



# Topic Models

Shantanu Jain



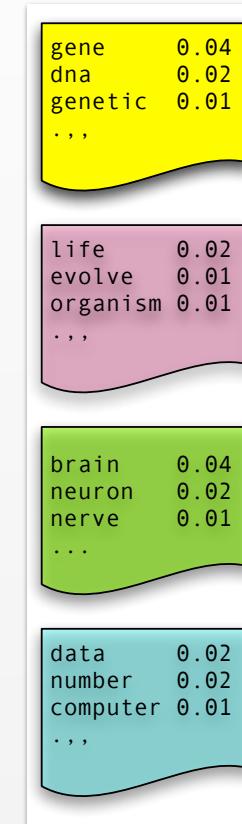
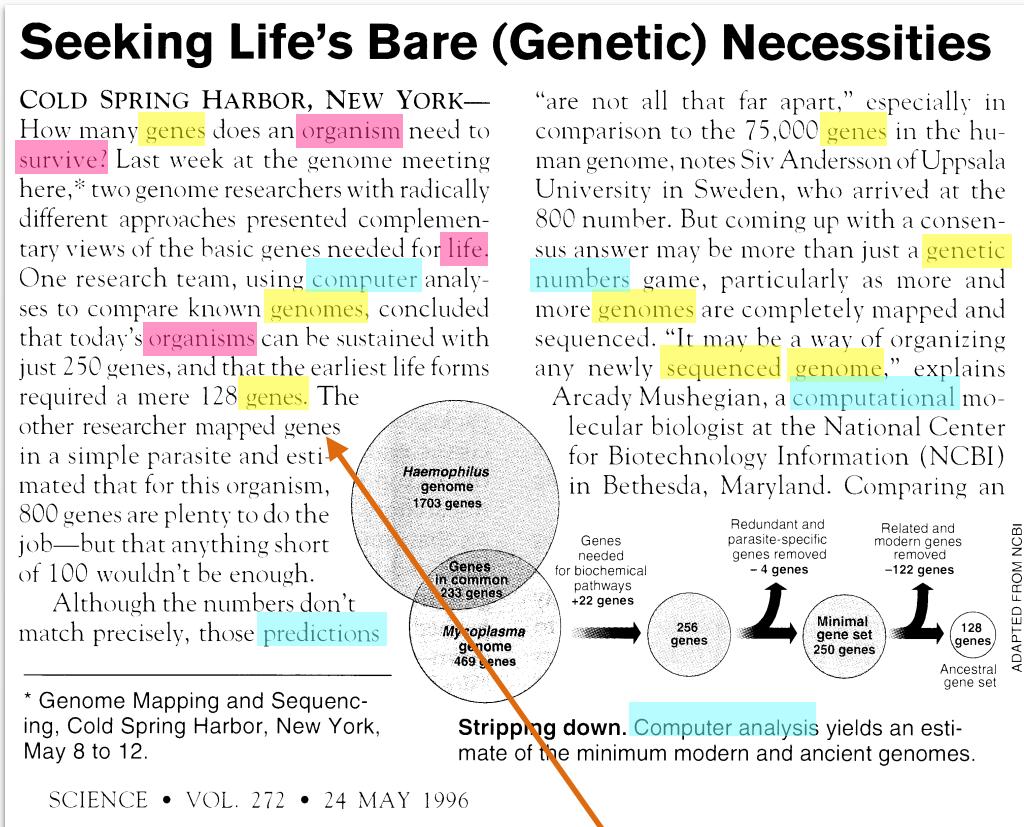
# Topic Modeling Basics



*Borrowing from:*  
David Blei  
(Columbia)

# Word Mixtures

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



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

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

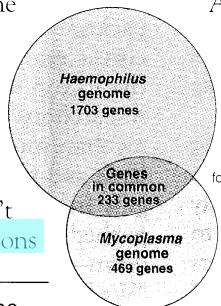
- Total  $N$  words in a document
- $n$  denotes a position in the document.
- $V$  is the number of words in the vocabulary.
- $K$  is the number of topics.

# 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.

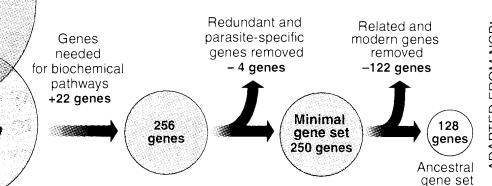
Although the numbers don't match precisely, those predictions



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

"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

ADAPTED FROM NCBI



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

$$\begin{aligned} \mathbf{z}_n &\sim \text{Discrete}(\boldsymbol{\theta}) \\ \mathbf{x}_n \mid \mathbf{z}_n = k &\sim \text{Discrete}(\boldsymbol{\beta}_k) \end{aligned}$$

*Pick a topic*  
*Pick a word given topic*

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
...	

$$p(\mathbf{x} \mid \mathbf{z}=1, \boldsymbol{\beta})$$

$$p(\mathbf{x} \mid \mathbf{z}=2, \boldsymbol{\beta})$$

$$p(\mathbf{x} \mid \mathbf{z}=3, \boldsymbol{\beta})$$

$$p(\mathbf{x} \mid \mathbf{z}=4, \boldsymbol{\beta})$$

$\boldsymbol{\theta}$ : topic proportions/probabilities, probability over the  $K$  topics

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

$$p(z_n = k \mid \boldsymbol{\theta}) = \theta_k$$

$\boldsymbol{\beta}_k$ :  $k^{th}$  topic's word probabilities over the vocabulary

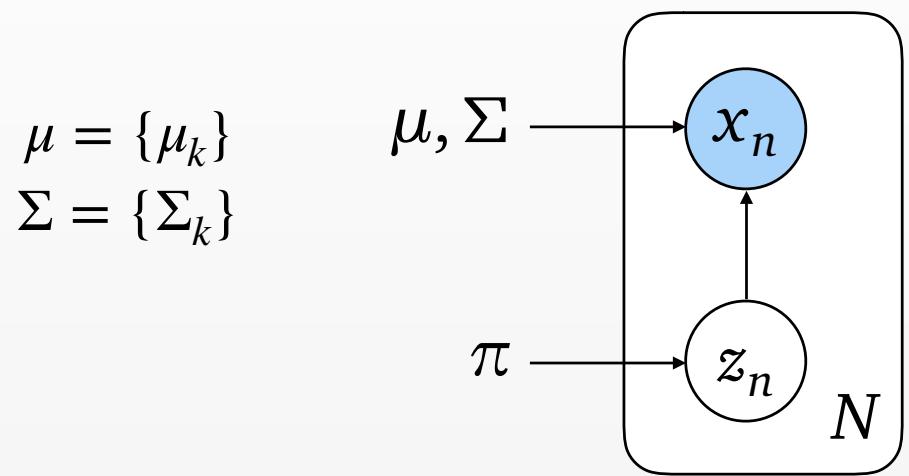
$$\boldsymbol{\beta}_k = [\beta_{k1}, \beta_{k2}, \dots, \beta_{kV}] \quad \sum_i \beta_{ki} = 1$$

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

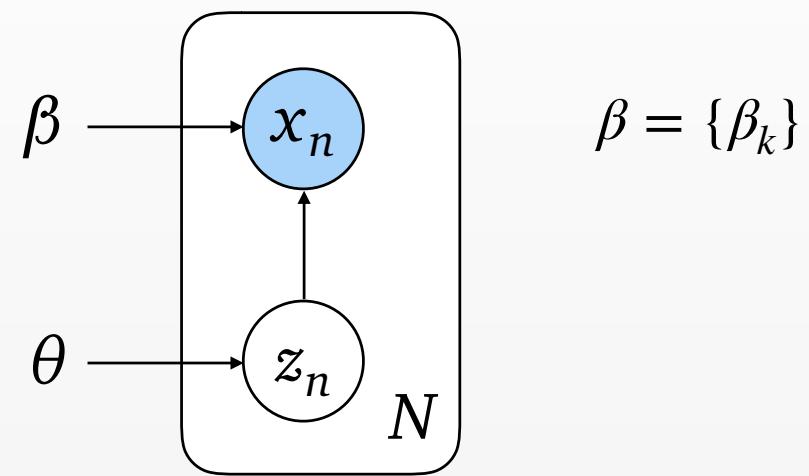
By Discrete distribution, we mean Categorical distribution

