



Topic Models

Shantanu Jain



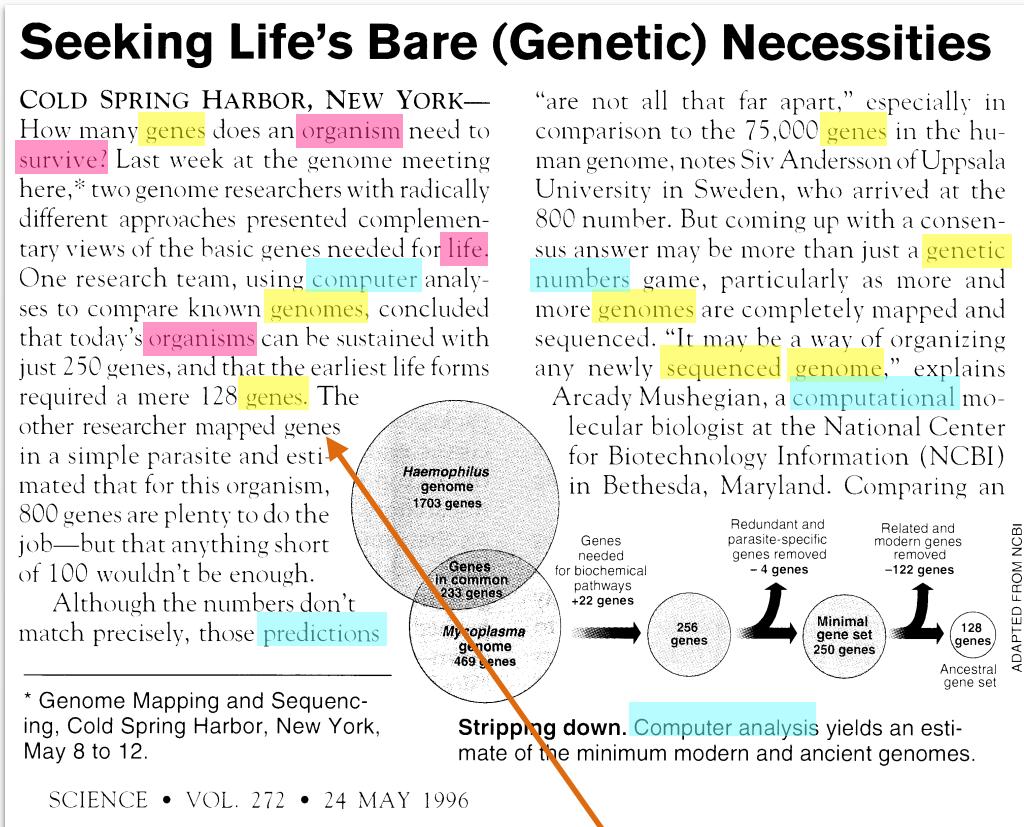
Topic Modeling Basics



Borrowing from:
David Blei
(Columbia)

Word Mixtures

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



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

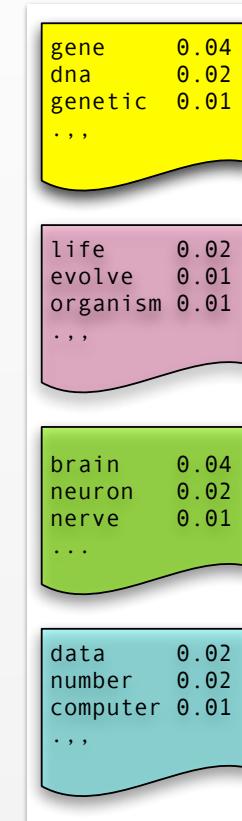
Topic assignment: $z_n \in \{1, \dots, K\}$

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

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

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

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



- 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.

