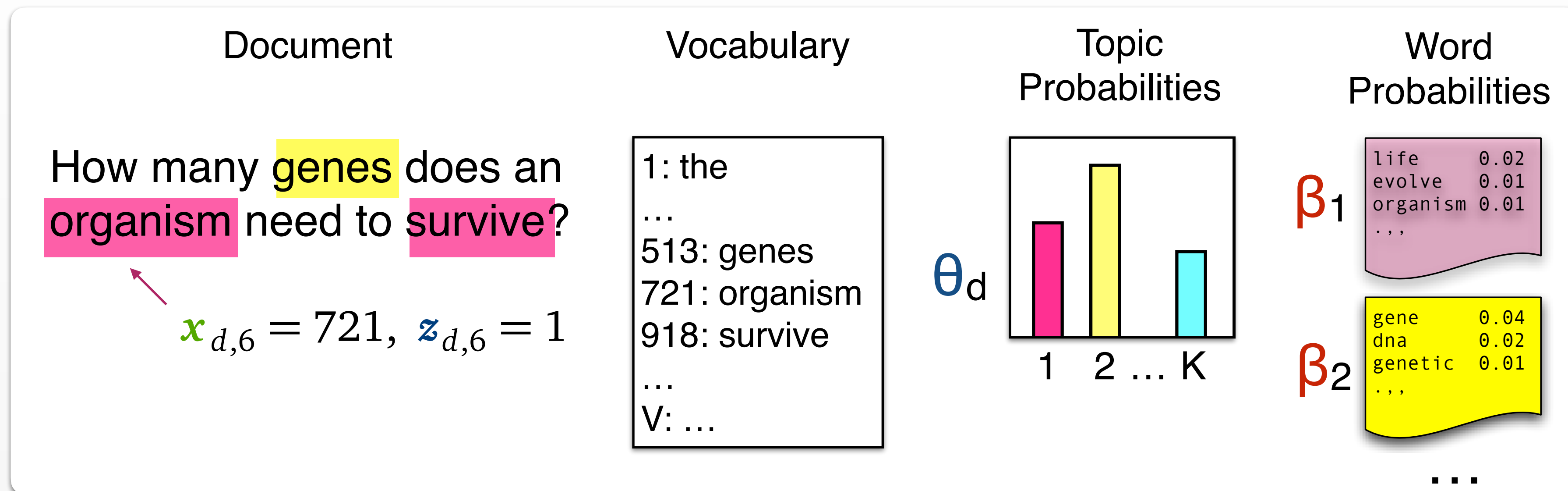
data. The Liapunov exponent is a measure of the average rate of divergence of two trajectories that start close together. A positive Liapunov exponent indicates that the system is chaotic, while a negative Liapunov exponent indicates that the system is periodic or quasi-periodic.

Ecologists have known since the pioneering work of May in the mid-1970s (1) that the population dynamics of animals and plants can be exceedingly complex. This complexity arises from two sources: The tangled web of interactions that constitute any natural community provide a myriad of different pathways for species to interact, both directly and indirectly. And even in isolated populations the nonlinear feedback processes present in all natural populations can result in complex dynamic behavior. Natural

The authors are in the Department of Biology, Imperial College at Silwood Park, Ascot, Berks, SL5 7PZ, UK. Email: m.hassell@ic.ac.uk

problem	model	selection	species
problems	rate	male	forest
mathematical	constant	males	ecology
number	distribution	females	fish
new	time	sex	ecological
mathematics	number	species	conservation
university	size	female	diversity
two	values	evolution	population
first	value	populations	natural
numbers	average	population	ecosystems
work	rates	sexual	populations
time	data	behavior	endangered
mathematicians	density	evolutionary	tropical
chaos	measured	genetic	forests
chaotic	models	reproductive	ecosystem

Estimating the Parameters

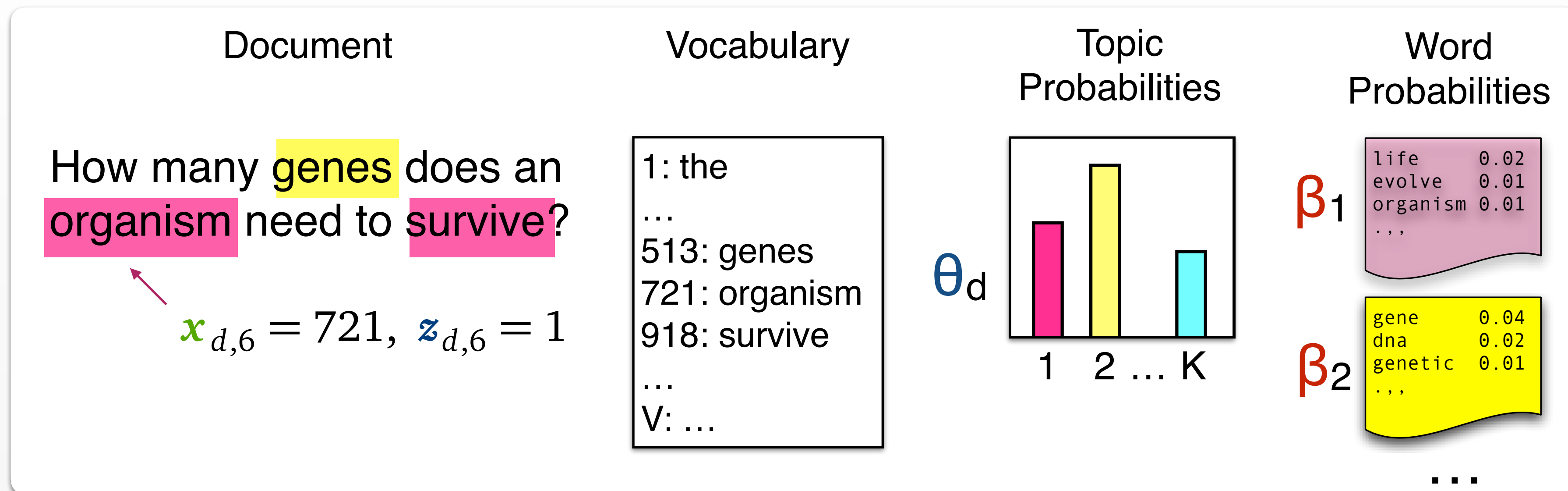


Maximum Likelihood: $\max_{\theta, \beta} \log p(\mathbf{x} | \theta, \beta)$

d=1	$[\mathbf{x}_{1,1}, \mathbf{x}_{1,2}, \dots, \mathbf{x}_{1,N_1}]$	$[\theta_{1,1}, \theta_{1,2}, \dots, \theta_{1,K}]$	$[\beta_{1,1}, \beta_{1,2}, \dots, \beta_{1,V}]$	k=1
d=2	$[\mathbf{x}_{2,1}, \mathbf{x}_{2,2}, \dots, \mathbf{x}_{2,N_2}]$	$[\theta_{2,1}, \theta_{2,2}, \dots, \theta_{2,K}]$	$[\beta_{2,1}, \beta_{2,2}, \dots, \beta_{2,V}]$	k=2
...
d=D	$[\mathbf{x}_{D,1}, \mathbf{x}_{D,2}, \dots, \mathbf{x}_{D,N_D}]$	$[\theta_{D,1}, \theta_{D,2}, \dots, \theta_{D,K}]$	$[\beta_{K,1}, \beta_{K,2}, \dots, \beta_{K,V}]$	k=K

(not a matrix) $D \times K$ $K \times V$

Calculating the Likelihood for each Word



Probability that word n is entry v in the vocabulary

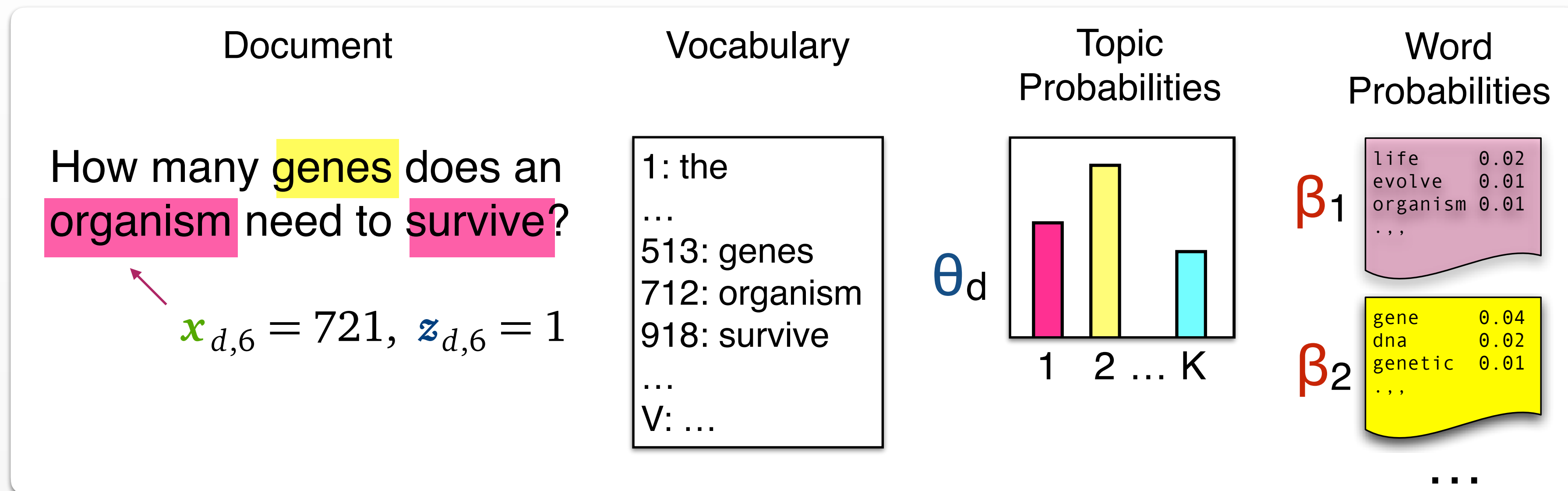
Probability of word v given topic k

Probability that word belongs to topic k

$$\begin{aligned}
 p(x_{d,n} = v \mid \beta, \theta_d) &= \sum_{k=1}^K p(x_{d,n} = v \mid \beta, z_{d,n} = k) p(z_{d,n} = k \mid \theta_d) \\
 &= \sum_{k=1}^K \beta_{k,v} \theta_{d,k}
 \end{aligned}$$

Two orange arrows point from the terms $\beta_{k,v}$ and $\theta_{d,k}$ in the second equation to the text descriptions above.

Computing the Likelihood



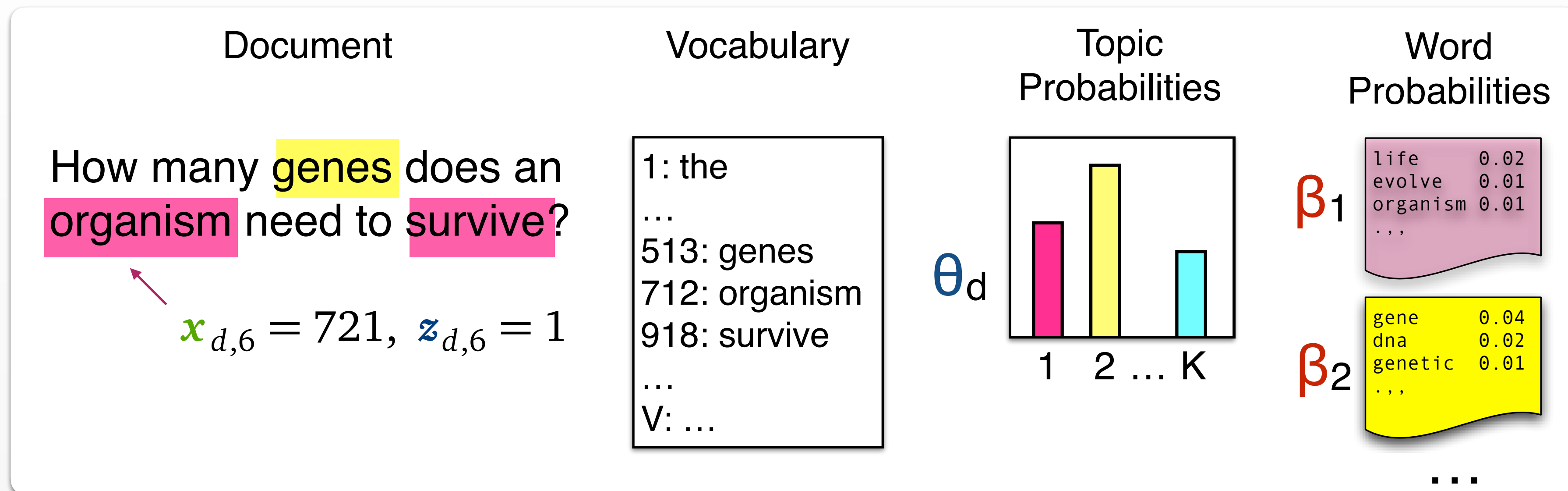
probability of all words $n = 1 \dots N_d$ in document d (use one-hot trick)

$$p(\mathbf{x}_d \mid \boldsymbol{\beta}, \boldsymbol{\theta}_d) = \prod_{n=1}^{N_d} \prod_{v=1}^V p(x_{d,n} = v \mid \boldsymbol{\beta}, \boldsymbol{\theta}_d)^{I[x_{d,n}=v]}$$

take log probability, substitute result from previous slide

$$\log p(\mathbf{x}_d \mid \boldsymbol{\beta}, \boldsymbol{\theta}_d) = \sum_{n=1}^{N_d} \sum_{v=1}^V I[x_{d,n} = v] \log \left(\sum_{k=1}^K \beta_{k,v} \theta_{d,k} \right)$$