# Gaussian Mixtures vs Word Mixtures

Gaussian Mixture Model



Discrete Mixture Model



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

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

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

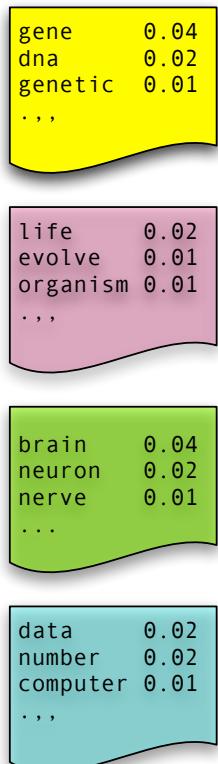
$$x_n | z_n=k \sim \text{Discrete}(\beta_{k,1}, \dots, \beta_{k,V})$$

Difference: Replace Gaussian with Discrete

Using the term Discrete distribution to mean Categorical distribution

# Topic Modeling

Topics  
(shared)



Words in Document  
(mixture over topics)

## 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

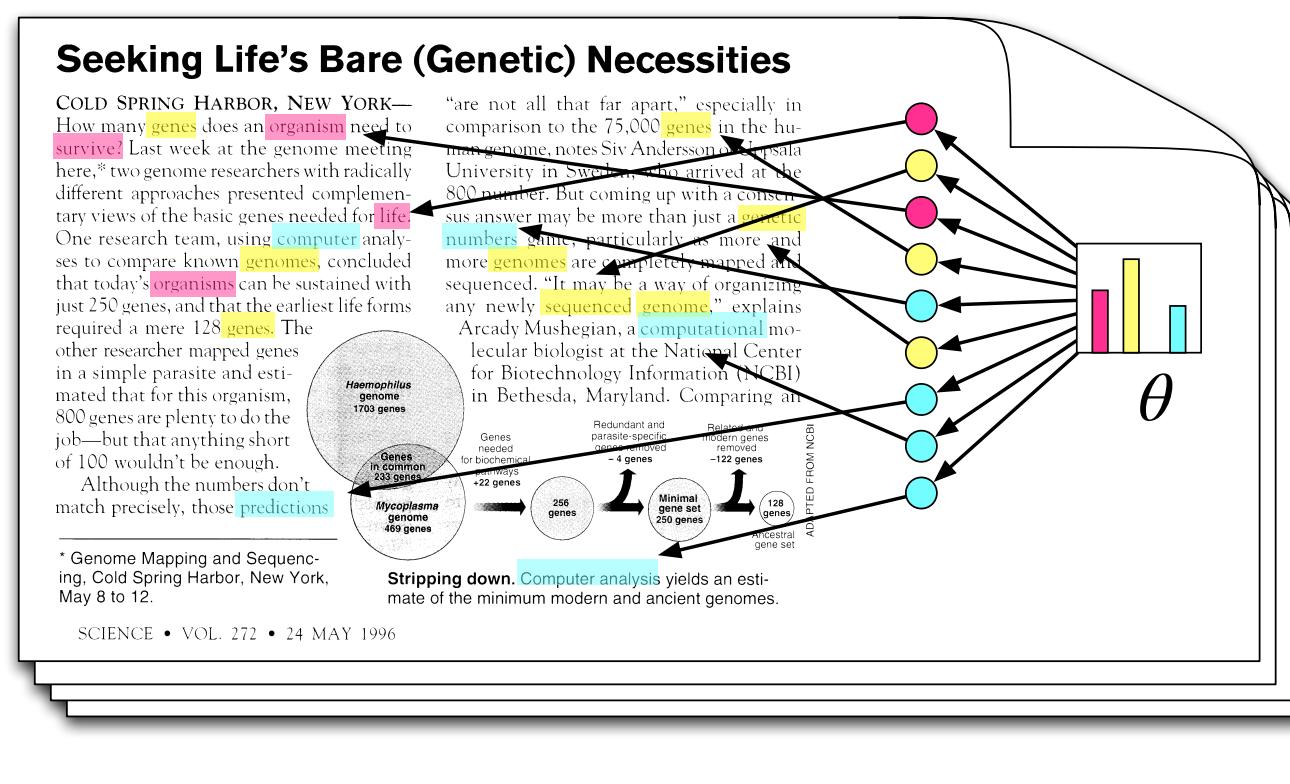
\* Genome Mapping and Sequencing, Cold Spring Harbor, New York, May 8 to 12.

"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

SCIENCE • VOL. 272 • 24 MAY 1996

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

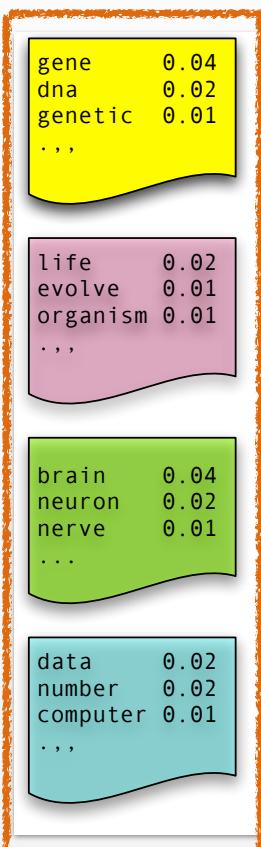
Topic Proportions  
(document-specific)



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

# Topic Modeling

$\beta_k$ : Topics



(shared across documents)

$x_d$ : Words

## 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.

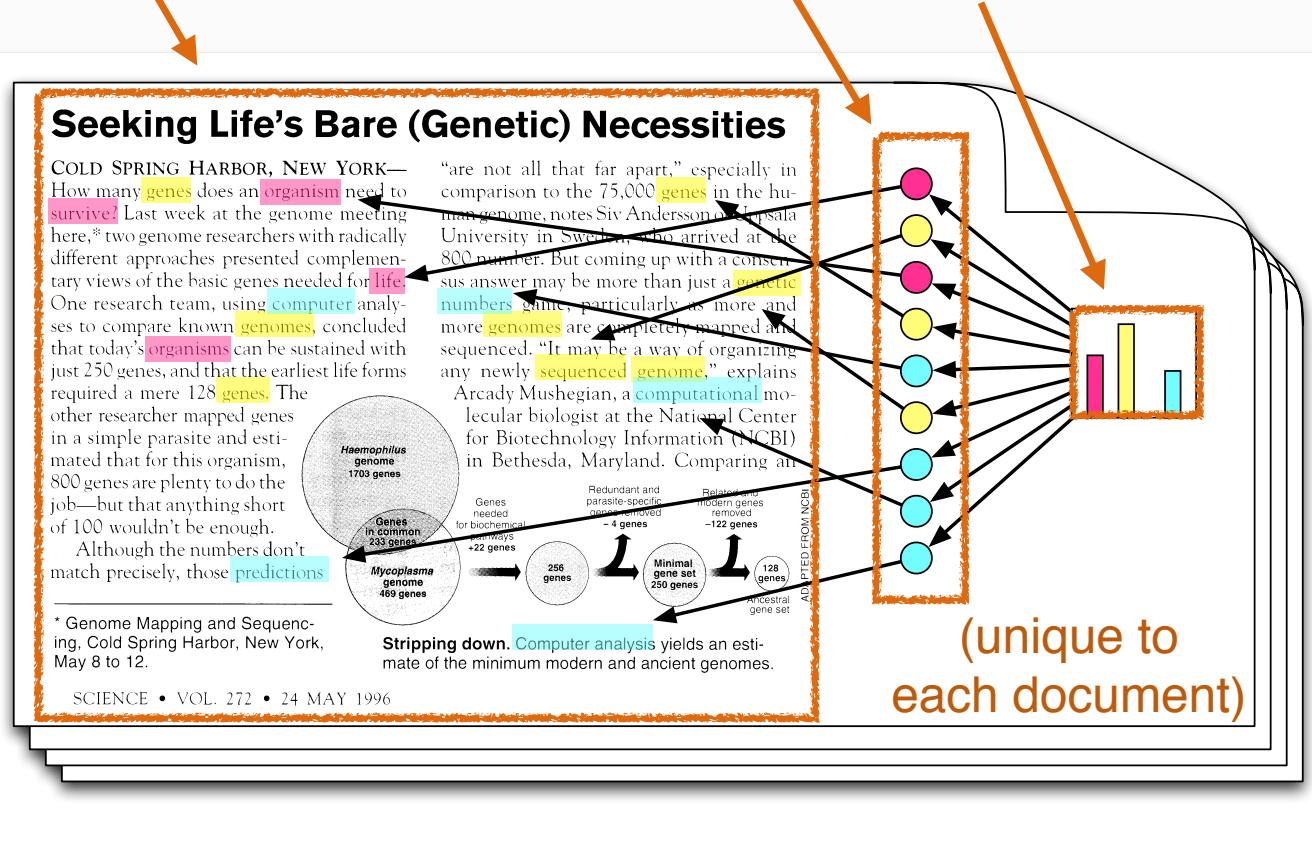
Although the numbers don't match precisely, those **predictions**