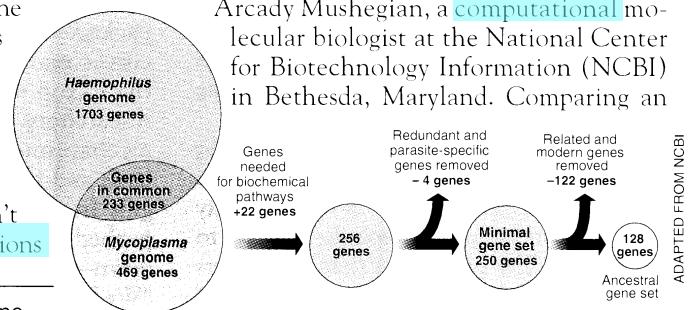
Word Mixtures

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

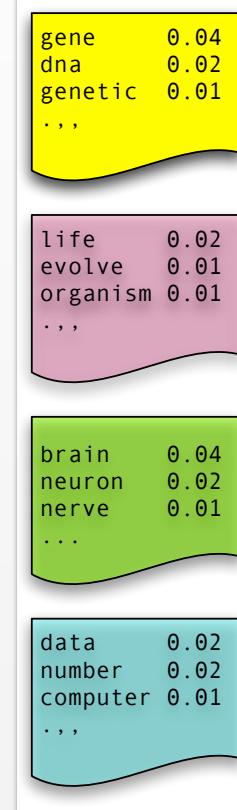
"are not all that far apart," especially in comparison to the 75,000 genes in the human genome, notes Siv Andersson of Uppsala University in Sweden, who arrived at the 800 number. But coming up with a consensus answer may be more than just a genetic numbers game, particularly as more and more genomes 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.

$$\begin{aligned} z_n &\sim \text{Discrete}(\theta) \\ x_n | z_n=k &\sim \text{Discrete}(\beta_k) \end{aligned}$$

Pick a topic
Pick a word given topic



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

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

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

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

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

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

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

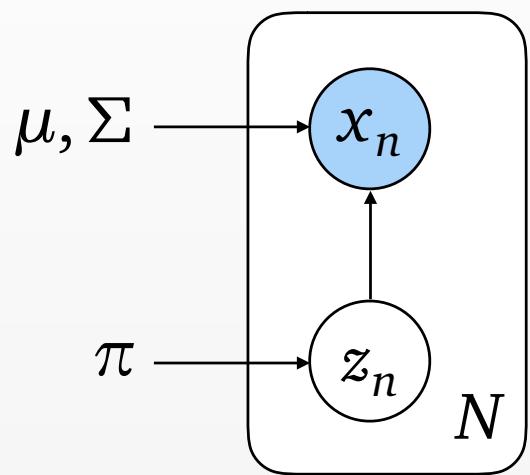
β_k : k^{th} topic's word probabilities over the vocabulary

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

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

Gaussian Mixtures vs Word Mixtures

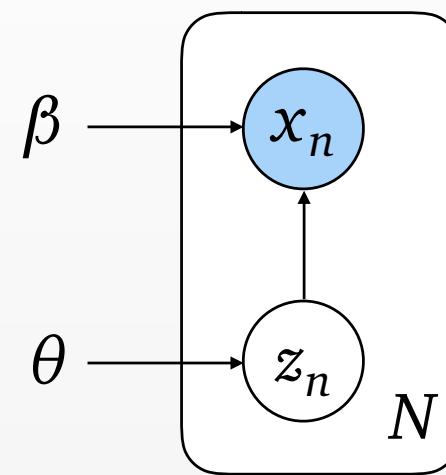
Gaussian Mixture Model



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

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

Discrete Mixture Model



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

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

Difference: Replace Gaussian with Discrete

Using the term Discrete distribution to mean Categorical distribution

Topic Modeling

Topics
(shared)

gene 0.04
dna 0.02
genetic 0.01
...