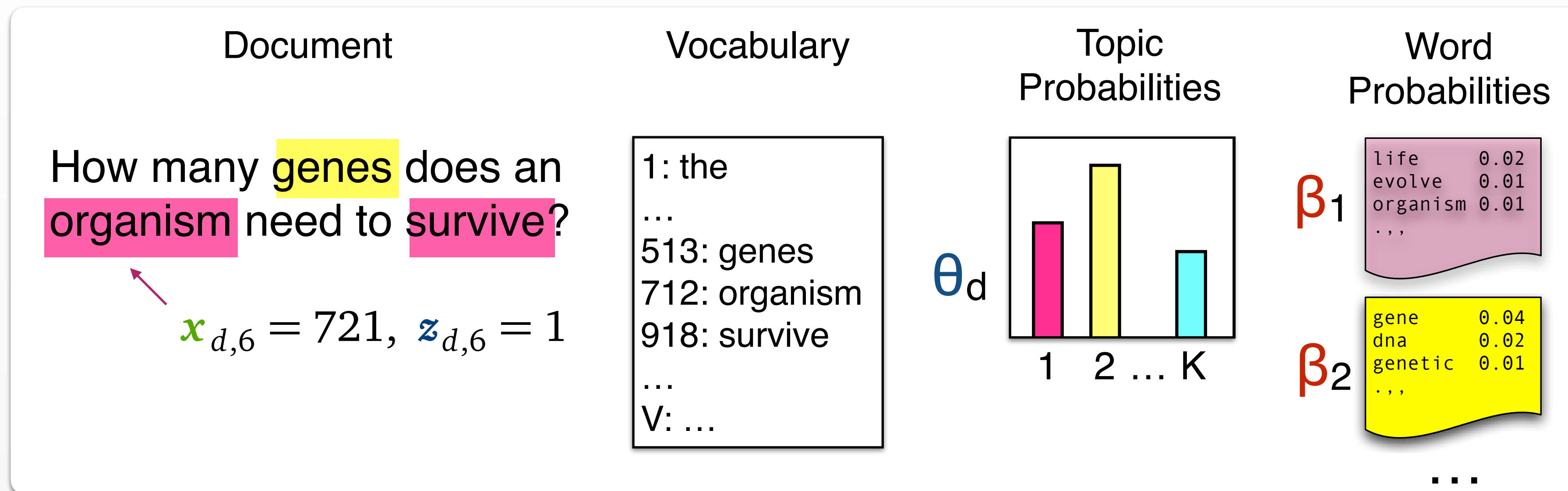
Calculating the Likelihood for all Words



log probability of all words in document d

$$\log p(\mathbf{x}_d \mid \boldsymbol{\beta}, \boldsymbol{\theta}_d) = \sum_{n=1}^{N_d} \sum_{v=1}^V I[\mathbf{x}_{d,n} = v] \log \left(\sum_{k=1}^K \boldsymbol{\beta}_{k,v} \boldsymbol{\theta}_{d,k} \right)$$

Calculating the Likelihood for all Words



log probability of all words in document d

$$\log p(\mathbf{x}_d \mid \boldsymbol{\beta}, \boldsymbol{\theta}_d) = \sum_{v=1}^V \sum_{n=1}^{N_d} I[\mathbf{x}_{d,n} = v] \log \left(\sum_{k=1}^K \boldsymbol{\beta}_{k,v} \boldsymbol{\theta}_{d,k} \right)$$

$$= \mathbf{X}_d \log(\boldsymbol{\theta}_d \boldsymbol{\beta})^\top$$

inner product between bag of word vector,
and log weighted average over topics

bag-of-word vector

$$\mathbf{X}_{d,v} = \sum_{n=1}^{N_d} I[\mathbf{x}_{d,n} = v]$$

Interpretation as Matrix Factorization

Log likelihood

$$\log(p(\mathbf{X}_d | \boldsymbol{\beta}, \boldsymbol{\theta}_d)) = \mathbf{X}_d \log(\boldsymbol{\theta}_d \boldsymbol{\beta})^\top$$

Bag of Word Vector

$$\mathbf{X}_{d,v} = \sum_{n=1}^{N_d} I[\mathbf{x}_{d,n} = v]$$

Word Counts

Topic Counts

Topic Word Probabilities

$$\mathbb{E} \begin{bmatrix} \begin{bmatrix} \mathbf{X}_{1,1} & \cdots & \mathbf{X}_{1,V} \\ \cdots & & \\ \mathbf{X}_{D,1} & \cdots & \mathbf{X}_{D,V} \end{bmatrix} \end{bmatrix} = \begin{bmatrix} N_1 \boldsymbol{\theta}_{1,1} & \cdots & N_1 \boldsymbol{\theta}_{1,K} \\ \cdots & & \\ N_D \boldsymbol{\theta}_{D,1} & \cdots & N_D \boldsymbol{\theta}_{D,K} \end{bmatrix} \begin{bmatrix} \boldsymbol{\beta}_{1,1} & \cdots & \boldsymbol{\beta}_{1,V} \\ \cdots & & \\ \boldsymbol{\beta}_{K,1} & \cdots & \boldsymbol{\beta}_{K,V} \end{bmatrix}$$

(D x V)
(D x K)
(K x V)

<p style="color: #8B4513; font-style: italic;">stocks</p> <p style="color: #8B4513; font-style: italic;">chairman</p> <p style="color: #8B4513; font-style: italic;">the</p> <p style="color: #8B4513; font-style: italic;">wins</p> <p style="color: #8B4513; font-style: italic;">game</p>	<p style="color: #8B4513; font-style: italic;">finance</p> <p style="color: #8B4513; font-style: italic;">Sports</p> <p style="color: #8B4513; font-style: italic;">stocks</p> <p style="color: #8B4513; font-style: italic;">game</p>
$\begin{bmatrix} 2 & 4 & 8 & \cdots & 0 & 1 \\ \cdots & & & & & \\ 0 & 1 & 7 & \cdots & 2 & 3 \end{bmatrix}$	$\approx \begin{bmatrix} 112 \cdot 0.91 & \cdots & 112 \cdot 0.01 \\ \cdots & & \\ 234 \cdot 0.02 & \cdots & 234 \cdot 0.86 \end{bmatrix} \begin{bmatrix} 0.0081 & \cdots & 0.0002 \\ \cdots & & \\ 0.0001 & \cdots & 0.0072 \end{bmatrix}$

Relationship to Latent Semantic Analysis

LSA: Factorize word counts (using PCA)

$$\mathbf{X} \ (V \times D) \approx \mathbf{U} \ (V \times K) \mathbf{Z} \ (K \times D)$$
$$\begin{pmatrix} \text{stocks: } 2 & \dots & 0 \\ \text{chairman: } 4 & \dots & 1 \\ \text{the: } 8 & \dots & 7 \\ \dots & \vdots & \vdots \\ \text{wins: } 0 & \dots & 2 \\ \text{game: } 1 & \dots & 3 \end{pmatrix} \approx \begin{pmatrix} 0.4 & \dots & -0.001 \\ 0.8 & \dots & 0.03 \\ 0.01 & \dots & 0.04 \\ \vdots & \dots & \vdots \\ 0.002 & \dots & 2.3 \\ 0.003 & \dots & 1.9 \end{pmatrix} \begin{pmatrix} | & & | \\ \mathbf{z}_1 & \dots & \mathbf{z}_n \\ | & & | \end{pmatrix}$$

Topic Models: Factorize word counts (using mixture model)

$$\mathbb{E}[\mathbf{X}^T] \ (V \times D) = \boldsymbol{\beta}^T \ (V \times K) \boldsymbol{\theta}^T \ (K \times D) \mathbf{N} \mathbf{I} \ (D \times D)$$

$$\mathbb{E}[\mathbf{X}] \ (D \times V) = \mathbf{N} \mathbf{I} \ (D \times D) \boldsymbol{\theta} \ (D \times K) \boldsymbol{\beta} \ (K \times V)$$

Topic Models: *Summary so far*

Core Idea:

Model documents as *mixtures* over topics

Model Parameters:

θ_d Topic probabilities for each document
(K-dimensional vector)

β_k Word probabilities for each topic
(V-dimensional vector)

Relationship to Dimensionality Reduction:

Similar to LSA, but assumes Discrete mixture instead of Gaussian distribution on word counts



Topic Models

Shantanu Jain



Estimating Parameters

Maximum Likelihood with EM

Review: Topic Modeling

β_k : Topics

x_d : Words

z_d : Assignments

θ_d : Topic Proportions

gene	0.04
dna	0.02
genetic	0.01
...	
life	0.02
evolve	0.01
organism	0.01
...	
brain	0.04
neuron	0.02
nerve	0.01
...	
data	0.02
number	0.02
computer	0.01
...	

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

"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 an

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

SCIENCE • VOL. 272 • 24 MAY 1996

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

(unique to each document)

(shared across documents)

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

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

Review: Interpretation as Matrix Factorization

Log marginal likelihood

$$\log(p(\mathbf{X}_d | \boldsymbol{\beta}, \boldsymbol{\theta}_d)) = \mathbf{X}_d \log(\boldsymbol{\theta}_d \boldsymbol{\beta})^\top$$

Bag of Word Vector

$$\mathbf{X}_{d,v} = \sum_{n=1}^{N_d} I[\mathbf{x}_{d,n} = v]$$

Word Counts

Topic Counts

Topic Word Probabilities

$$\mathbb{E} \begin{bmatrix} \begin{bmatrix} \mathbf{X}_{1,1} & \cdots & \mathbf{X}_{1,V} \\ \cdots & & \\ \mathbf{X}_{D,1} & \cdots & \mathbf{X}_{D,V} \end{bmatrix} \end{bmatrix} = \begin{bmatrix} N_1 \boldsymbol{\theta}_{1,1} & \cdots & N_1 \boldsymbol{\theta}_{1,K} \\ \cdots & & \\ N_D \boldsymbol{\theta}_{D,1} & \cdots & N_D \boldsymbol{\theta}_{D,K} \end{bmatrix} \begin{bmatrix} \boldsymbol{\beta}_{1,1} & \cdots & \boldsymbol{\beta}_{1,V} \\ \cdots & & \\ \boldsymbol{\beta}_{K,1} & \cdots & \boldsymbol{\beta}_{K,V} \end{bmatrix}$$

(D x V)
(D x K)
(K x V)

	chairman	the	wins	game		"sports"	"finance"	stocks	game			
stocks	2	4	8	...	0	1	112 · 0.91	...	112 · 0.01	0.0081	...	0.0002
chairman	4	8	...	0	1
the	0	1	7	...	2	3	234 · 0.02	...	234 · 0.86	0.0001	...	0.0072
wins	≈											
game												

Relationship to Latent Semantic Analysis

LSA: Factorize matrix of word counts (using PCA)

$$\begin{matrix} \mathbf{X} & (V \times D) & \approx & \mathbf{U} & (V \times K) & \mathbf{Z} & (K \times D) \\ \left(\begin{array}{l} \text{stocks: } 2 \dots\dots\dots 0 \\ \text{chairman: } 4 \dots\dots\dots 1 \\ \text{the: } 8 \dots\dots\dots 7 \\ \dots \vdots \dots\dots\dots \vdots \\ \text{wins: } 0 \dots\dots\dots 2 \\ \text{game: } 1 \dots\dots\dots 3 \end{array} \right) & \approx & \left(\begin{array}{l} 0.4 \dots -0.001 \\ 0.8 \dots 0.03 \\ 0.01 \dots 0.04 \\ \vdots \dots \vdots \\ 0.002 \dots 2.3 \\ 0.003 \dots 1.9 \end{array} \right) & \left(\begin{array}{l} | \qquad | \\ \mathbf{z}_1 \dots \mathbf{z}_n \\ | \qquad | \end{array} \right) \end{matrix}$$

LSA: Assume Gaussian distribution

Topic Models: Assume mixture of Discrete distributions

Estimating Model Parameters

Question: How can we estimate β_k and θ_d ?

1. Expectation Maximization
(this video)
2. Variational Inference
(will discuss at a high level)
3. Gibbs Sampling
(not in this module)

PLSI/PLSA*: EM for Topic Models

Generative Model

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

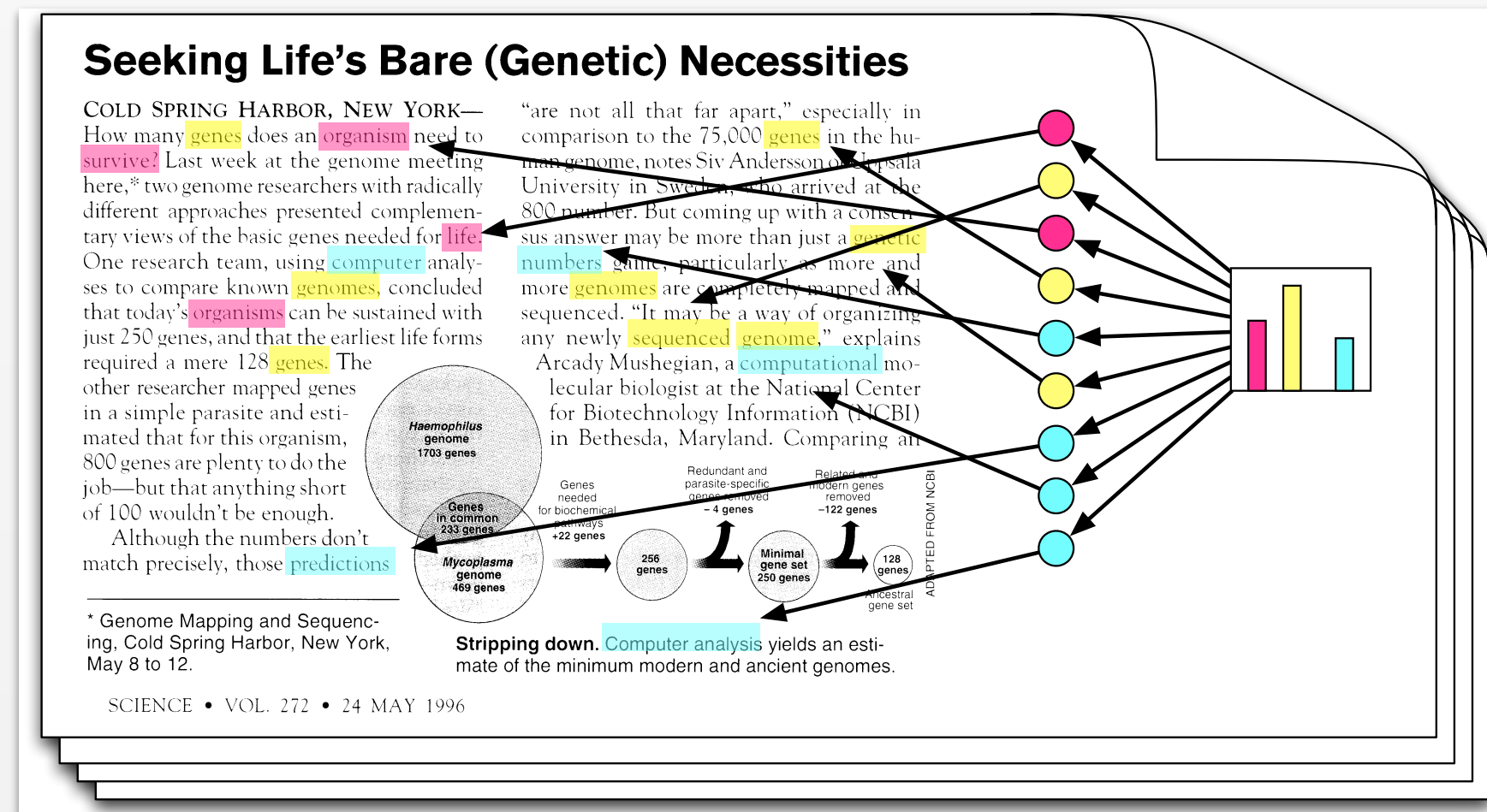
E-step: Update assignments

Calculate probability that word n in document d belongs to topic k

$$\phi_{d,n,k} = p(\mathbf{z}_{d,n} = k \mid \mathbf{x}_{d,n}, \boldsymbol{\beta}, \boldsymbol{\theta}_d)$$