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

SCIENCE • VOL. 272 • 24 MAY 1996

$z_d$ : Assignments

$\theta_d$ : Topic Proportions



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

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

# 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 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

"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 view of the bare necessities is not an easy job—but that of 100 would do it.

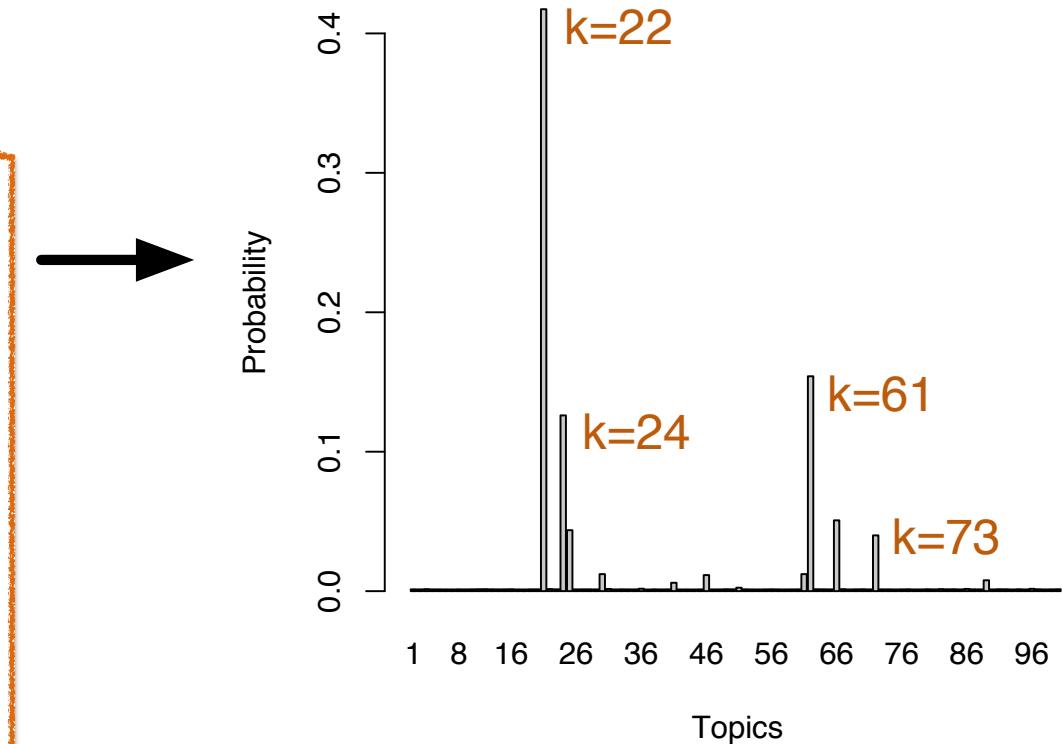
Although the two teams did not match precise

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

SCIENCE

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

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

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 can result in complex dynamic behavior. Natural

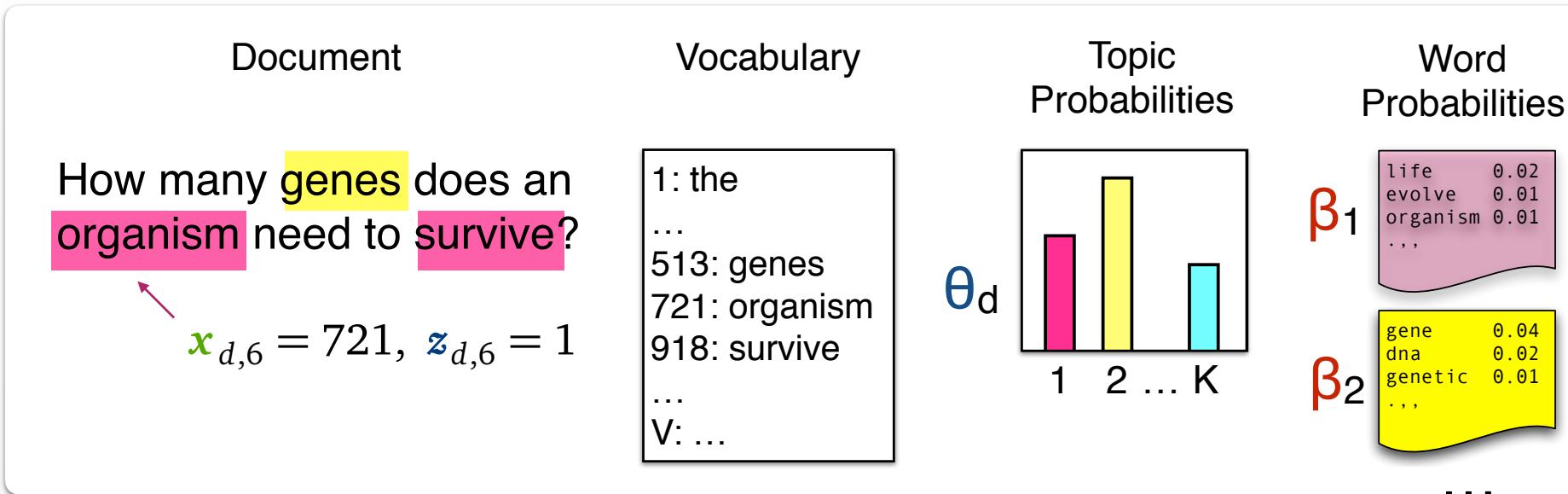
convincing evidence to date of complex dynamics and chaos in a biological population—of the flour beetle, *Tribolium castaneum* (see figure). It has proven extremely difficult to obtain such evidence because the beetles move over the surface of the attractor, sets of adjacent trajectories are pulled 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.



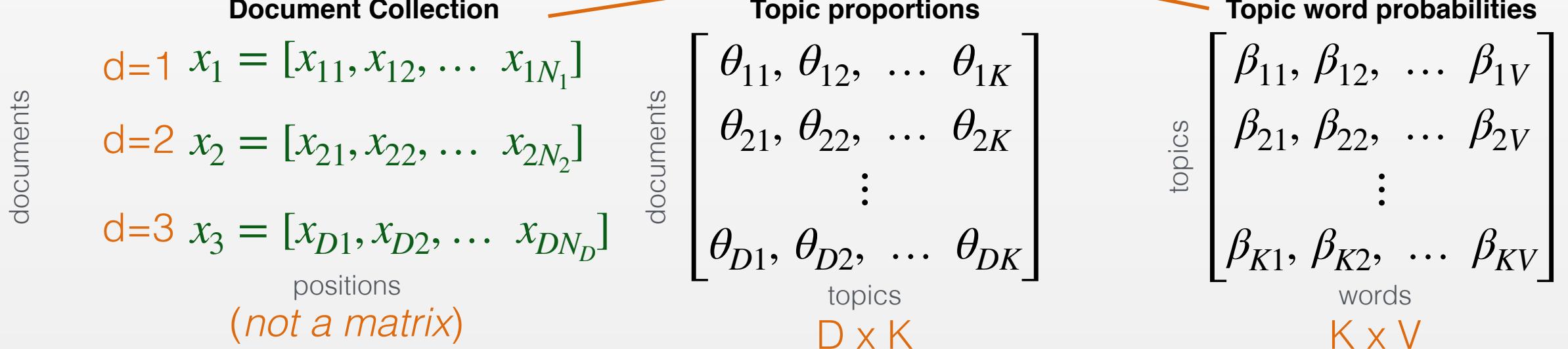
The authors are in the Department of Biology, Imperial College at Silwood Park, Ascot, Berks, SL5 7PZ; mail: m.hassell@ic.ac.uk.