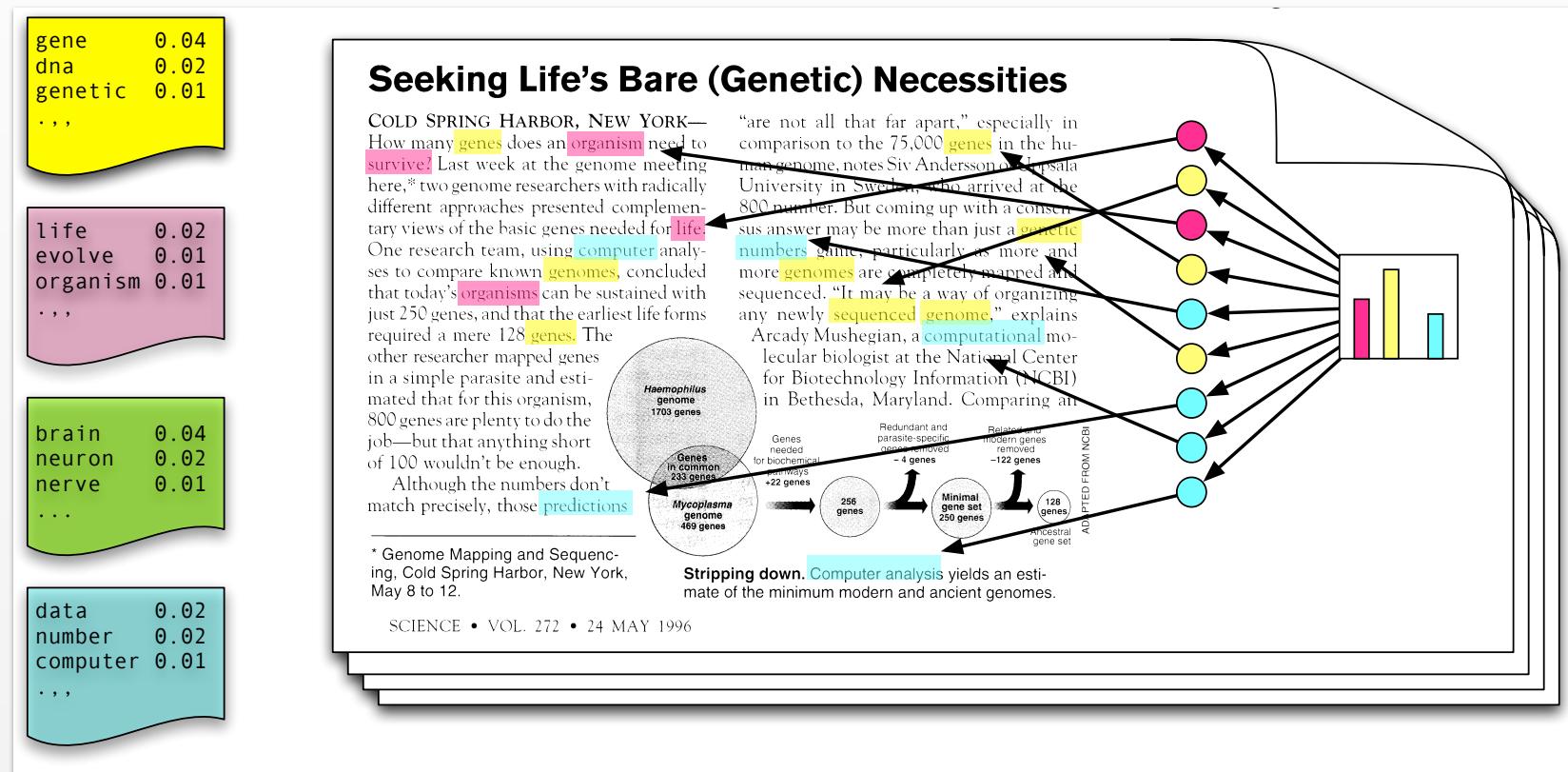
life 0.02
evolve 0.01
organism 0.01
...

brain 0.04
neuron 0.02
nerve 0.01
...

data 0.02
number 0.02
computer 0.01
...

Words in Document
(mixture over topics)

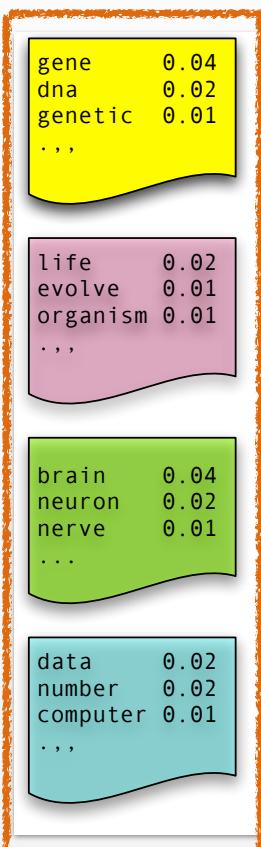
Topic Proportions
(document-specific)



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

Topic Modeling

β_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.