M-step: Update parameters

Use assignment probabilities ϕ_d to update topics assignment probabilities $\boldsymbol{\theta}_d$ and topic word probabilities $\boldsymbol{\beta}_k$



*(Probabilistic Latent Semantic Indexing, a.k.a. Probabilistic Latent Semantic Analysis)

PLSI/PLSA: E-step

$$\begin{aligned}\phi_{d,n,k} &= p(\mathbf{z}_{d,n} = k \mid \mathbf{x}_{d,n} = v, \boldsymbol{\beta}, \boldsymbol{\theta}_d) \\ &= \frac{p(\mathbf{x}_{d,n} = v, \mathbf{z}_{d,n} = k \mid \boldsymbol{\beta}, \boldsymbol{\theta}_d)}{p(\mathbf{x}_{d,n} = v \mid \boldsymbol{\beta}, \boldsymbol{\theta}_d)} && \text{(Apply Bayes' Rule)} \\ &= \frac{\theta_{d,k} \boldsymbol{\beta}_{k,v}}{\sum_{l=1}^K \theta_{d,l} \boldsymbol{\beta}_{l,v}} && \text{(Substitute results from previous slides)}\end{aligned}$$

General Form, with One-hot Indexing Trick

$$\phi_{d,n,k} = \frac{\theta_{d,k} \left(\sum_{v=1}^V \boldsymbol{\beta}_{k,v} I[\mathbf{x}_{d,n} = v] \right)}{\sum_{l=1}^K \theta_{d,l} \left(\sum_{v=1}^V \boldsymbol{\beta}_{l,v} I[\mathbf{x}_{d,n} = v] \right)}$$

PLSI/PLSA*: EM for Topic Models

Generative Model

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

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

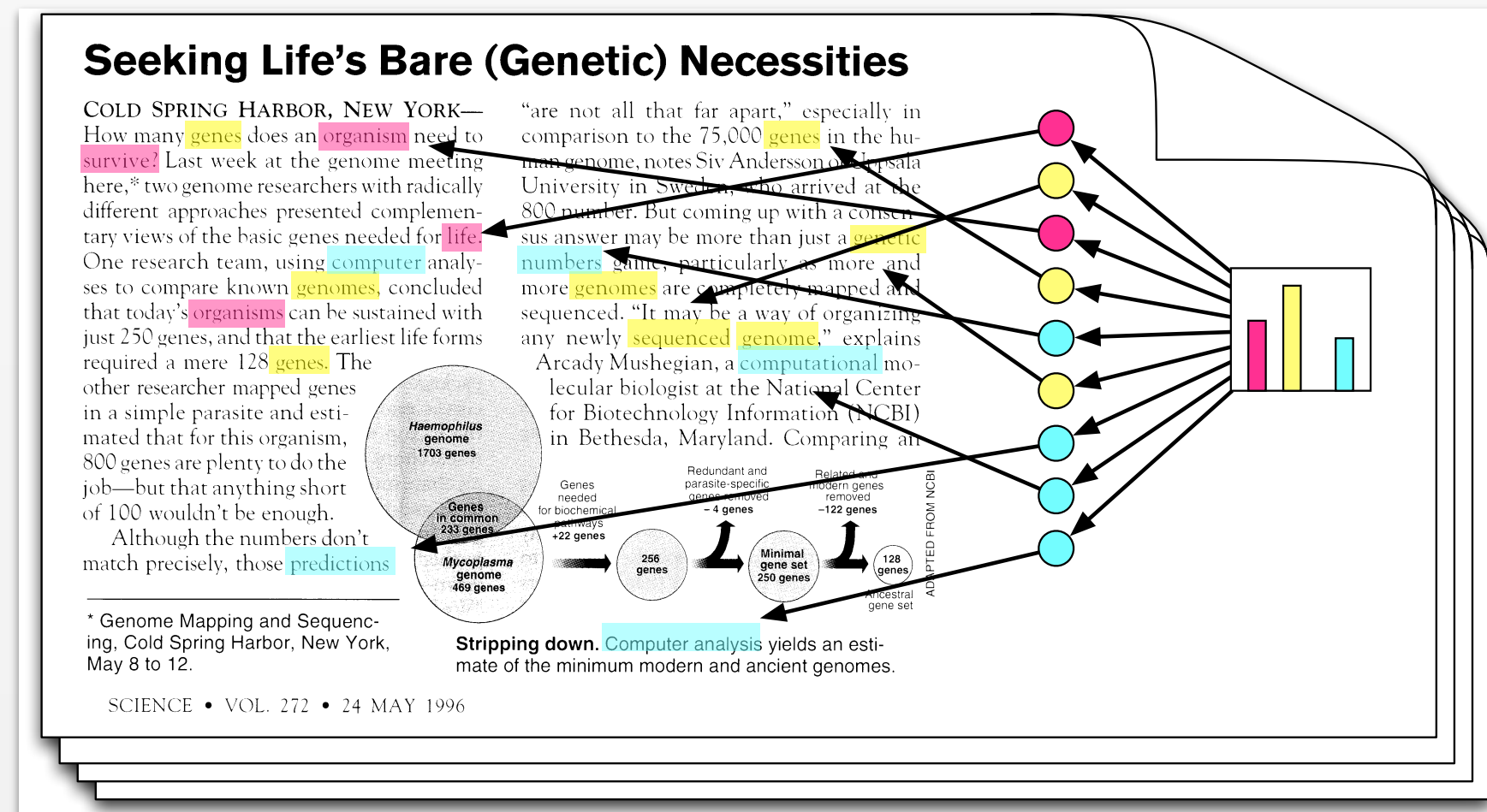
E-step: Update assignments

$$\phi_{d,n,k} = p(\mathbf{z}_{d,n} = k \mid \mathbf{x}_{d,n} = \mathbf{v}, \boldsymbol{\beta}, \boldsymbol{\theta}_d)$$

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

M-step: Update parameters

Use assignment probabilities ϕ_d to update topics assignment probabilities θ_d and topic word probabilities β_k



*(Probabilistic Latent Semantic Indexing, a.k.a. Probabilistic Latent Semantic Analysis)

PLSI/PLSA: M-Step

Idea: Compute (expected) sufficient statistics

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

Probability that word n in document d belongs to topic k

$$N_{d,k}^{\theta} = \sum_{n=1}^{N_d} \phi_{d,n,k}$$

Number of words in document d that belong to topic k

$$N_{k,v}^{\beta} = \sum_{d=1}^D \sum_{n=1}^{N_d} \phi_{d,n,k} I[x_{d,n} = v]$$

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

M-Step: Update parameters using sufficient statistics

$$\theta_{d,k} = \frac{N_{d,k}^{\theta}}{N_d}$$

Fraction of topic k in document d

$$\beta_{k,v} = \frac{N_{k,v}^{\beta}}{\sum_{d=1}^D N_{d,k}^{\theta}}$$

Fraction of word v in topic k

PLSI/PLSA*: EM for Topic Models

Generative Model

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

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

E-step: Update assignments

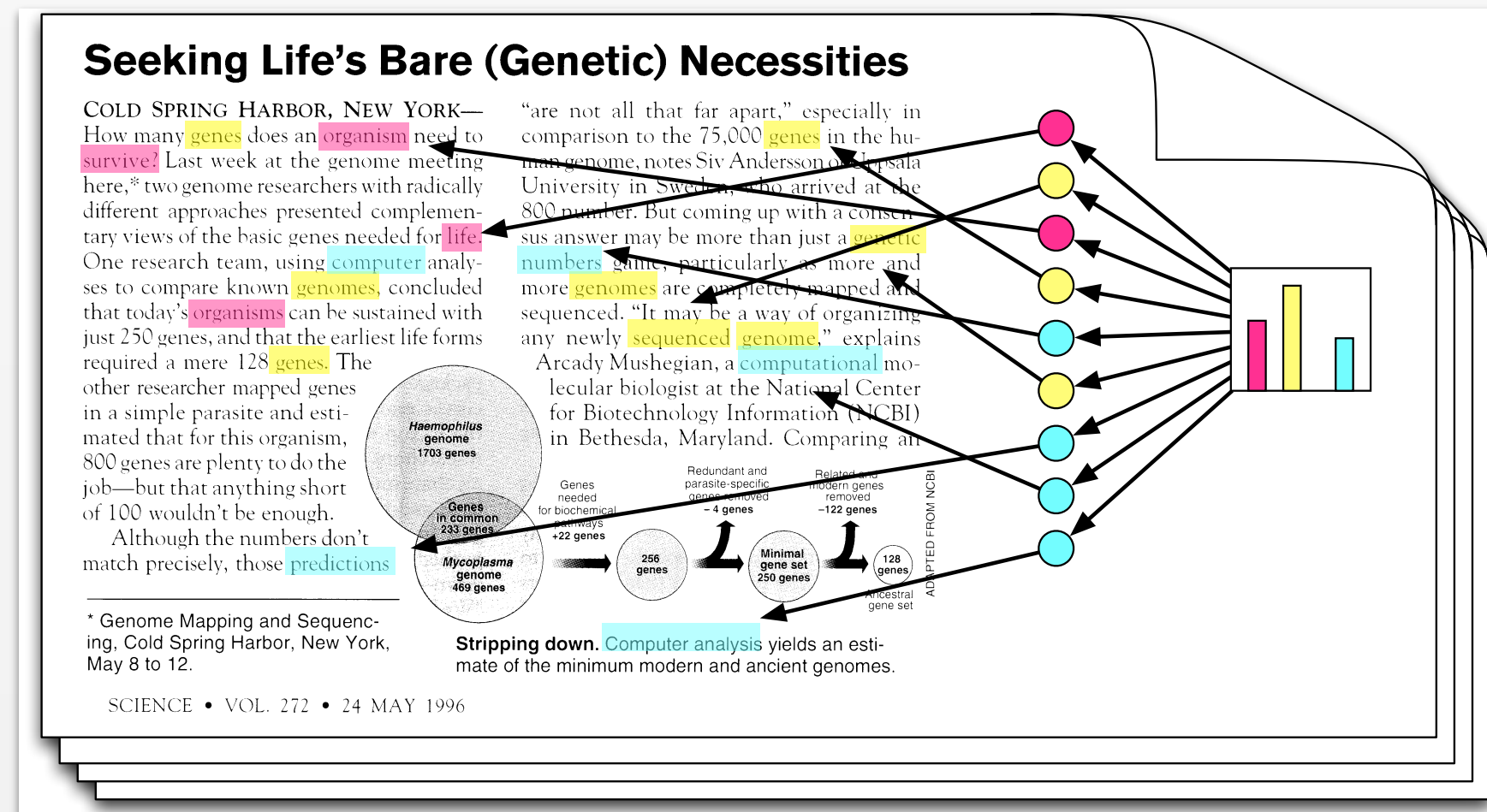
$$\phi_{d,n,k} = p(\mathbf{z}_{d,n} = k \mid \mathbf{x}_{d,n} = \mathbf{v}, \boldsymbol{\beta}, \boldsymbol{\theta}_d)$$

$$= \frac{\boldsymbol{\theta}_{d,k} \left(\sum_{v=1}^V \boldsymbol{\beta}_{k,v} I[\mathbf{x}_{d,n} = \mathbf{v}] \right)}{\sum_{l=1}^K \boldsymbol{\theta}_{d,l} \left(\sum_{v=1}^V \boldsymbol{\beta}_{l,v} I[\mathbf{x}_{d,n} = \mathbf{v}] \right)}$$

M-step: Update parameters

$$\boldsymbol{\beta}_{k,v} = \frac{N_{k,v}^{\beta}}{\sum_{d=1}^D N_{d,k}^{\theta}} \quad N_{k,v}^{\beta} = \sum_{d=1}^D \sum_{n=1}^{N_d} \phi_{d,n,k} I[\mathbf{x}_{d,n} = \mathbf{v}]$$

$$\boldsymbol{\theta}_{d,k} = \frac{N_{d,k}^{\theta}}{N_d} \quad N_{d,k}^{\theta} = \sum_{n=1}^{N_d} \phi_{d,n,k}$$



*(Probabilistic Latent Semantic Indexing, a.k.a. Probabilistic Latent Semantic Analysis)