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



# Interpretation as Matrix Factorization

Bag of word Vector

$$X_d = [X_{d1}, X_{d2}, \dots, X_{dV}], \\ X_{dv} = \sum_{n=1}^{N_d} I[x_{dn} = v]$$

$$\mathbb{E}[X_{dv}] = N_d p(x_{dn} = v | \theta, \beta) = N_d \sum_{k=1}^K \theta_{dk} \beta_{kv}$$

$$\mathbb{E}[X] = N \times \theta \times \beta$$

**Document word counts**      **Topic counts**      **Topic word probabilities**  
 $\mathbb{E} \begin{bmatrix} X_{11}, X_{12}, \dots, X_{1V} \\ X_{21}, X_{22}, \dots, X_{2V} \\ \vdots \\ X_{D1}, X_{D2}, \dots, X_{DV} \end{bmatrix} = \begin{bmatrix} N_1, 0, \dots, 0 \\ 0, N_2, \dots, 0 \\ \vdots \\ 0, 0, \dots, N_D \end{bmatrix} \times \begin{bmatrix} \theta_{11}, \theta_{12}, \dots, \theta_{1K} \\ \theta_{21}, \theta_{22}, \dots, \theta_{2K} \\ \vdots \\ \theta_{D1}, \theta_{D2}, \dots, \theta_{DK} \end{bmatrix} \times \begin{bmatrix} \beta_{11}, \beta_{12}, \dots, \beta_{1V} \\ \beta_{21}, \beta_{22}, \dots, \beta_{2V} \\ \vdots \\ \beta_{K1}, \beta_{K2}, \dots, \beta_{KV} \end{bmatrix}$

$$X \approx N \times \theta \times \beta$$

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

# Relationship to Latent Semantic Analysis

LSA: Factorize word counts (using PCA)

$$\mathbf{X} \text{ (V x D)} \approx \mathbf{U} \text{ (V x K)} \mathbf{Z} \text{ (K x D)}$$

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

PCA as matrix factorization:  
•  $Z = U^T X$   
•  $X$  is approximated as  $X \approx UZ$   
Minimize Reconstruction error  
 $\min_{U,Z} \| X - UZ \|^2$

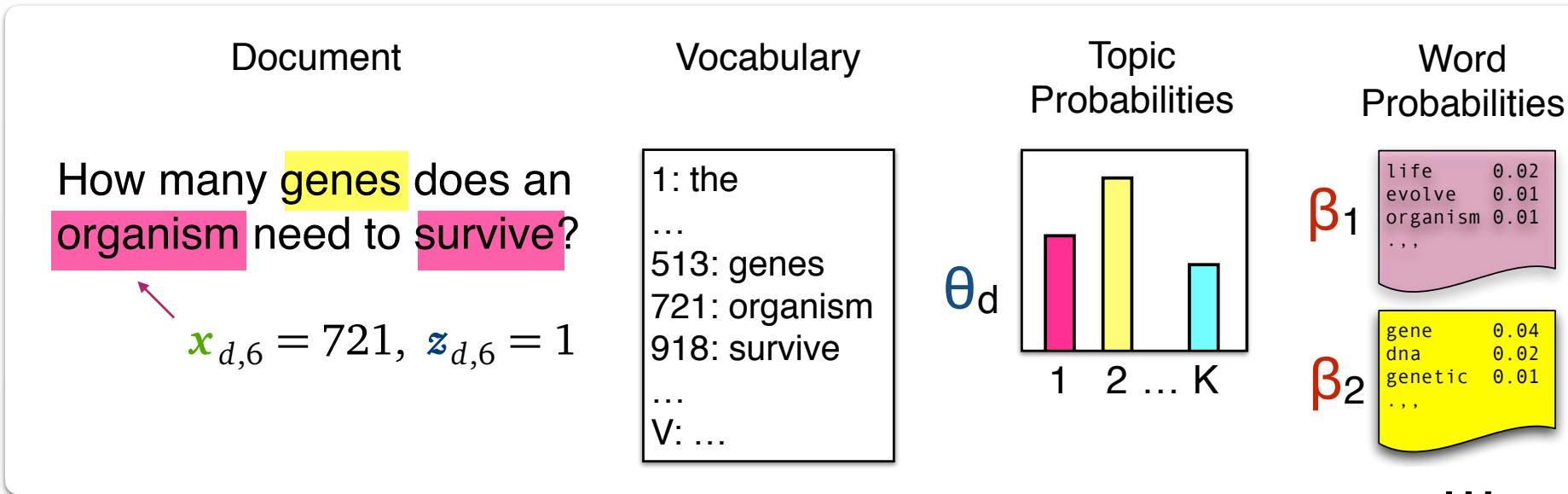
Topic Models: Factorize word counts (using mixture model)

$$\mathbf{X} \approx N \times \boldsymbol{\theta} \times \boldsymbol{\beta} \quad N^{-1} \mathbf{X}_{(D \times V)} \approx \boldsymbol{\theta}_{(D \times K)} \quad \boldsymbol{\beta}_{(K \times V)}$$

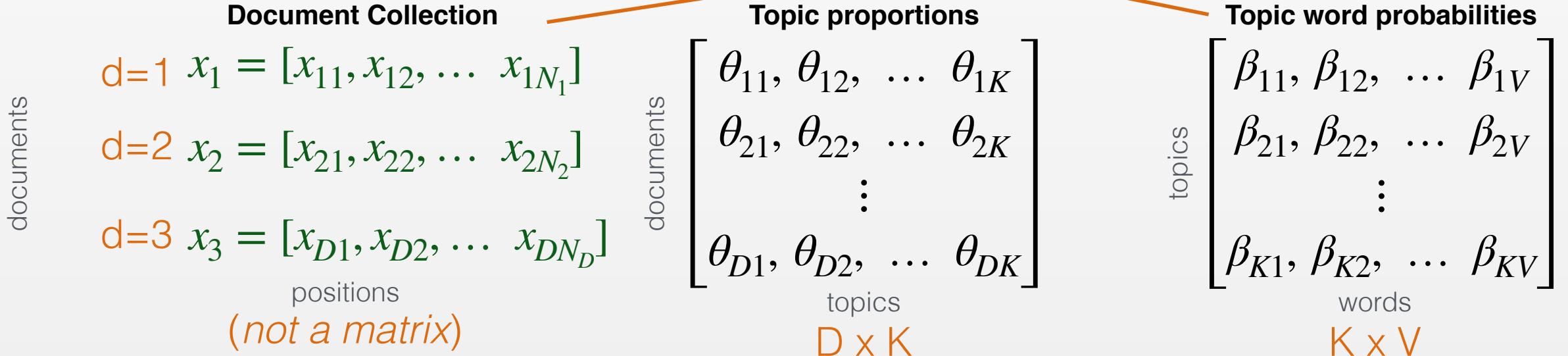
$$\mathbf{X}^T N^{-1}_{(V \times D)} \approx \boldsymbol{\beta}^T_{(V \times K)} \quad \boldsymbol{\theta}^T_{(K \times D)}$$

Topic models can also be tackled as matrix factorization problem, where rows of  $\boldsymbol{\beta}, \boldsymbol{\theta}$  are constrained to be probability vectors