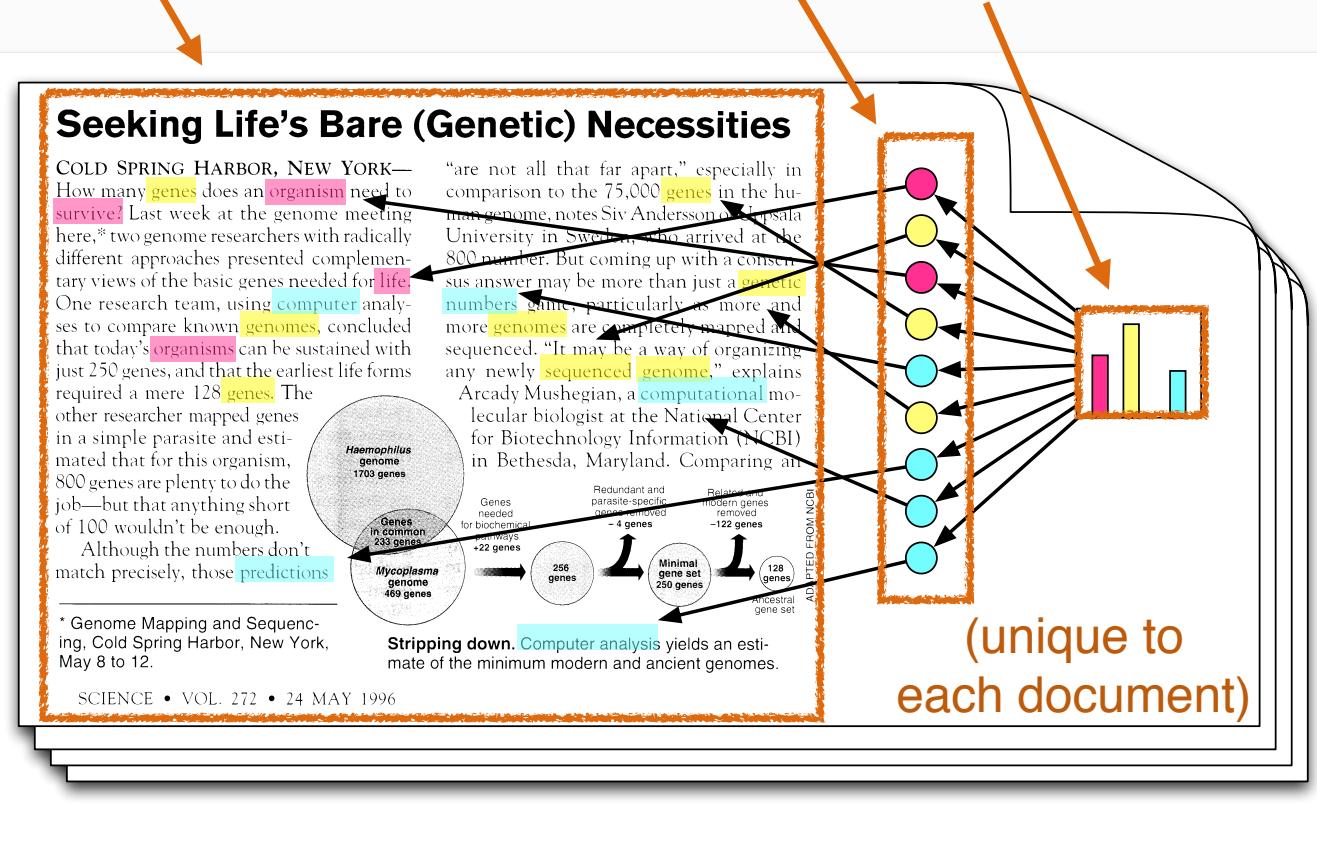
Although the numbers don't match precisely, those **predictions**

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

SCIENCE • VOL. 272 • 24 MAY 1996

z_d : Assignments

θ_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.