Topic Models

Jan-Willem van de Meent

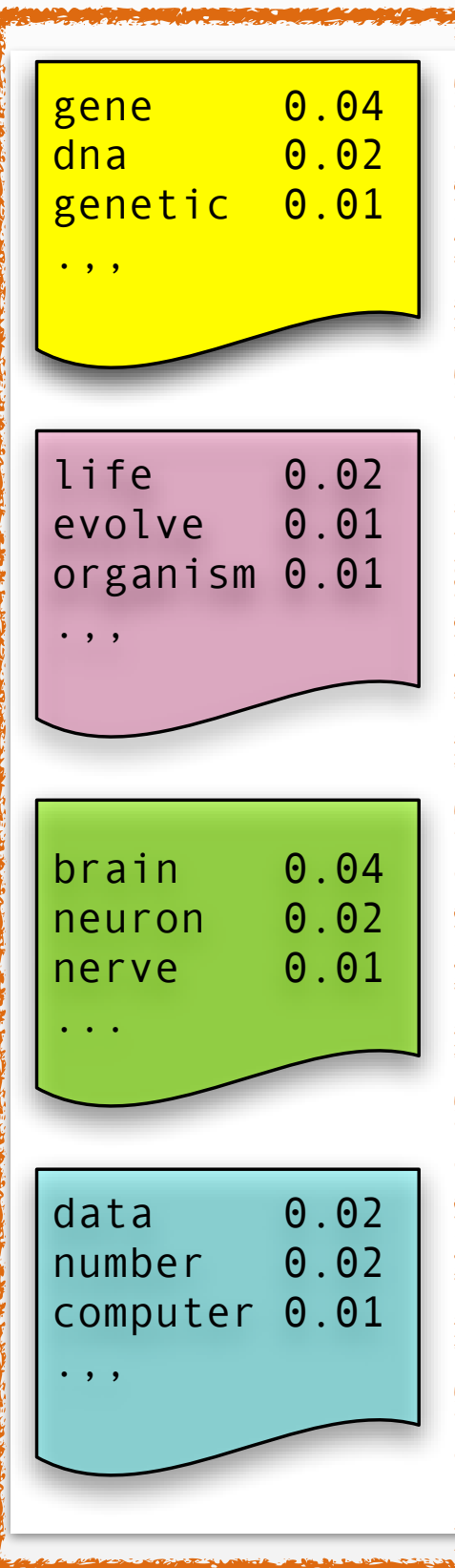


Latent Dirichlet Allocation

Topic Models with Dirichlet Priors

Review: Topic Modeling with PLSA/PLSI

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

"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 an

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

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

SCIENCE • VOL. 272 • 24 MAY 1996

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

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

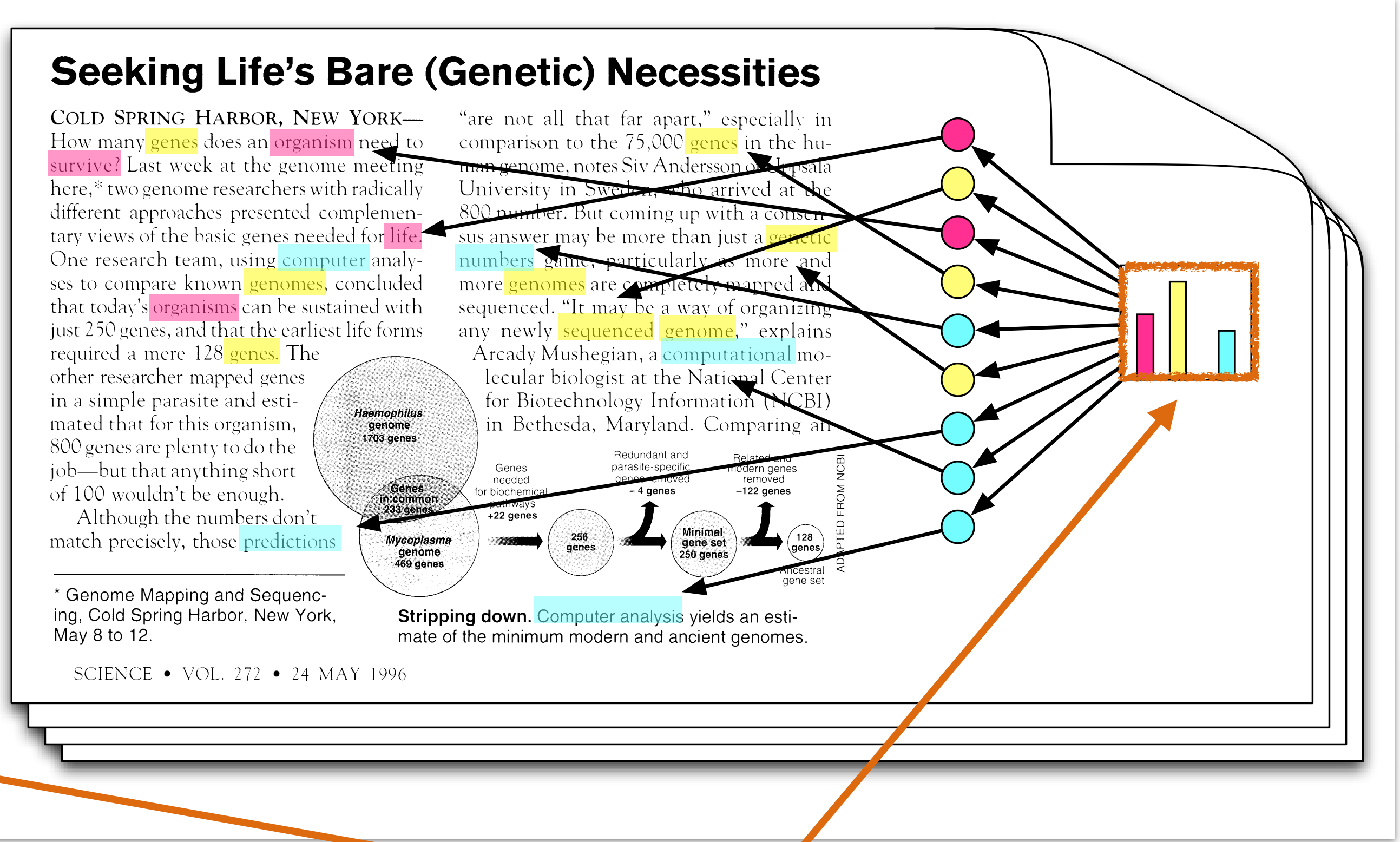
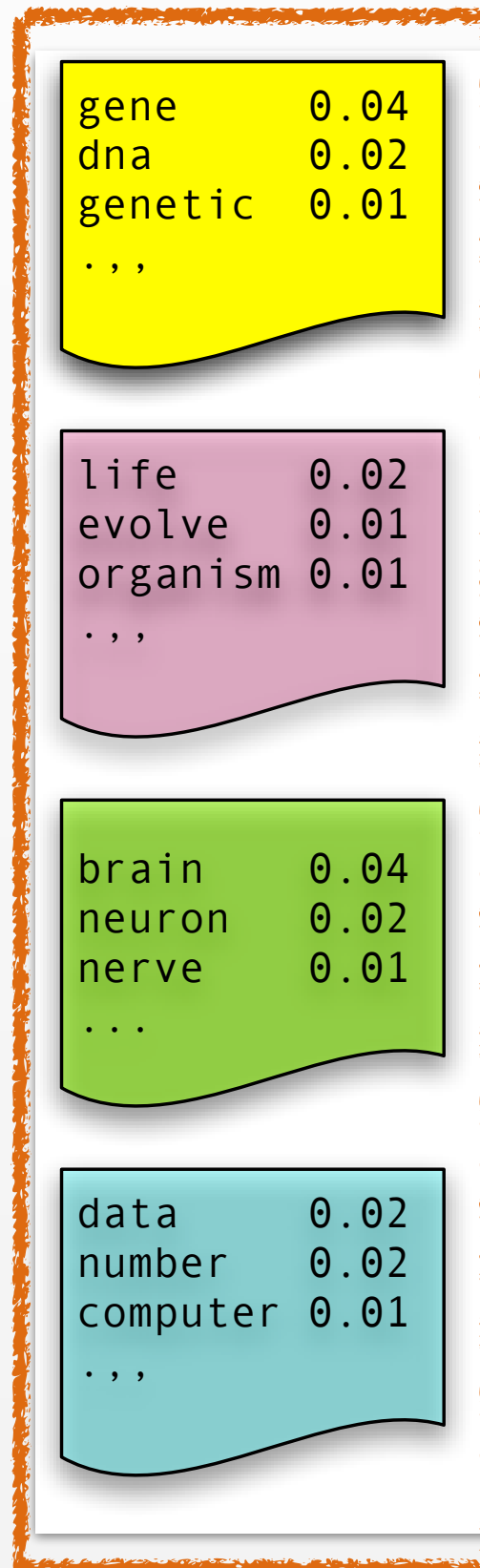
LDA: Add *Dirichlet* Priors

β_k : Topics
(shared)

x_d : Words

z_d : Assignments
(document-specific)

θ_d : Topic Proportions



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

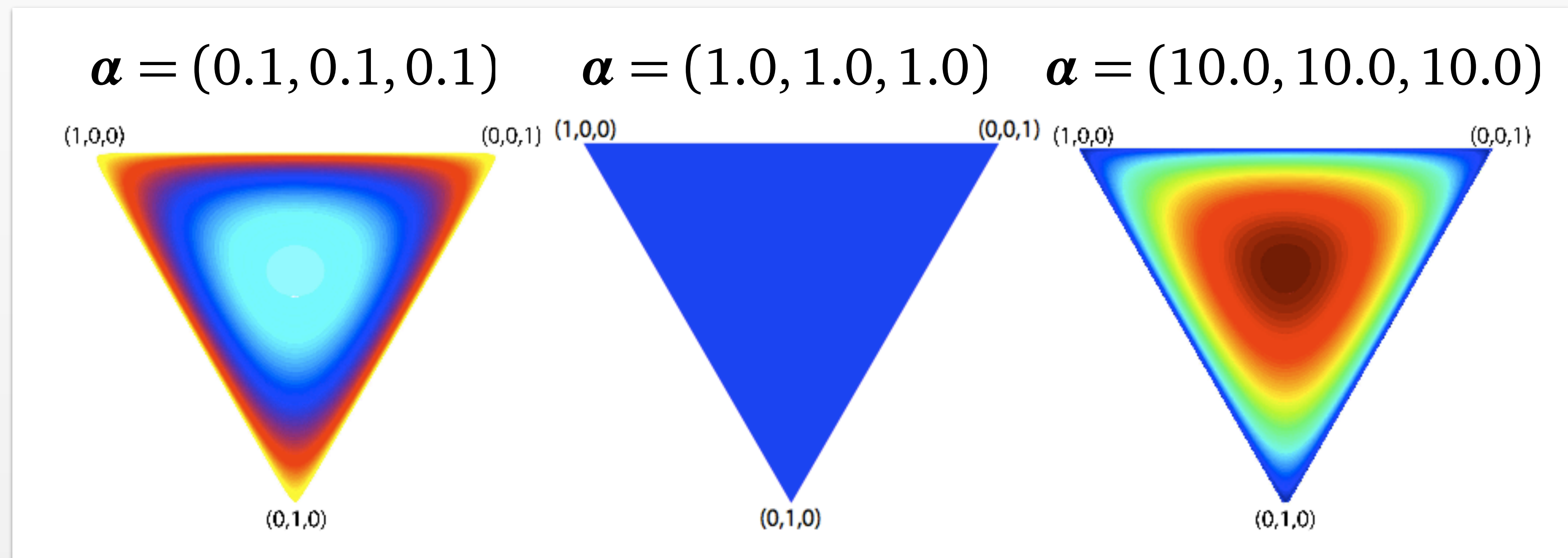
$$x_{d,n} \mid z_{d,n} = k \sim \text{Discrete}(\beta_k)$$

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

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

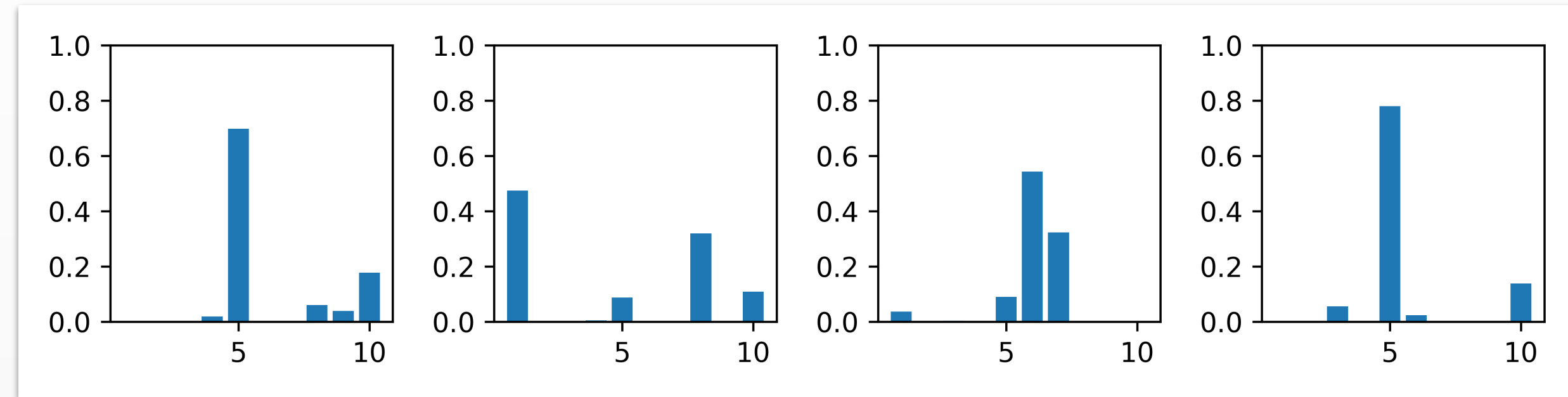
Review: Dirichlet Distribution

$$p(\boldsymbol{\theta}) = \frac{1}{B(\boldsymbol{\alpha})} \prod_{k=1}^K \theta_k^{\alpha_k - 1} \quad B(\boldsymbol{\alpha}) := \frac{\prod_{k=1}^K \Gamma(\alpha_k)}{\Gamma\left(\sum_{k=1}^K \alpha_k\right)}$$