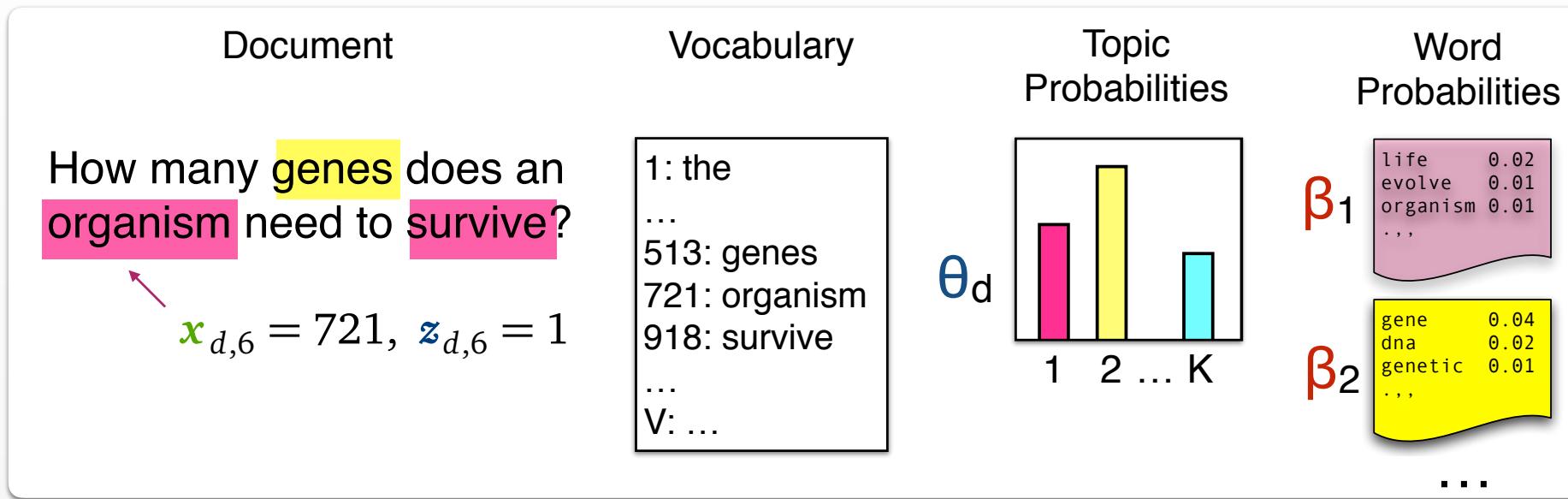
# Estimating the Parameters



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



# Calculating the Likelihood for each Word



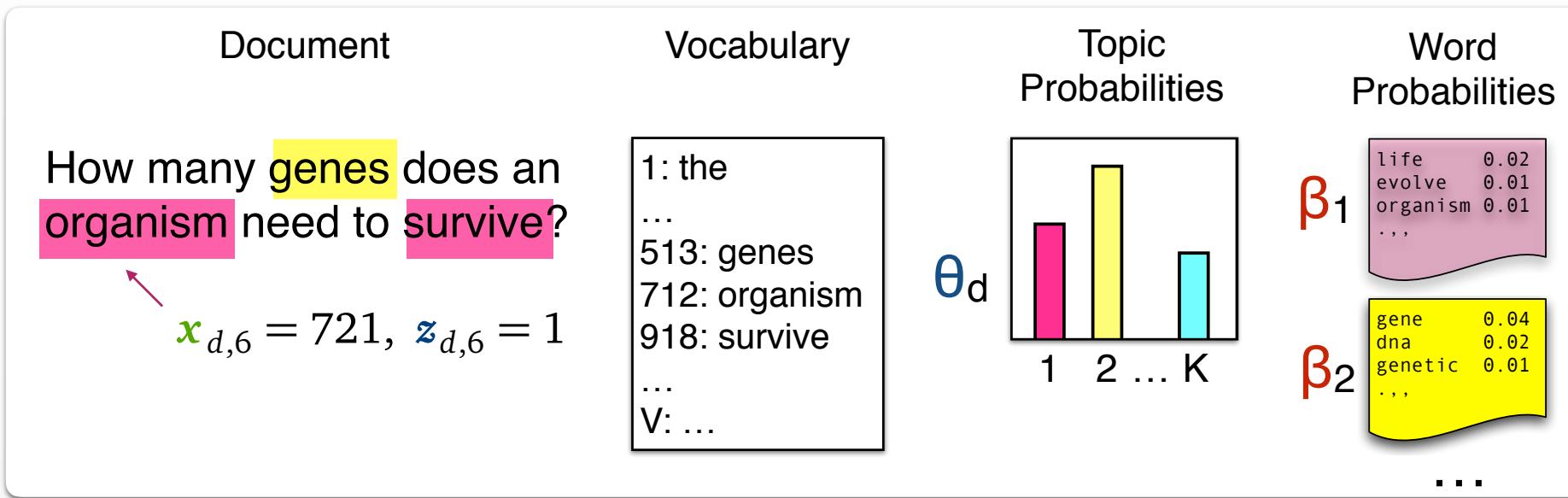
Probability that word  $n$  is entry  $v$  in the vocabulary

Probability of word  $v$  given topic  $k$

Probability that word belongs to topic  $k$

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

# Computing the Likelihood



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

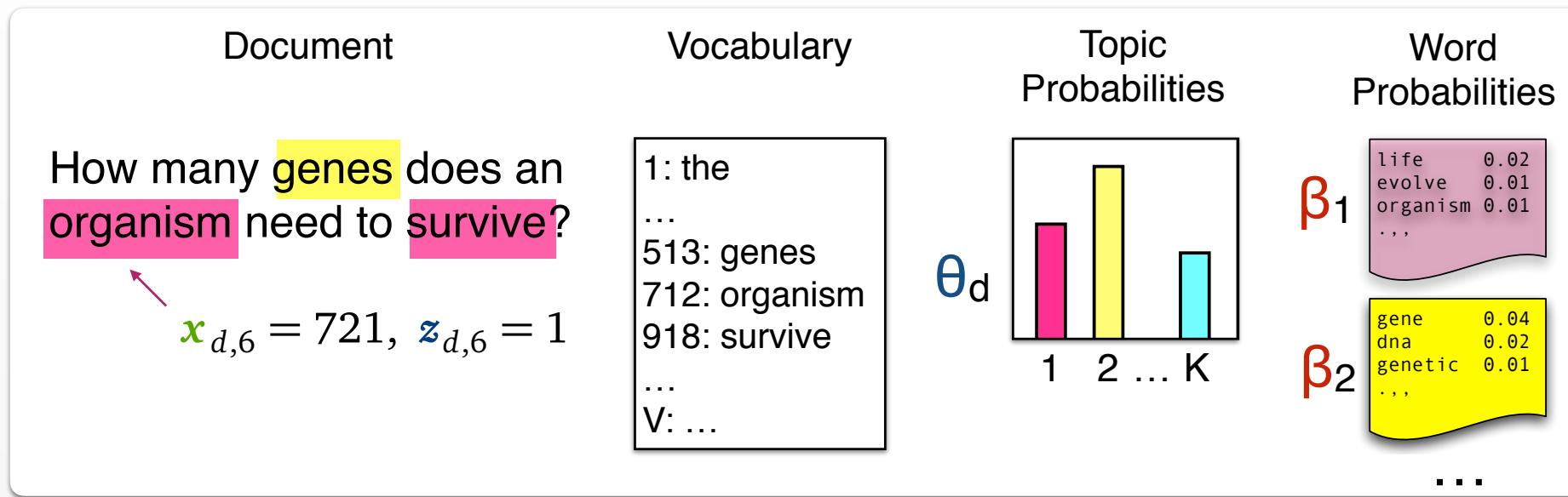
$$x_d = [x_{d1}, x_{d2}, \dots, x_{dN_d}]$$