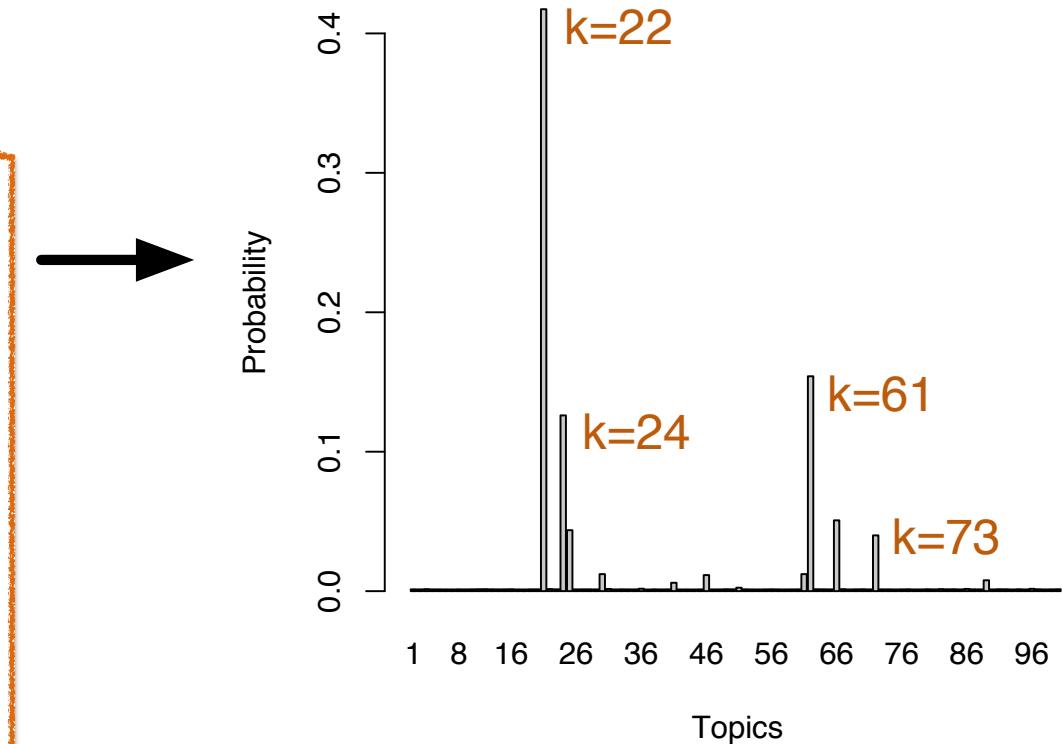
Although the two teams did not match precise

* 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 are small and their life cycle is long. The

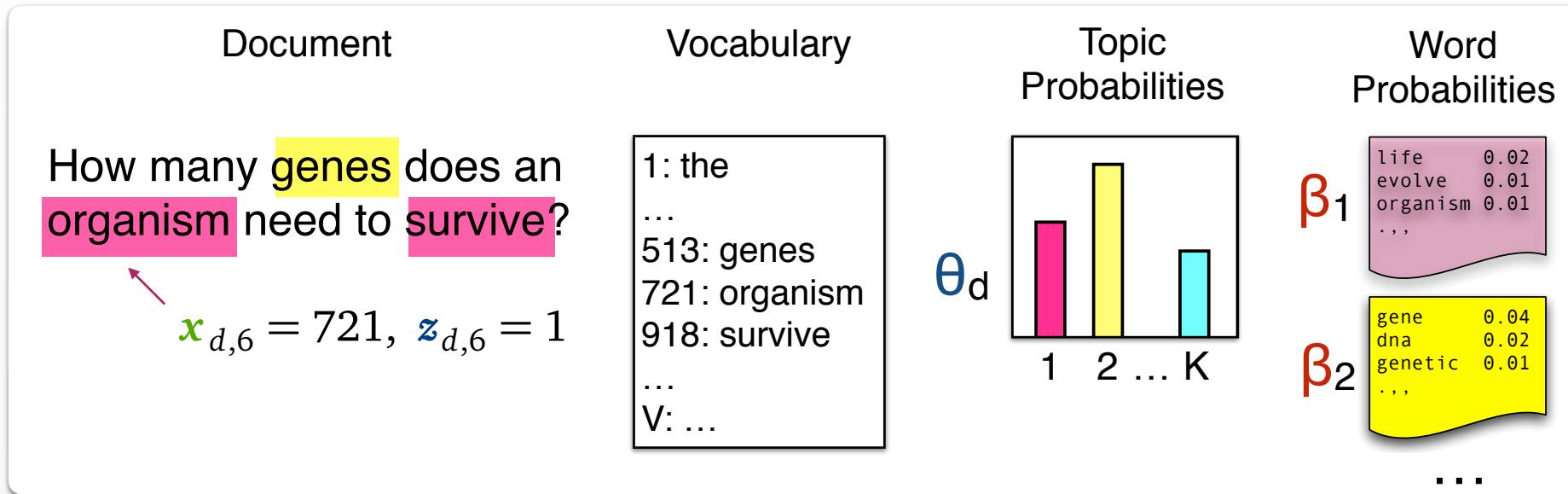
move over the surface of the attractor, sets of adjacent trajectories are pulled apart, then stretched and folded, so that it becomes impossible to predict exact population densities into the future. The strength of the mixing that gives rise to the extreme sensitivity to initial conditions can be measured 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