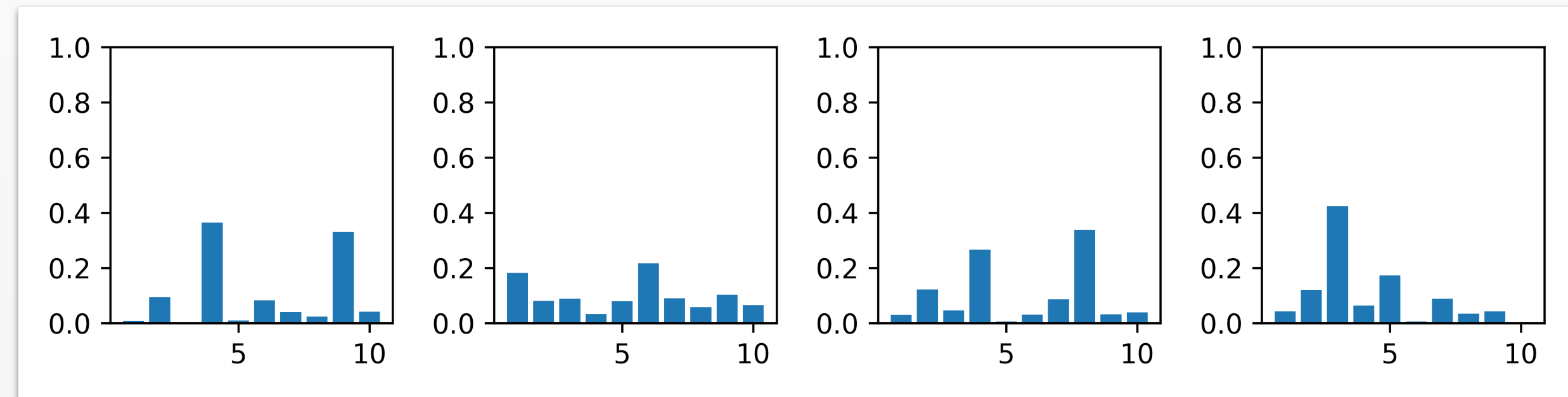


Review: Dirichlet Distribution

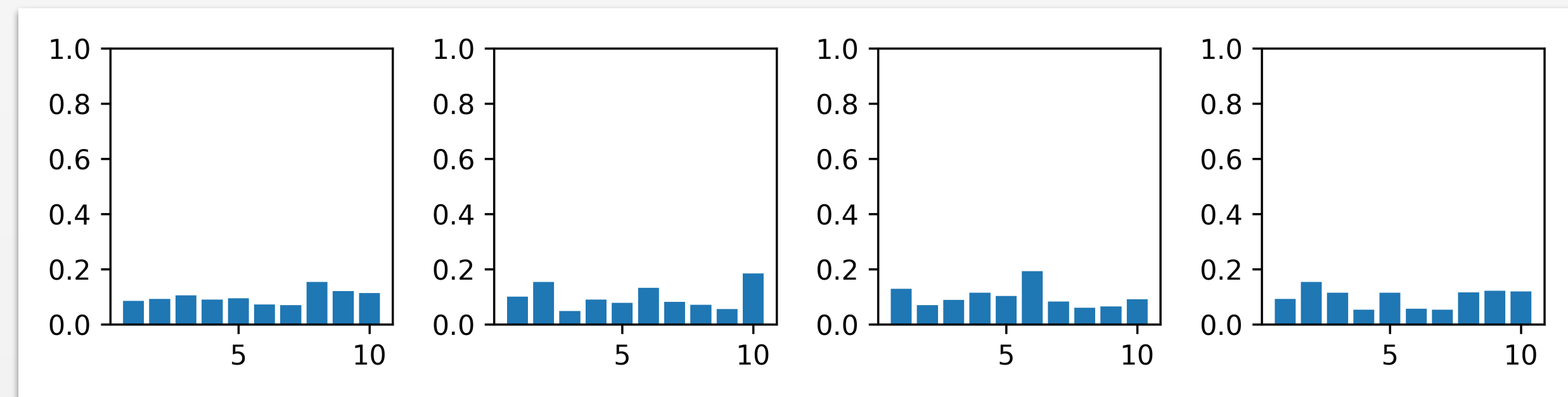
$$\alpha_k = 0.1$$



$$\alpha_k = 1.0$$



$$\alpha_k = 10.0$$



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

LDA: *Summary so far*

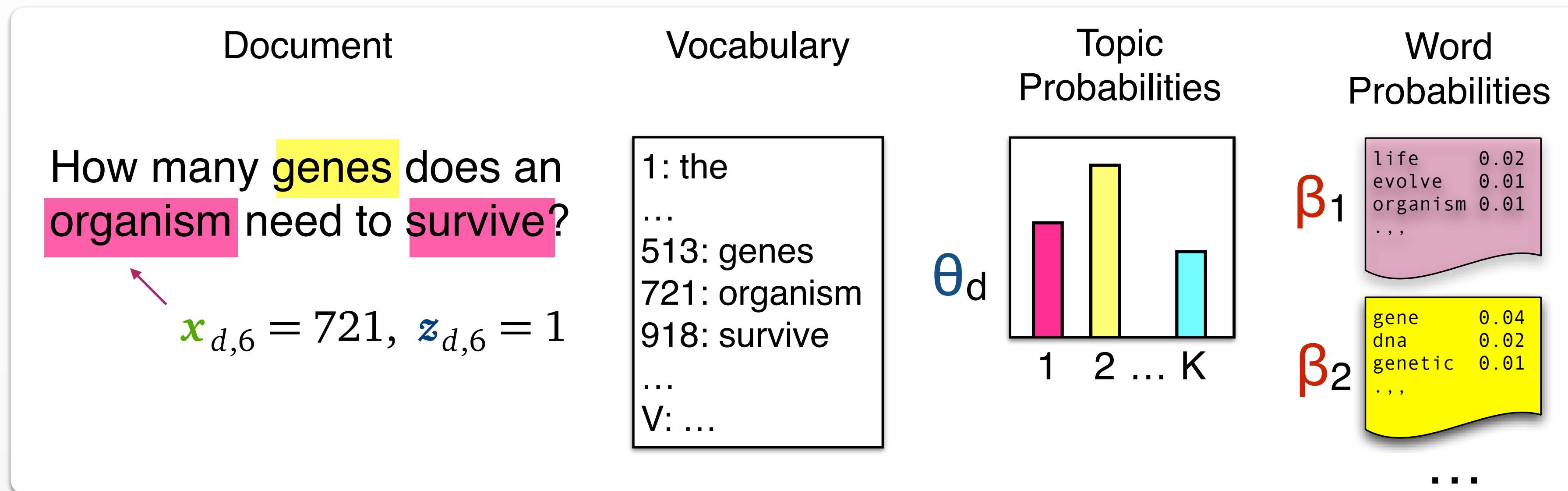
- **Idea:** Model documents as *mixtures* over topics
- **Model parameters:**
 - θ_d Topic probabilities for each document
(K-dimensional vector for each document)
 - β_k Word probabilities for each topic
(V-dimensional vector for each topic)
- **Interpretation Dimensionality Reduction:**
Similar to LSA, but assumes Discrete mixture instead of Gaussian distribution on word counts
- **Dirichlet Priors:** Enforce sparsity, associate a small number of topics which each document

Estimating Model Parameters

Question: How can we estimate β_k and θ_d ?

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

Estimating the Parameters



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

Maximum a Posteriori: $\max_{\theta, \beta} \log p(\theta, \beta \mid \mathbf{x})$

Review: Conjugate Priors for Coin Flips



Likelihood

$$\text{Bern}(x \mid \mu) = \mu^{N_1} (1 - \mu)^{N_0}$$

Sufficient Statistics

$$N_1 = \sum_{n=1}^N x_n, \quad N_0 = N - N_1$$

Conjugate Prior

$$\text{Beta}(\mu \mid a, b) = \frac{1}{B(a, b)} \mu^{a-1} (1 - \mu)^{b-1}$$

Posterior

$$\text{Beta}(\mu \mid N_1 + a, N_0 + b)$$

MAP Estimate

$$\mu^* = \frac{N_1 + a - 1}{N + a + b - 2}$$

Generalization: Dirichlet and Discrete



Likelihood

$$p(z | \theta) = \prod_n p(z_n | \theta) = \prod_k \theta_k^{N_k}$$

Sufficient Statistics

$$N_k = \sum_n I[z_n = k]$$

Conjugate Prior

$$\text{Dir}(\theta | \alpha_1, \dots, \alpha_K) = \frac{1}{B(\alpha)} \prod_{k=1}^K \theta_k^{\alpha_k - 1}$$

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

Posterior

$$\text{Dir}(\theta | N_1 + \alpha_1, \dots, N_K + \alpha_K)$$

MAP Estimate

$$\theta_k^* = \frac{N_k + \alpha_k - 1}{\sum_k N_k + \sum_k \alpha_k - K}$$

MAP estimation for LDA with EM

Generative Model

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

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

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

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

E-step: Update assignments

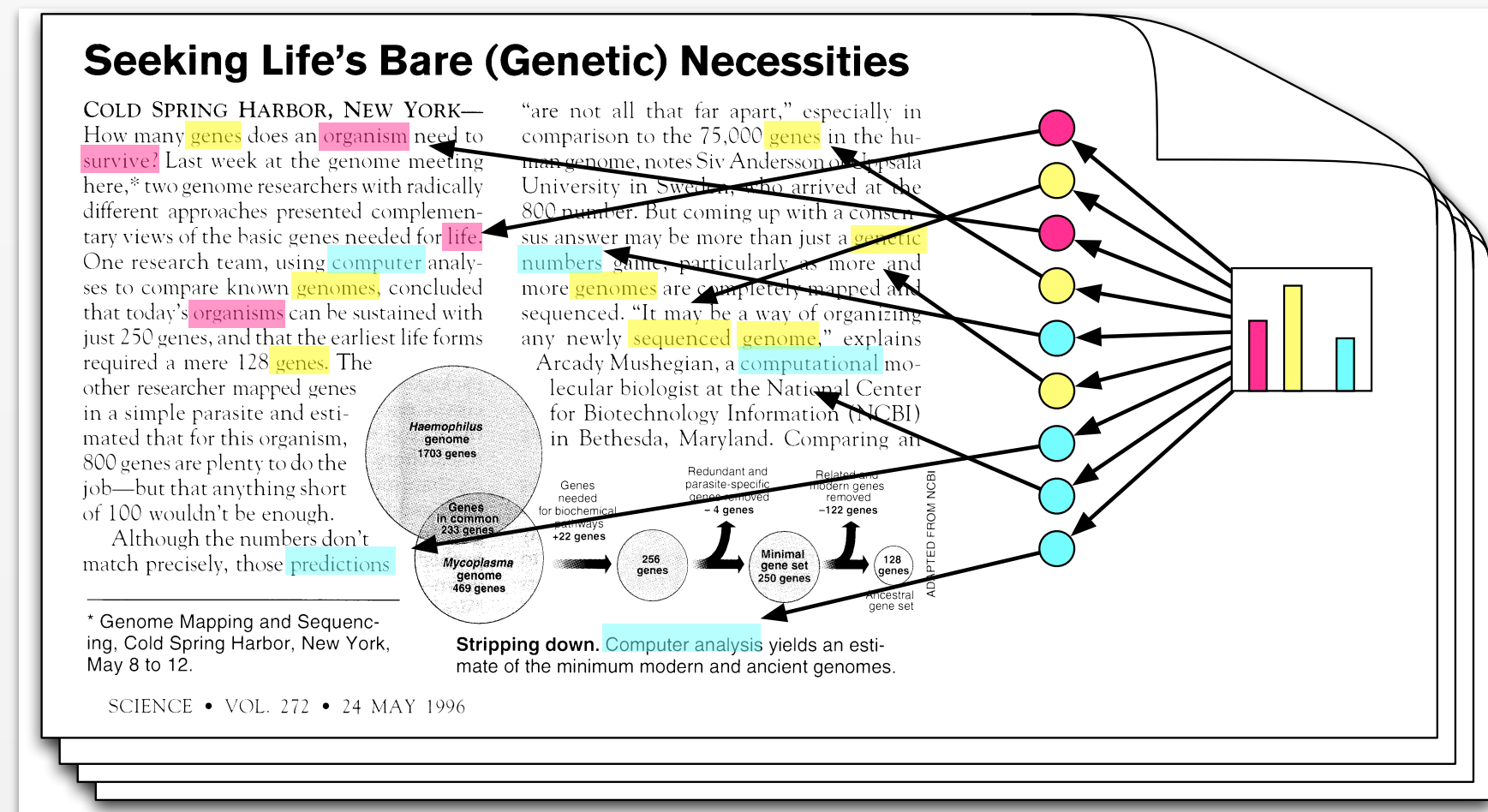
$$\phi_{d,n,k} = p(z_{d,n} = k \mid x_{d,n} = v, \beta, \theta_d)$$

$$\begin{aligned} & \theta_{d,k} \left(\sum_{v=1}^V \beta_{k,v} I[x_{d,n} = v] \right) \\ &= \frac{\theta_{d,k} \left(\sum_{v=1}^V \beta_{k,v} I[x_{d,n} = v] \right)}{\sum_{l=1}^K \theta_{d,l} \left(\sum_{v=1}^V \beta_{l,v} I[x_{d,n} = v] \right)} \end{aligned}$$

M-step: Update parameters

$$\beta_{k,v} = \frac{N_{k,v}^{\beta} + \eta_{k,v} - 1}{\sum_{d=1}^D N_{d,k}^{\theta} + \sum_v \eta_{k,v} - V}$$

$$\theta_{d,k} = \frac{N_{d,k}^{\theta} + \alpha_k - 1}{N_d + \sum_k \alpha_k - K}$$



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

Variational Expectation Maximization (high-level)

Idea: Approximate $p(\mathbf{z}, \boldsymbol{\theta}, \boldsymbol{\beta} \mid \mathbf{x})$ with $q(\mathbf{z}, \boldsymbol{\theta}, \boldsymbol{\beta})$

$$\boldsymbol{\phi}, \boldsymbol{\gamma}, \boldsymbol{\lambda} = \underset{\boldsymbol{\phi}, \boldsymbol{\gamma}, \boldsymbol{\lambda}}{\operatorname{argmin}} \operatorname{KL}(q(\mathbf{z}, \boldsymbol{\theta}, \boldsymbol{\beta}) \parallel p(\mathbf{z}, \boldsymbol{\theta}, \boldsymbol{\beta} \mid \mathbf{x}))$$

$$q(\mathbf{z}, \boldsymbol{\theta}, \boldsymbol{\beta}) = q(\mathbf{z}; \boldsymbol{\phi}) q(\boldsymbol{\theta}; \boldsymbol{\gamma}) q(\boldsymbol{\beta}; \boldsymbol{\lambda})$$

Discrete Dirichlet Dirichlet

Variational E-step: Update $\boldsymbol{\phi}$

$$\boldsymbol{\phi}_{d,n,k} = \exp \left(\mathbb{E}_q \left[\log \boldsymbol{\theta}_{d,k} + \sum_{v=1}^V I[\mathbf{x}_{d,n} = v] \log \boldsymbol{\beta}_{k,v} \right] \right)$$

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

Variational M-step: Update $\boldsymbol{\gamma}$ and $\boldsymbol{\lambda}$

$$\boldsymbol{\gamma}_{d,k} = \boldsymbol{\alpha}_k + N_{d,k}^{\boldsymbol{\theta}} \quad \boldsymbol{\lambda}_{k,v} = \boldsymbol{\eta}_{k,v} + N_{k,v}^{\boldsymbol{\beta}}$$

(analogous to MAP estimation – need to know this)

EM vs. Variational EM

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

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

$$\text{M-step:} \quad \theta_{d,k} = \frac{N_{d,k}^{\theta}}{N_d} \quad \beta_{k,v} = \frac{N_{k,v}^{\beta}}{\sum_{d=1}^D N_{d,k}^{\theta}}$$

$$\text{Variational EM} \quad \phi, \gamma, \lambda = \underset{\phi, \gamma, \lambda}{\operatorname{argmin}} \operatorname{KL}(q(\mathbf{z}, \theta, \beta) || p(\mathbf{z}, \theta, \beta | \mathbf{x}))$$

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

$$\text{M-step:} \quad \gamma_{d,k} = \alpha_k + N_{d,k}^{\theta} \quad \lambda_{k,v} = \eta_{k,v} + N_{k,v}^{\beta}$$