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

take log probability, substitute result from previous slide

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

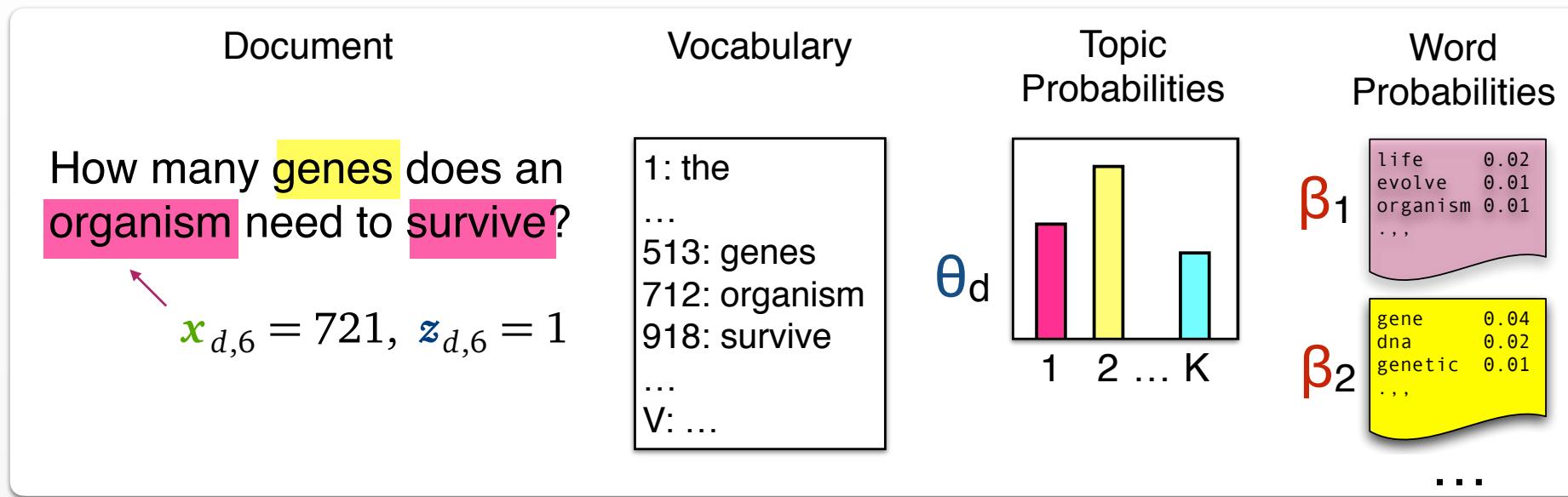
# Calculating the Likelihood for all Words



log probability of all words in document  $d$

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

# Calculating the Likelihood for all Words



log probability of all words in document  $d$

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

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

bag-of-word vector

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

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

# Interpretation as Matrix Factorization

Log likelihood

$$\log p(x_d | \beta, \theta_d) = X_d \log(\theta_d \beta)^T$$

$$\log p(x | \beta, \theta) = \sum_{d=1}^D X_d \log(\theta_d \beta)^T$$

Bag of word Vector

$$X_d = [X_{d1}, X_{d2}, \dots, X_{dV}],$$
$$X_{dv} = \sum_{n=1}^{N_d} I[x_{dn} = v]$$

Difficult to optimize with gradient ascent.  
Need to use EM!

# Sketch of EM derivation

## Topic models

### Observed variable

- $x = \{x_{nd}\}$ : all words in the corpus.

### Unobserved variable

- $z = \{z_{nd}\}$  : the topics of all words in the corpus

### Model Parameters

- $\beta = \{\beta_k\}$ : word probabilities, one per topic
- $\theta = \{\theta_d\}$ : topic probabilities, one per document

$$\eta = \{\theta, \beta\}$$

$$l(\eta) = \log p(x | \eta)$$

- 1) Initialize  $\eta$  as  $\eta_0$
- 2) For  $t = 0, 1, \dots$ , repeat until convergence

2a) E-Step:

$$Q^t(\eta) = \mathbf{E}_{q^t(z)}[\log p(x, z | \eta)],$$

where  $q^t(z) = p(z | x, \eta^t)$

2b) M-Step:

$$\eta^{t+1} \leftarrow \operatorname{argmax}_\eta Q^t(\eta)$$

## Clustering: GMM

### Observed variable

- $x = \{x_n\}$ : all datapoint

### Unobserved variable

- $z = \{z_n\}$ : the cluster index of all datapoint

### Model Parameters

- $\mu = \{\mu_k\}$ : the means of all gaussian clusters.
- $\Sigma = \{\Sigma_k\}$ : the covariances of all gaussian clusters
- $\alpha = \{\alpha_k\}$ : the cluster proportions

$$\theta = \{\alpha, \mu, \Sigma\}$$

- 1) Initialize  $\theta$  as  $\theta_0$

- 2) For  $t = 0, 1, \dots$ , repeat until convergence

$$l(\theta) = \log p(x | \theta)$$

2a) E-Step:

$$Q^t(\theta) = \mathbf{E}_{q^t(z)}[\log p(x, z | \theta)],$$

where  $q^t(z) = p(z | x, \theta^t)$

2b) M-Step:

$$\theta^{t+1} \leftarrow \operatorname{argmax}_\theta Q^t(\theta)$$

# PLSI/PLSA\*: EM for Topic Models

## Generative Model

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

Need to compute  
 $q^t(z) = p(z|x, \theta^t, \beta^t)$  for the E-Step.  
This is a joint distribution over all the topic assignment variables, one per word in the corpus. It can be decomposed as a product of  $\phi_{d,n,k}$ .

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

# PLSI/PLSA: E-step

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

General Form, with One-hot Indexing Trick

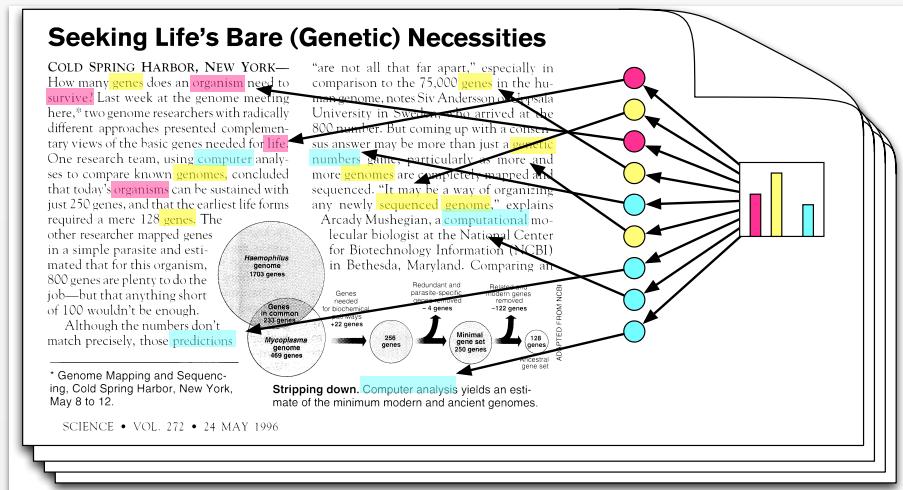
$$\phi_{d,n,k} = \frac{\theta_{d,k} \left( \sum_{v=1}^V \beta_{k,v} I[\mathbf{x}_{d,n} = v] \right)}{\sum_{l=1}^K \theta_{d,l} \left( \sum_{v=1}^V \beta_{l,v} I[\mathbf{x}_{d,n} = v] \right)}$$

# 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

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

## M-step: Update parameters

Use assignment probabilities  $\boldsymbol{\phi}_d$  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: 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 as

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

Fraction of topic  $k$  in document  $d$

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

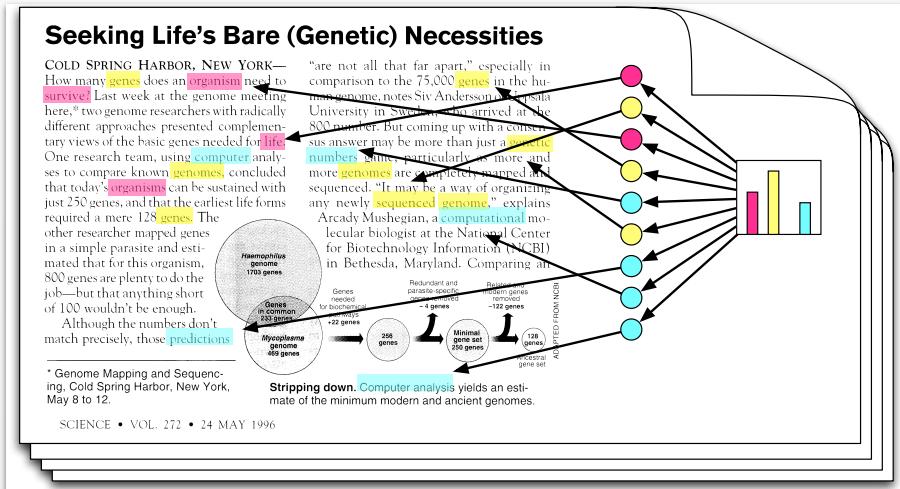
Fraction of word  $v$  in topic  $k$

# 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

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

## M-step: Update parameters

$$\begin{aligned} \boldsymbol{\beta}_{k,v} &= \frac{N_{k,v}^\beta}{\sum_{d=1}^D N_{d,k}^\theta} & N_{k,v}^\beta &= \sum_{d=1}^D \sum_{n=1}^{N_d} \boldsymbol{\phi}_{d,n,k} I[\mathbf{x}_{d,n} = v] \\ \boldsymbol{\theta}_{d,k} &= \frac{N_{d,k}^\theta}{N_d} & N_{d,k}^\theta &= \sum_{n=1}^{N_d} \boldsymbol{\phi}_{d,n,k} \end{aligned}$$