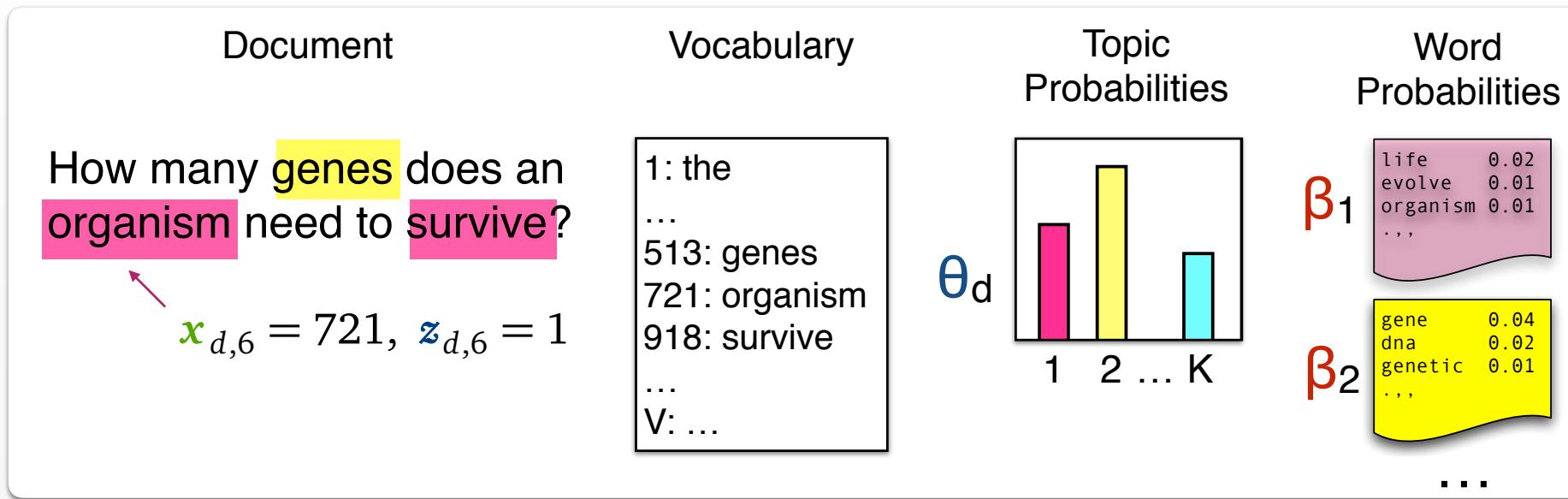
Maximum Likelihood: $\max_{\theta, \beta} \log p(\mathbf{x} | \theta, \beta)$

$d=1$	$[x_{1,1}, x_{1,2}, \dots, x_{1,N_1}]$	$[\theta_{1,1}, \theta_{1,2}, \dots, \theta_{1,K}]$	$[\beta_{1,1}, \beta_{1,2}, \dots, \beta_{1,V}]$	$k=1$
$d=2$	$[x_{2,1}, x_{2,2}, \dots, x_{2,N_2}]$	$[\theta_{2,1}, \theta_{2,2}, \dots, \theta_{2,K}]$	$[\beta_{2,1}, \beta_{2,2}, \dots, \beta_{2,V}]$	$k=2$

$d=D$	$[x_{D,1}, x_{D,2}, \dots, x_{D,N_D}]$ <i>(not a matrix)</i>	$[\theta_{D,1}, \theta_{D,2}, \dots, \theta_{D,K}]$	$[\beta_{K,1}, \beta_{K,2}, \dots, \beta_{K,V}]$	$k=K$