EM vs. Variational EM

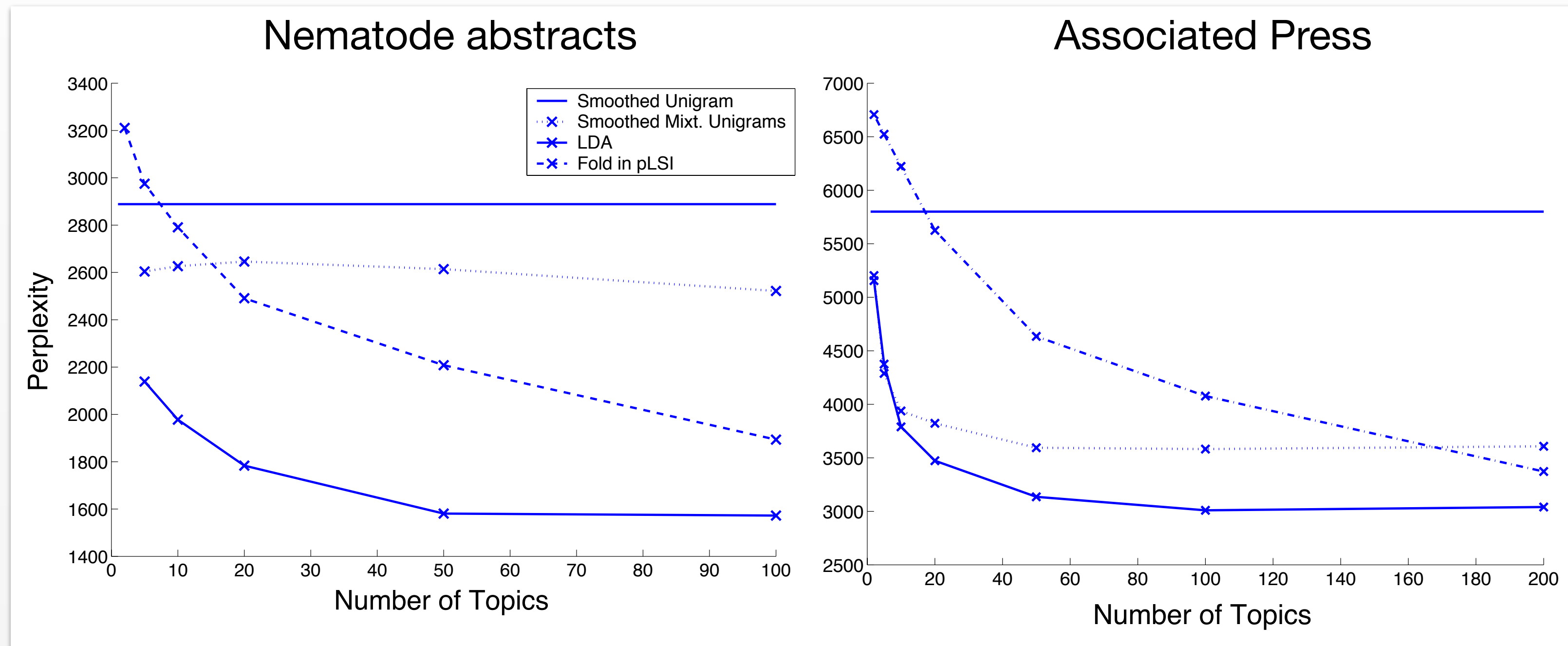
Commonalities: Both compute sufficient statistics

$\phi_{d,n,k}$	Probability that word n in document d belongs to topic k
$N_{d,k}^\theta$	Number of words in document d that belong to topic k
$N_{k,v}^\beta$	Number of times word v appears in topic k (across <u>all</u> 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 <i>Posterior</i> 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

Jan-Willem van de Meent



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](#)

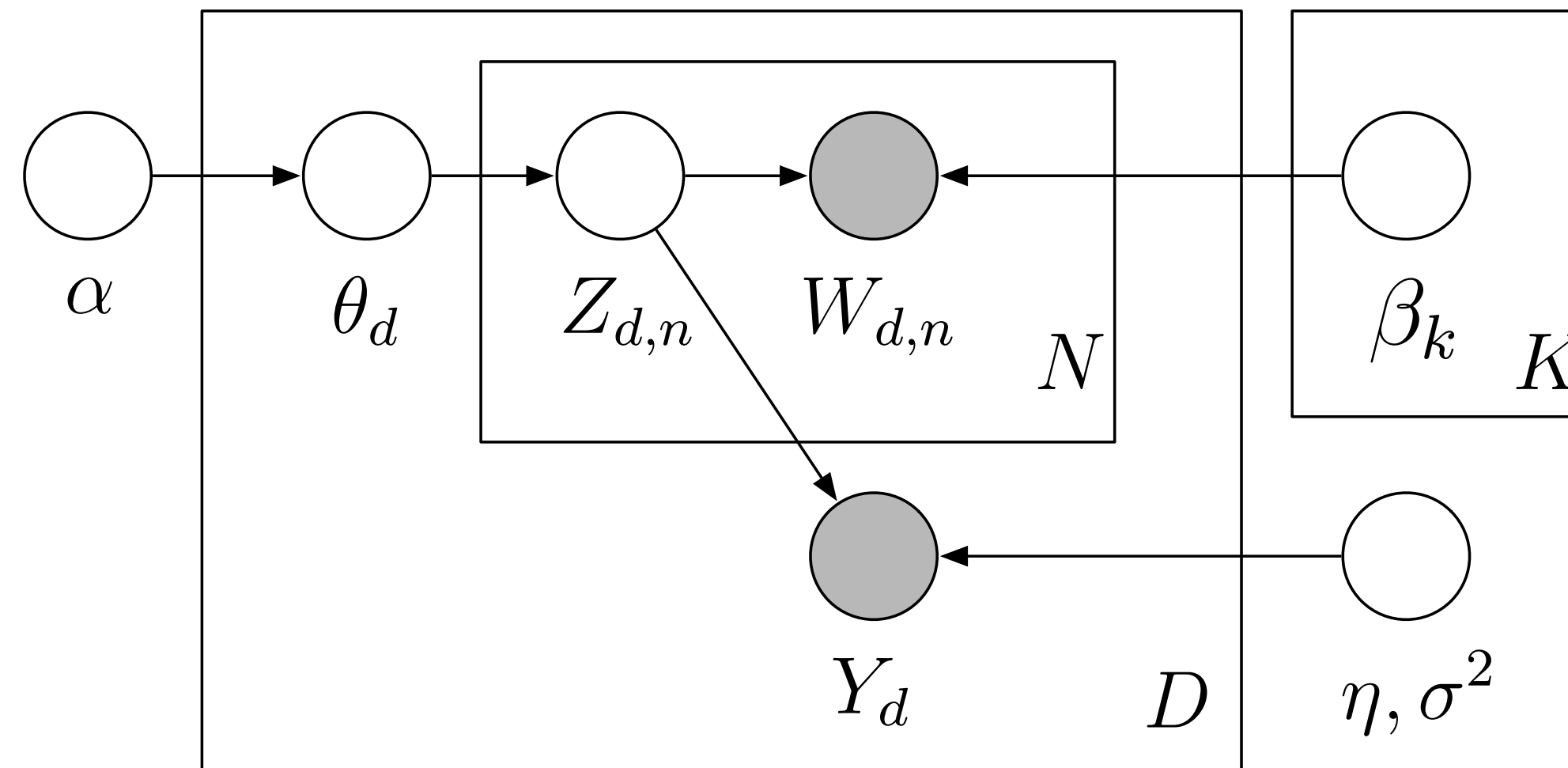
Abstract We describe latent Dirichlet allocation (LDA), a generative probabilistic model for collections of discrete data such as text corpora. LDA is a three-level hierarchical Bayesian model, in which each item of a collection is modeled as a finite mixture over an underlying ...

[Cited by 15971](#) [Related articles](#) [All 124 versions](#) [Cite](#) [Save](#)

Reasons for popularity of LDA:

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

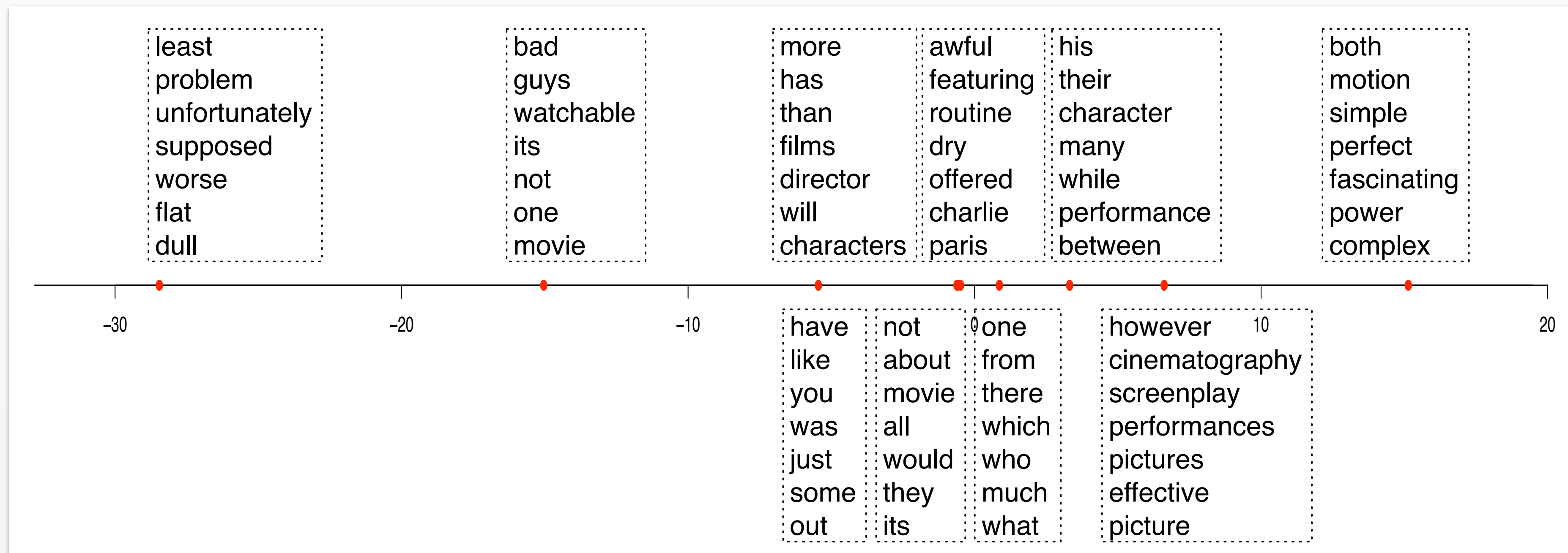
Extensions: Supervised LDA



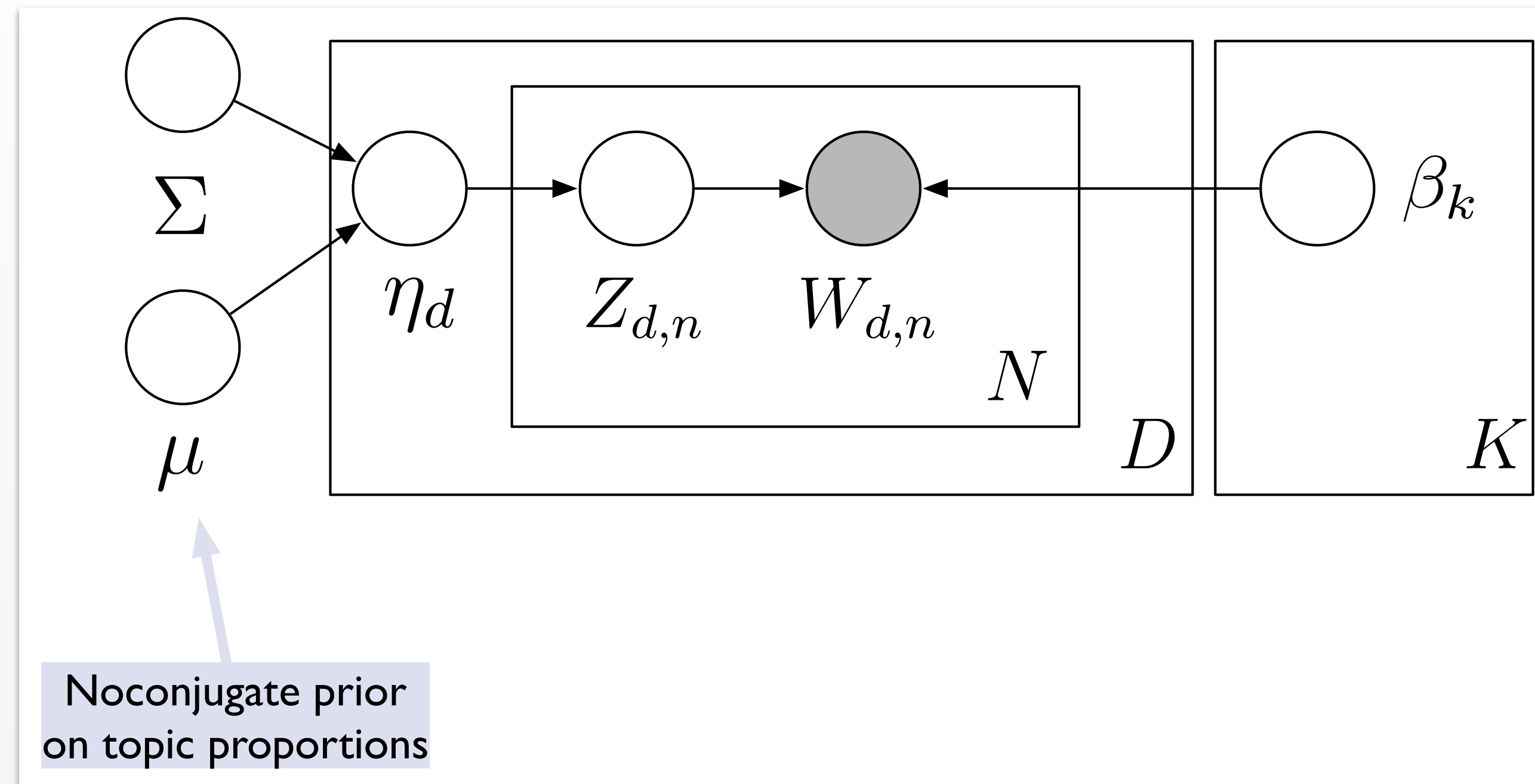
- 1 Draw topic proportions $\theta \mid \alpha \sim \text{Dir}(\alpha)$.
- 2 For each word
 - Draw topic assignment $z_n \mid \theta \sim \text{Mult}(\theta)$.
 - Draw word $w_n \mid z_n, \beta_{1:K} \sim \text{Mult}(\beta_{z_n})$.
- 3 Draw response variable $y \mid z_{1:N}, \eta, \sigma^2 \sim \text{N}(\eta^\top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