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

# Topic Models: *Summary so far*

Core Idea:

Model documents as *mixtures* over topics

Model Parameters:

$\theta_d$  Topic probabilities for each document  
(K-dimensional vector)

$\beta_k$  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



# Latent Dirichlet Allocation

## Topic Models with Dirichlet Priors

# Review: Topic Modeling with PLSA/PLSI

$\beta_k$ : Topics  
(shared)

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
...	

$x_d$ : Words

$z_d$ : Assignments  
(document-specific)

$\theta_d$ : Topic Proportions

## Seeking Life's Bare (Genetic) Necessities

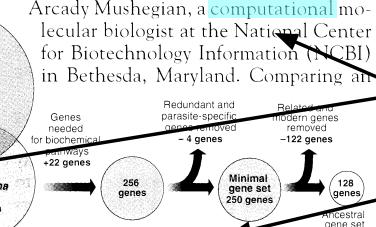
COLD SPRING HARBOR, NEW YORK—How many genes does an organism need to survive? Last week at the genome meeting here,\* two genome researchers with radically different approaches presented complementary views of the basic genes needed for life. One research team, using computer analyses to compare known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



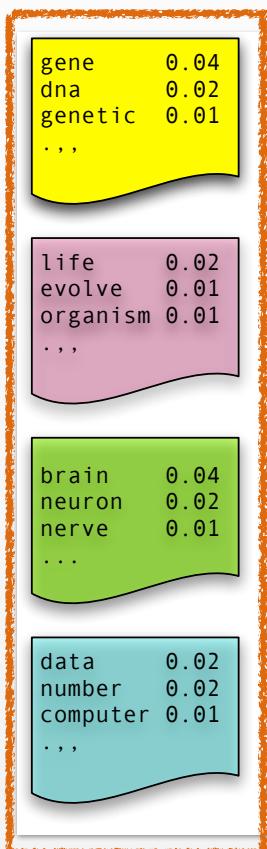
Stripping down. Computer analysis yields an estimate of the minimum modern and ancient genomes.

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

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

# LDA: Add Dirichlet Priors

$\beta_k$ : Topics  
(shared)



$x_d$ : Words

$z_d$ : Assignments  
(document-specific)

$\theta_d$ : Topic Proportions

## Seeking Life's Bare (Genetic) Necessities

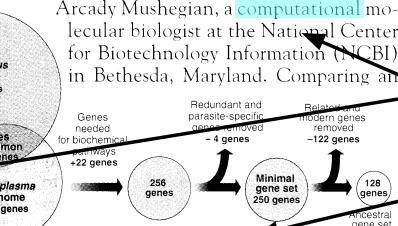
COLD SPRING HARBOR, NEW YORK—How many genes does an organism need to survive? Last week at the genome meeting here,\* two genome researchers with radically different approaches presented complementary views of the basic genes needed for life. One research team, using computer analyses to compare known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

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

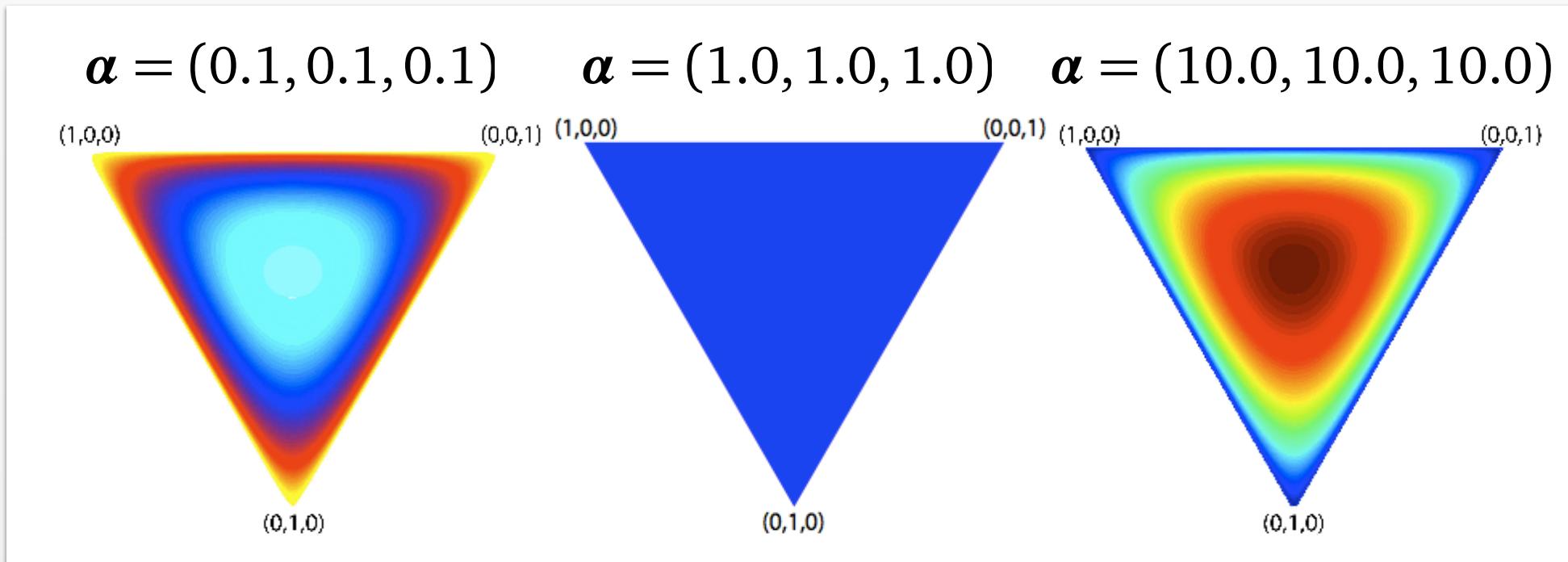
$$x_{d,n} | z_{d,n}=k \sim \text{Discrete}(\beta_k)$$