$D \times K$ $K \times V$

Calculating the Likelihood for each Word



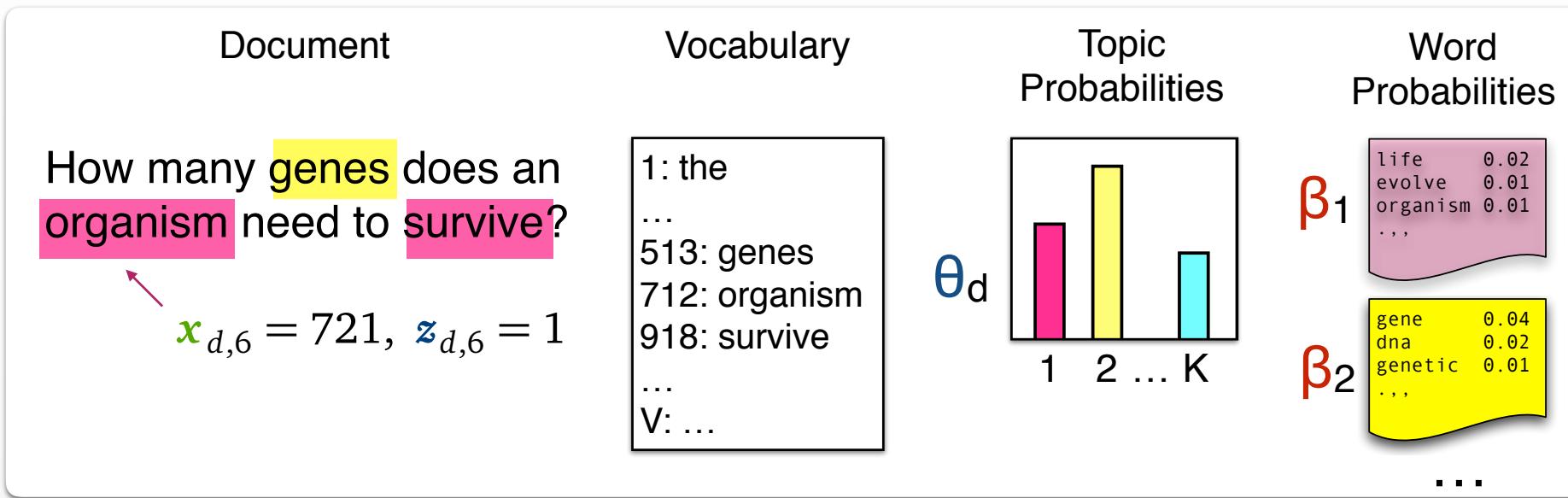
Probability that word n is entry v in the vocabulary

Probability of word v given topic k

Probability that word belongs to topic k

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

Computing the Likelihood



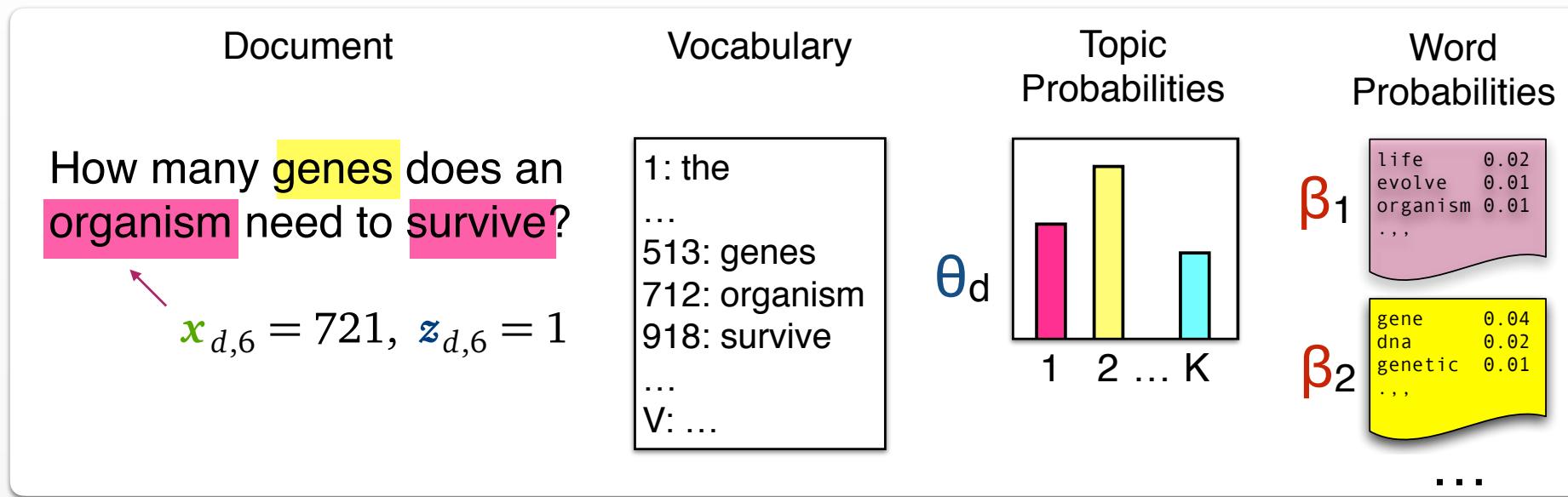
probability of all words $n = 1 \dots N_d$ in document d (use one-hot trick)