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

Extensions: Supervised LDA

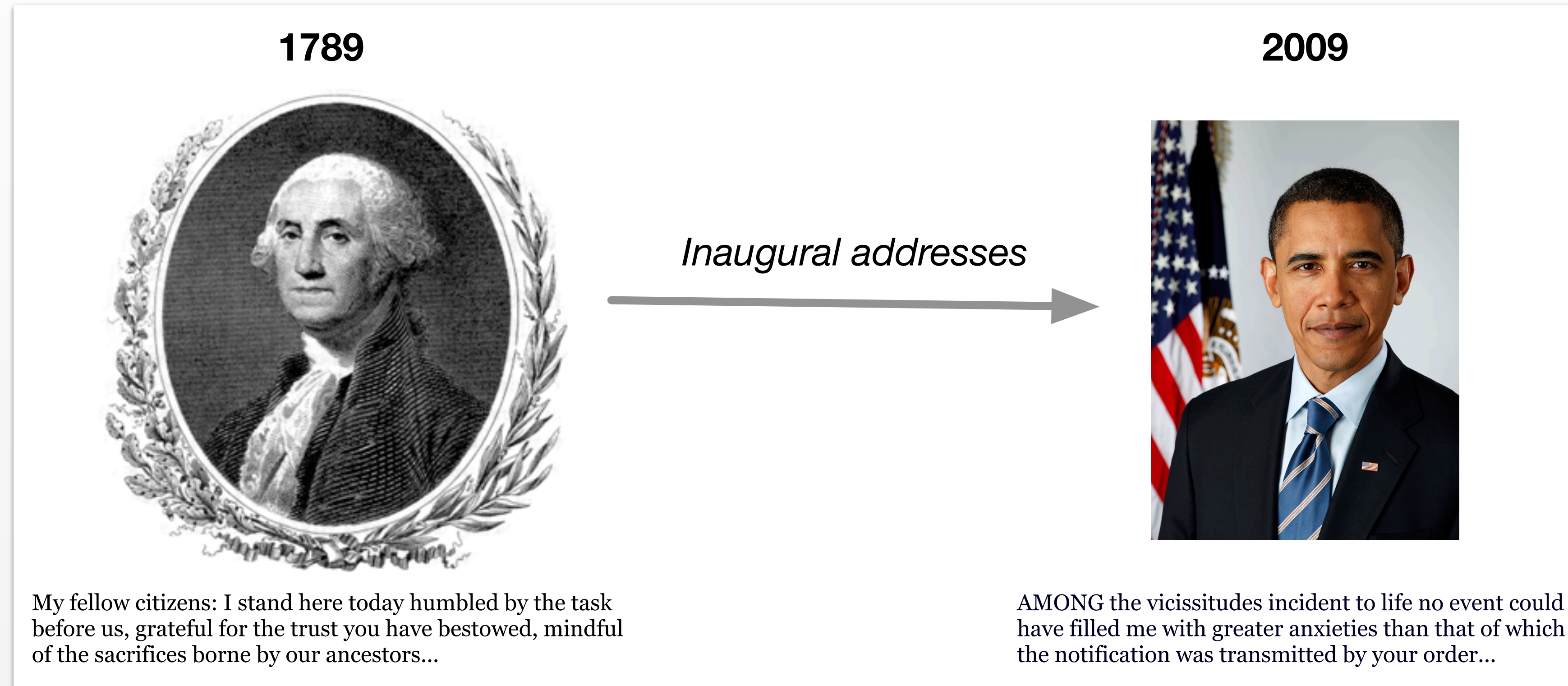


Extensions: Correlated Topic Model



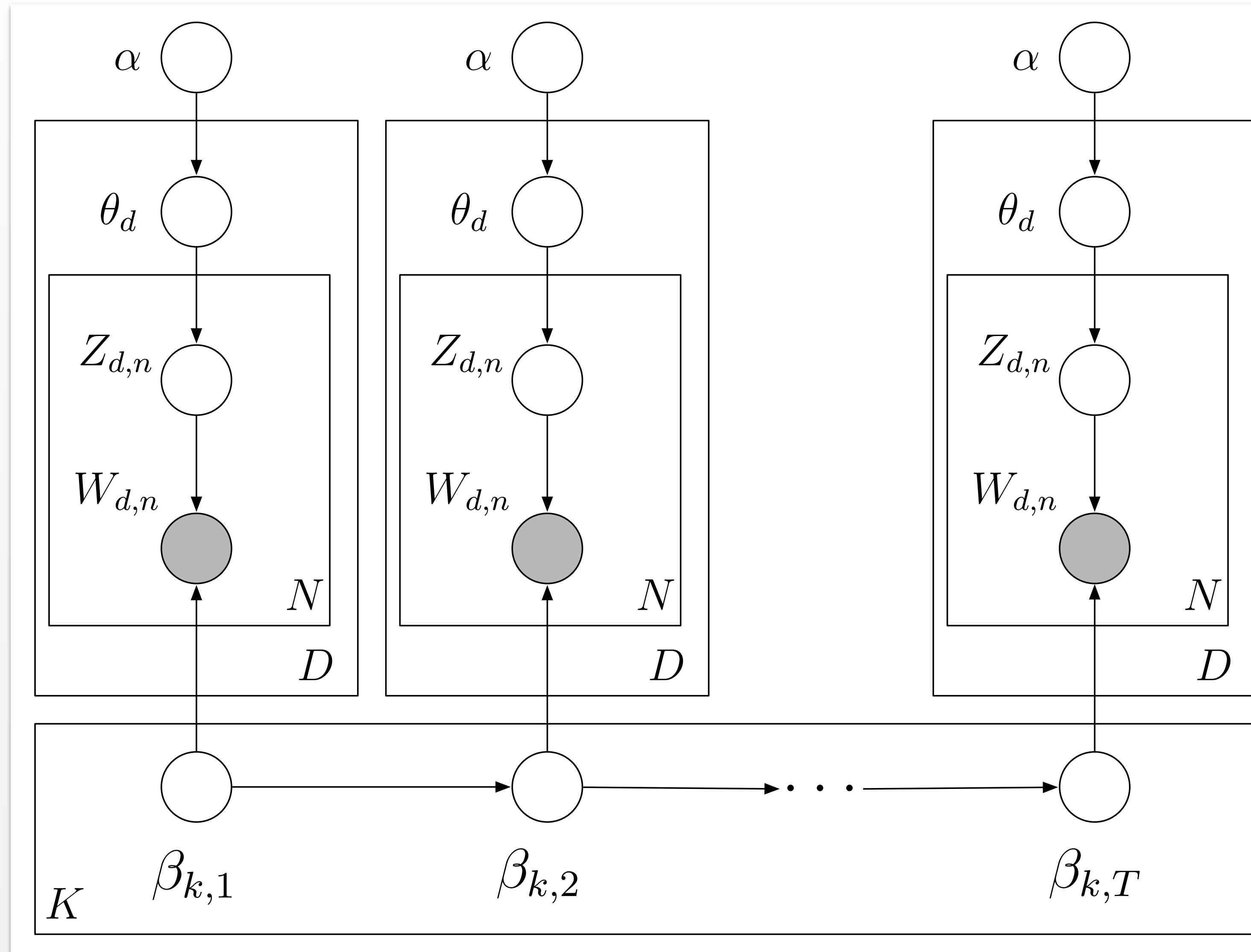
Estimate a covariance matrix Σ that parameterizes correlations between topics in a document

Extensions: Dynamic Topic Models

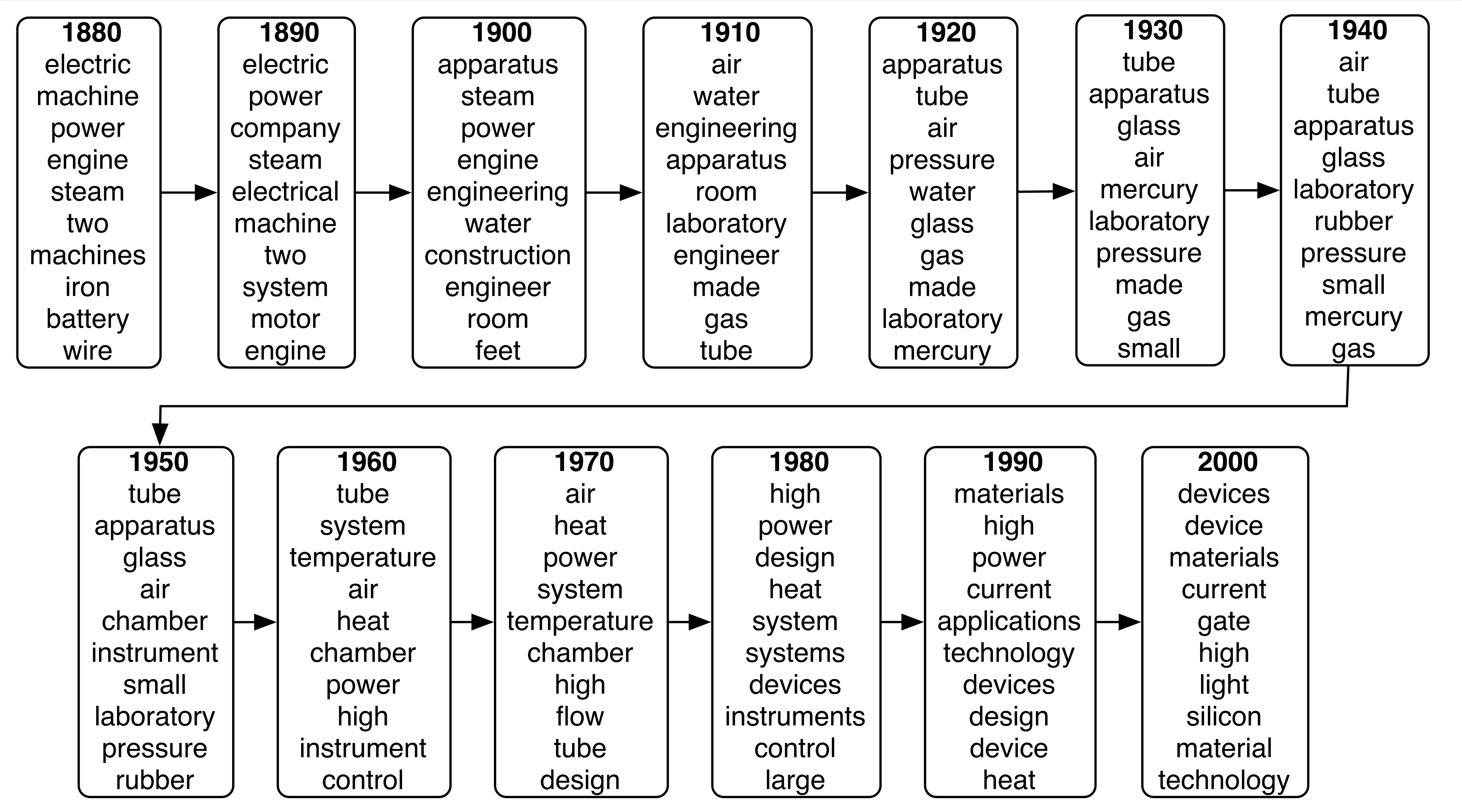


Track changes in word distributions associated with a topic over time.