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

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

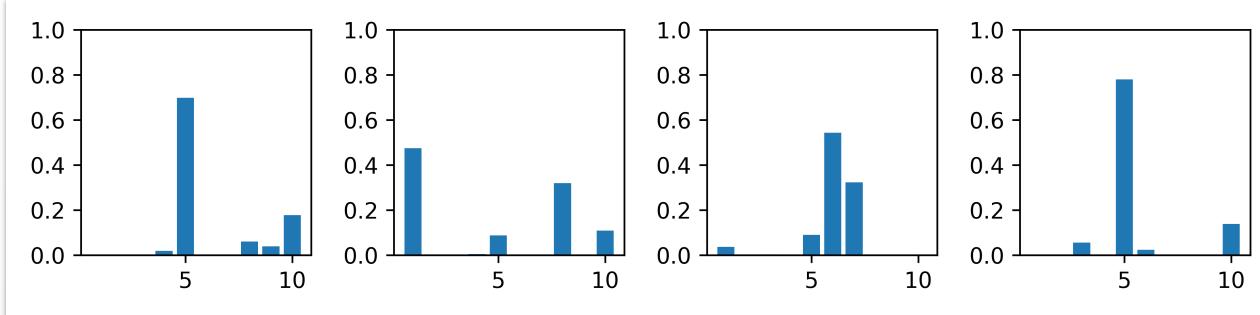
# Review: Dirichlet Distribution

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

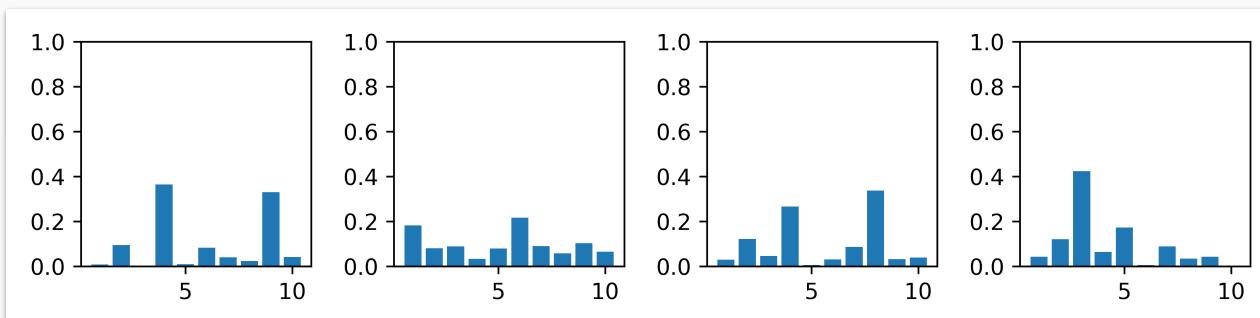


# Review: Dirichlet Distribution

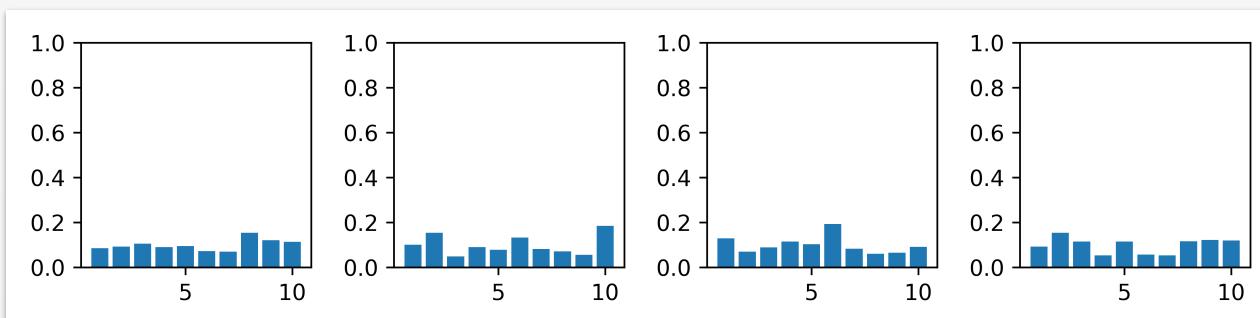
$\alpha_k = 0.1$



$\alpha_k = 1.0$



$\alpha_k = 10.0$



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

# 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  $\theta$  (unknown) defined on the sample space the space  $\mathcal{X}$

Maximum likelihood estimate of  $\theta$ .

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

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.

# 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  $\theta$

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  $\theta$

$$\theta \sim P_\alpha$$

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

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

# Maximum a posteriori (MAP) estimation

Likelihood

$$\begin{aligned} p(X|\mu) &= \prod_{n=1}^N p(x_n|\mu) \\ &= \prod_{n=1}^N \mu^{x_n} (1-\mu)^{1-x_n} \end{aligned}$$

$$p(\mu|X, a, b) \propto p(X|\mu)p(\mu|a, b)$$

Posterior

$$\begin{aligned} &\propto \left( \prod_{n=1}^N \mu^{x_n} (1-\mu)^{1-x_n} \right) \frac{\mu^{a-1} (1-\mu)^{b-1}}{B(a, b)} \quad \text{Not a function of } \mu \\ &\propto \mu \left[ \sum_{n=1}^N x_n \right] (1-\mu)^{N - \sum_{n=1}^N x_n} \mu^{a-1} (1-\mu)^{b-1} \\ &\propto \mu^{N_1 + a - 1} (1-\mu)^{N_0 + b - 1} \quad N_1 = \sum_{n=1}^N x_n \\ &= \text{Beta}(\mu | N_1 + a, N_0 + b) \quad N_0 = N - N_1 \end{aligned}$$

MAP Estimate

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

Prior

$$p(\mu|a, b) = \frac{\mu^{a-1} (1-\mu)^{b-1}}{B(a, b)}$$

$$X = \{x_n\}_{n=1}^N \quad x_n \in \{0, 1\}$$

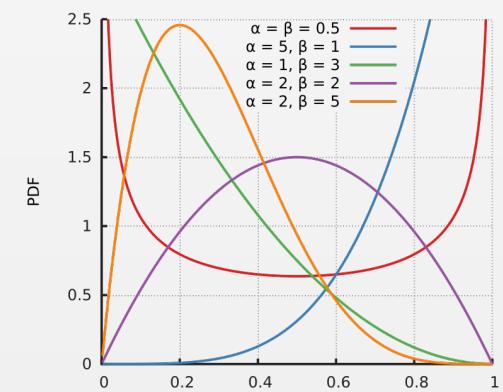
$$x_n \sim \text{Bernoulli}(\mu)$$

$$p(x_n|\mu) = \begin{cases} \mu & \text{if } x_n = 1 \\ 1-\mu & \text{if } x_n = 0 \end{cases}$$

$$p(x_n|\mu) = \mu^{x_n} (1-\mu)^{1-x_n}$$

$$\mu \sim \text{Beta}(a, b)$$

$$p(\mu|a, b) = \frac{\mu^{a-1} (1-\mu)^{b-1}}{B(a, b)}$$



# MAP generalization to Discrete distribution

Likelihood

$$p(X|\mu) = \prod_{n=1}^N p(x_n|\mu)$$

$$= \prod_{n=1}^N \prod_{k=1}^K \mu_k^{I[x_n=k]}$$

$$p(\mu|X, a, b) \propto p(X|\mu)p(\mu|a, b)$$

Posterior

$$\propto \left( \prod_{n=1}^N \prod_{k=1}^K \mu_k^{I[x_n=k]} \right) \frac{\prod_{k=1}^K \mu_k^{\alpha_k}}{B(\alpha)} \quad \text{Not a function of } \mu$$

$$\propto \prod_{k=1}^K \mu_k^{\left[ \sum_{n=1}^N I[x_n=k] \right]} \prod_{k=1}^K \mu_k^{\alpha_k}$$

$$\propto \prod_{k=1}^K \mu_k^{N_k + \alpha_k} \quad N_k = \sum_{n=1}^N I[x_n = k]$$

$$= \text{Dirichlet}(\mu | N_1 + \alpha_1, N_2 + \alpha_2, \dots, N_K + \alpha_K)$$

MAP Estimate

$$\mu_k^{MAP} = \frac{N_k + \alpha_k - 1}{N + \sum_{k=1}^K \alpha_k - K}$$

Prior

$$p(\mu|\alpha) = \frac{\prod_{k=1}^K \mu_k^{\alpha_k}}{B(\alpha)}$$

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

$$x_n \sim \text{Discrete}(\mu) \quad \sum_{k=1}^K \mu_k = 1$$

$$p(x_n|\mu) = \begin{cases} \mu_1 & \text{if } x_n = 1 \\ \mu_2 & \text{if } x_n = 2 \\ \vdots & \\ \mu_K & \text{if } x_n = K \end{cases}$$

$$p(x_n|\mu) = \prod_{k=1}^K \mu_k^{I[x_n=k]}$$

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

$$p(\mu|\alpha) = \frac{\prod_{k=1}^K \mu_k^{\alpha_k}}{B(\alpha)}$$

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

# LDA: *Summary so far*

- Idea: Model documents as *mixtures* over topics
- Model parameters:
  - $\Theta_d$  Topic probabilities for each document  
(K-dimensional vector for each document)
  - $\beta_k$  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