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

take log probability, substitute result from previous slide

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

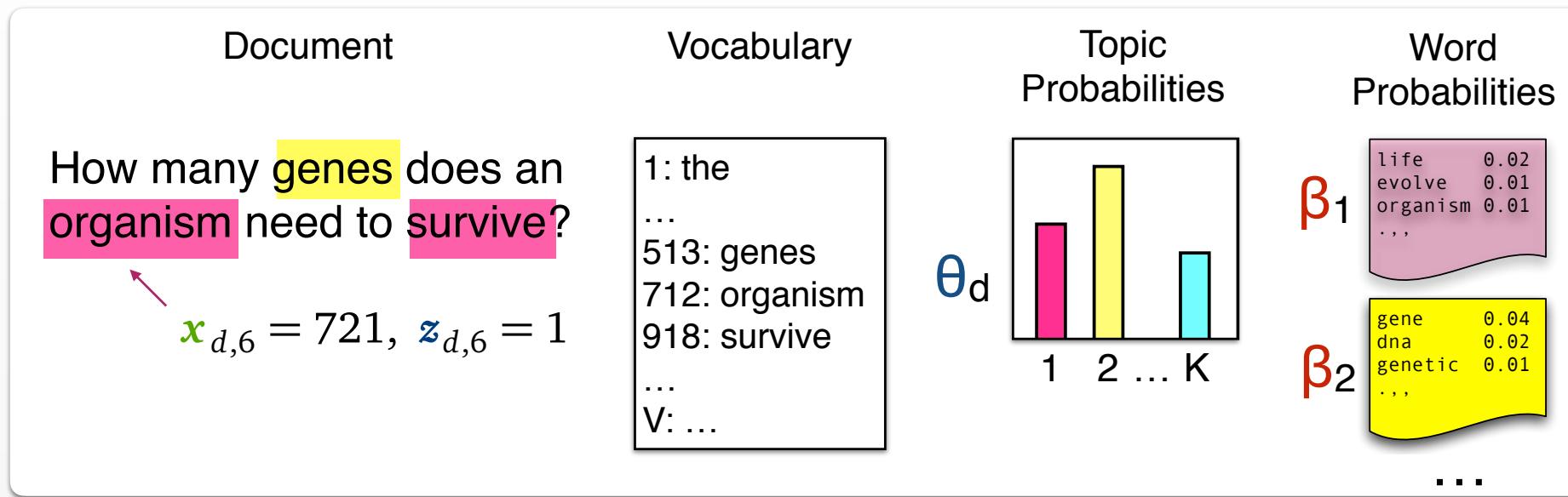
Calculating the Likelihood for all Words



log probability of all words in document d

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

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

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

Bag of Word Vector

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

Word Counts

$$\mathbb{E} \left[\begin{bmatrix} \mathbf{X}_{1,1} & \dots & \mathbf{X}_{1,V} \\ \dots & & \dots \\ \mathbf{X}_{D,1} & \dots & \mathbf{X}_{D,V} \end{bmatrix}_{(D \times V)} \right] = \begin{bmatrix} N_1 \boldsymbol{\theta}_{1,1} & \dots & N_1 \boldsymbol{\theta}_{1,K} \\ \dots & & \dots \\ N_D \boldsymbol{\theta}_{D,1} & \dots & N_D \boldsymbol{\theta}_{D,K} \end{bmatrix}_{(D \times K)} \begin{bmatrix} \boldsymbol{\beta}_{1,1} & \dots & \boldsymbol{\beta}_{1,V} \\ \dots & & \dots \\ \boldsymbol{\beta}_{K,1} & \dots & \boldsymbol{\beta}_{K,V} \end{bmatrix}_{(K \times V)}$$

Topic Counts

Topic Word Probabilities

stocks chairman the wins game	$\begin{bmatrix} 2 & 4 & 8 & \dots & 0 & 1 \\ \dots & & & & & \\ 0 & 1 & 7 & \dots & 2 & 3 \end{bmatrix}$	\approx	finance Sports stocks game
			$\begin{bmatrix} 112 \cdot 0.91 & \dots & 112 \cdot 0.01 \\ \dots & & \dots \\ 234 \cdot 0.02 & \dots & 234 \cdot 0.86 \end{bmatrix}$
			$\begin{bmatrix} 0.0081 & \dots & 0.0002 \\ \dots & & \dots \\ 0.0001 & \dots & 0.0072 \end{bmatrix}$

Relationship to Latent Semantic Analysis

LSA: Factorize word counts (using PCA)

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

Topic Models: Factorize word counts (using mixture model)

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

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

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