Extensions: Dynamic Topic Models

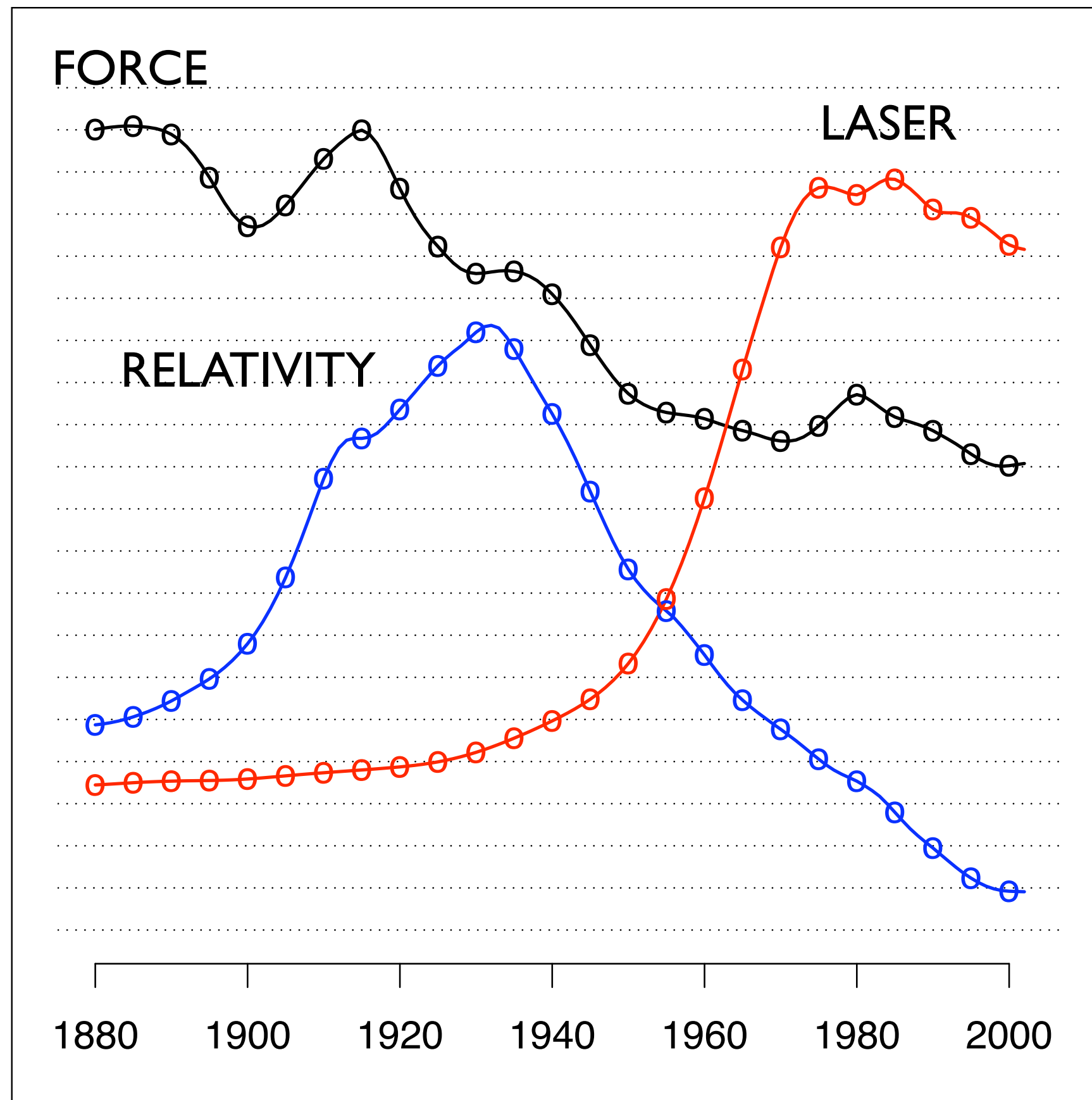


Extensions: Dynamic Topic Models

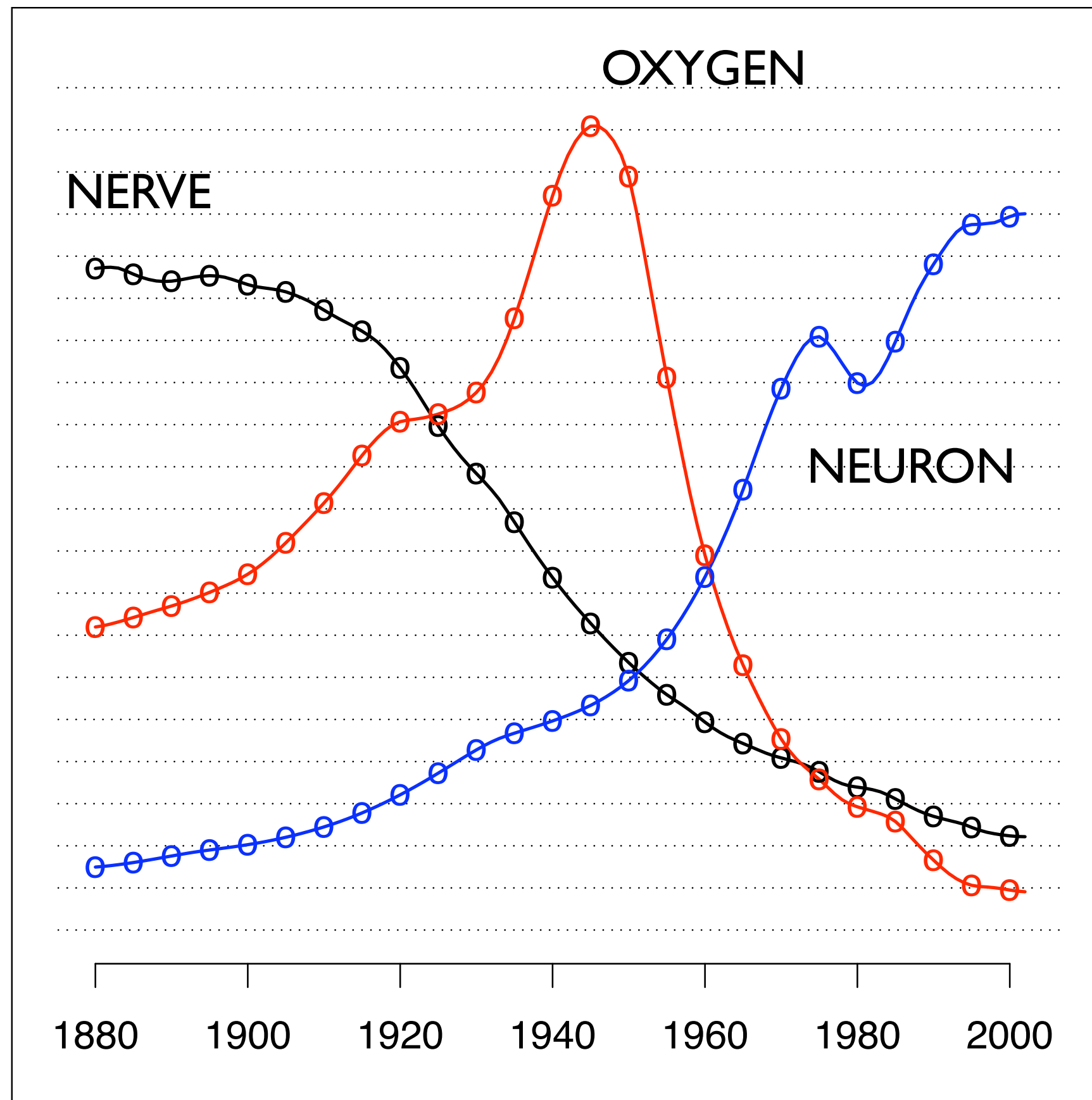


Extensions: Dynamic Topic Models

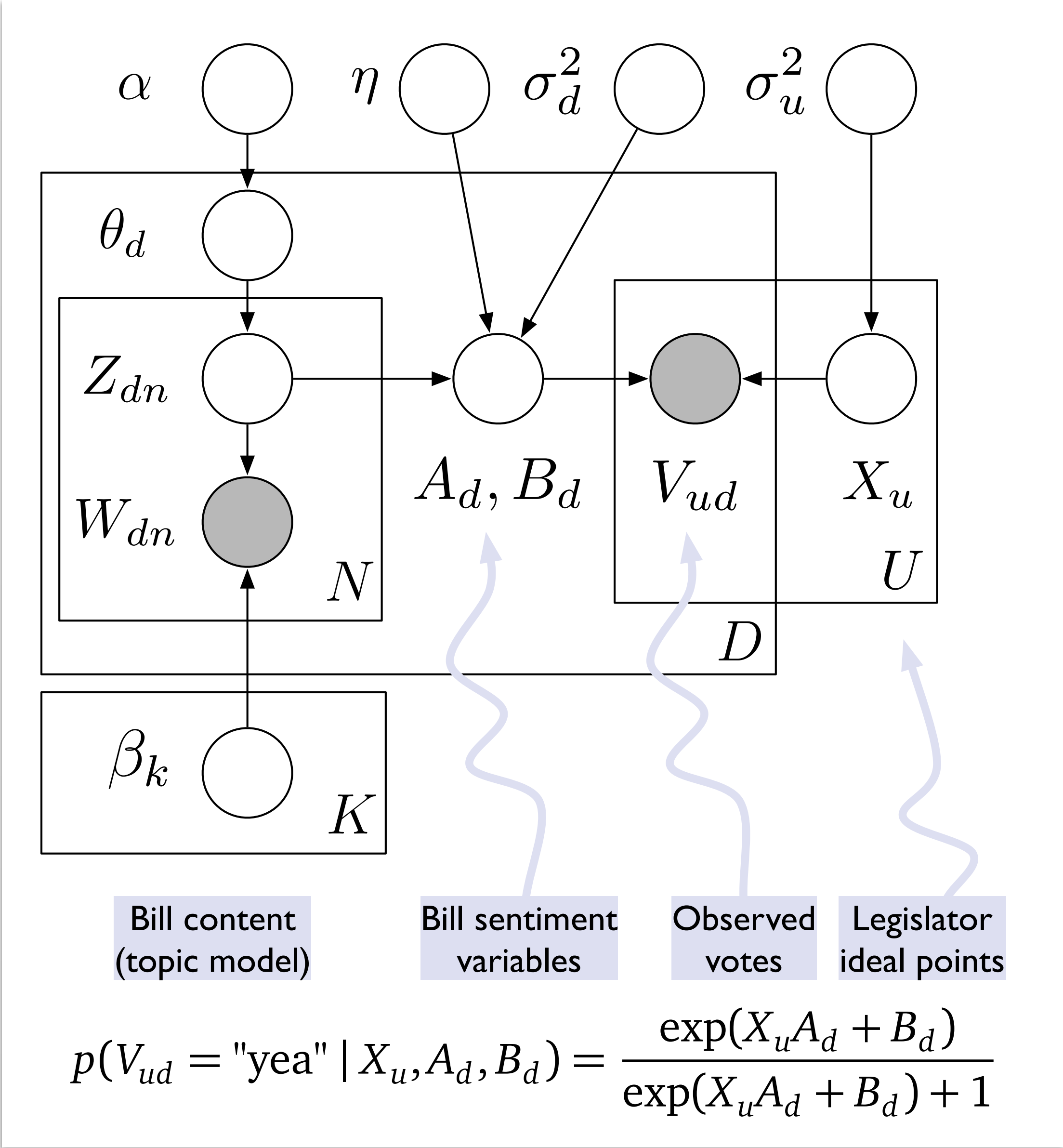
"Theoretical Physics"



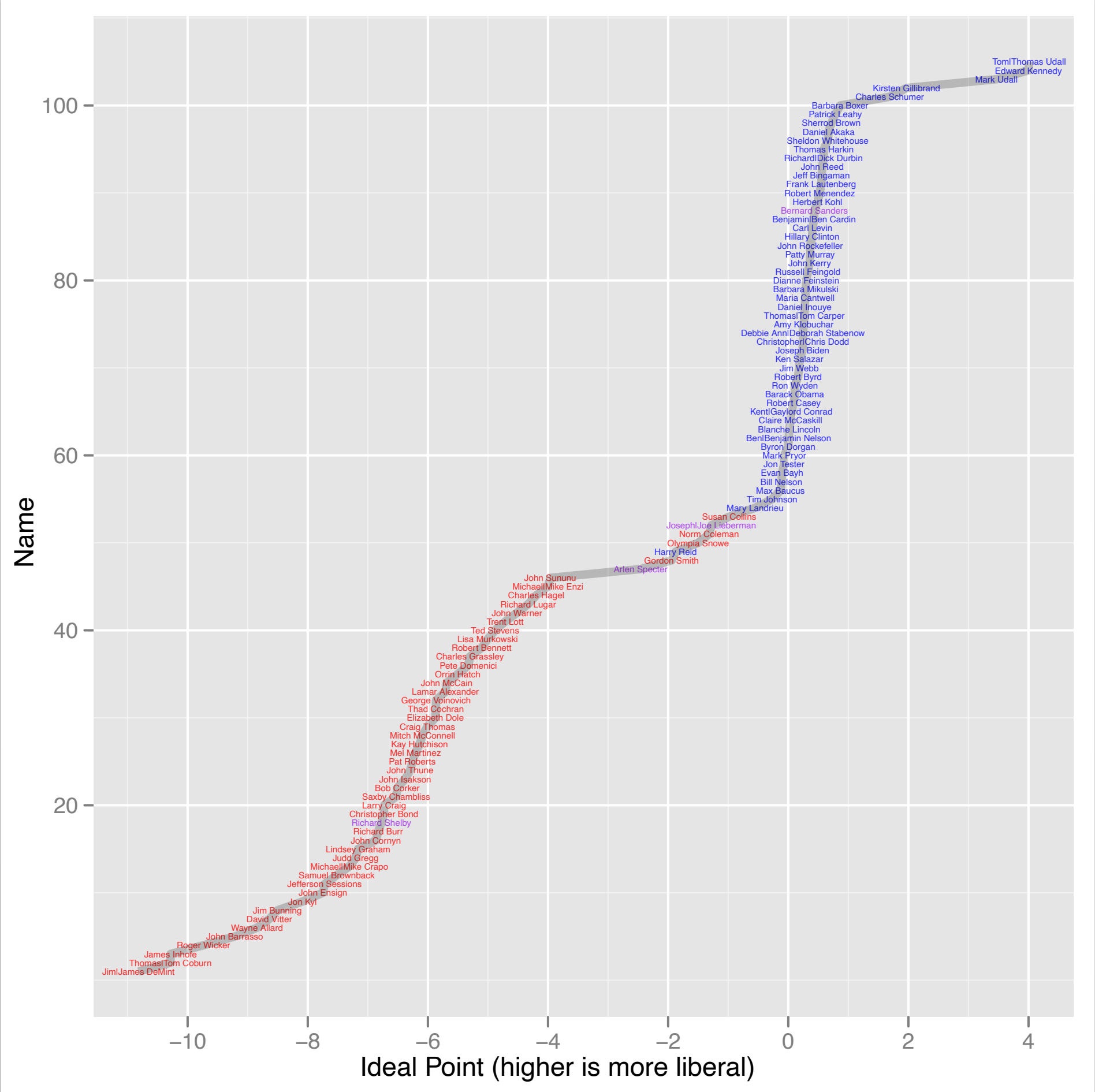
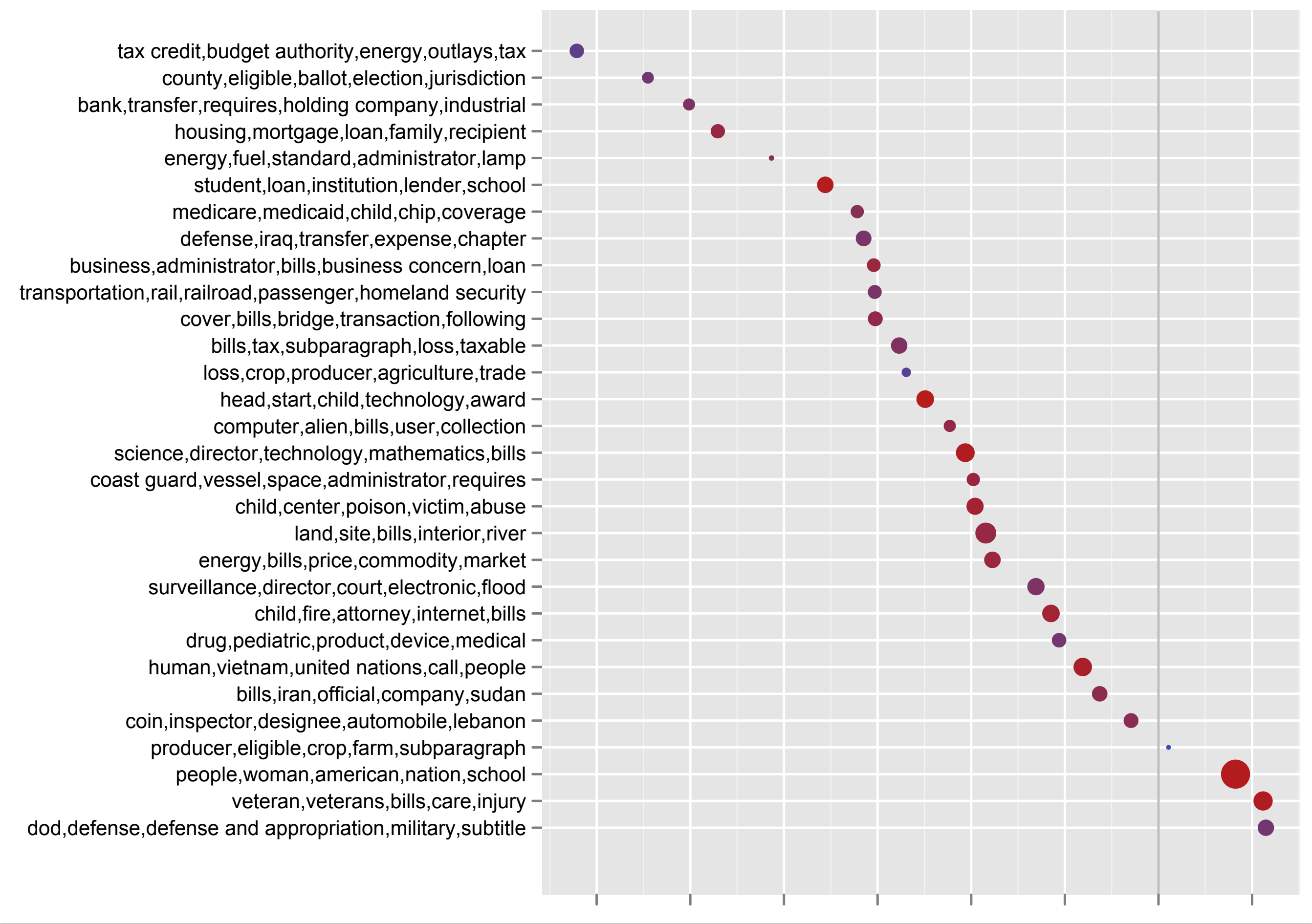
"Neuroscience"



Extensions: Ideal Point Topic Models



Extensions: Ideal Point Topic Models



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

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