

Latent Dirichlet Allocation

Authors: David M. Blei, Andrew Y. Ng, and Michael I. Jordan
Presenter: Leah Li





measure sustain **peak gust** **sustain wind** **wind gust** **gust mph** **flood stage** **wind mph** **high peak** **peak wind** **mph highwaydat sensor**

The figure is a word cloud where the size of each word represents its frequency or importance in the context of emergency management reports. The most frequent words are 'emergency management', 'rain fall', 'inch rain', 'heavy rain', 'management report', 'reduce visibility', 'wind speed', 'condition continue', 'trees down', and 'county produce'. Other significant terms include 'quarter mile', 'two inch', 'four inch', 'half inch is', 'mile less', 'fall thunderstorms', 'severe drought', 'tough situation', 'water level', 'down areas', 'street flood', 'several peaks', 'drought condition', 'thundershower move', and 'report trees'.

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
...	

Documents

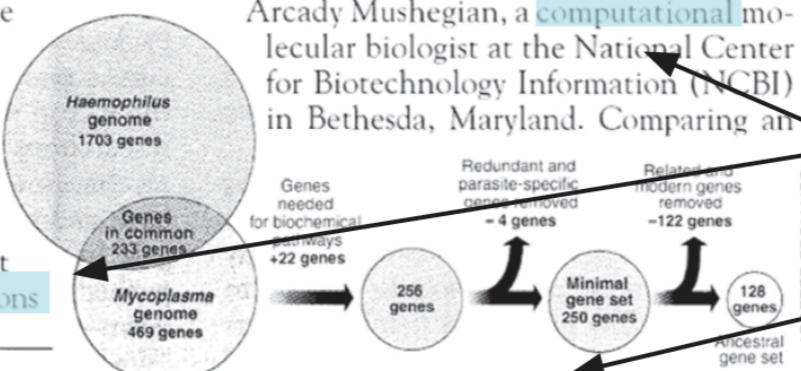
Seeking Life's Bare (Genetic) Necessities

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.

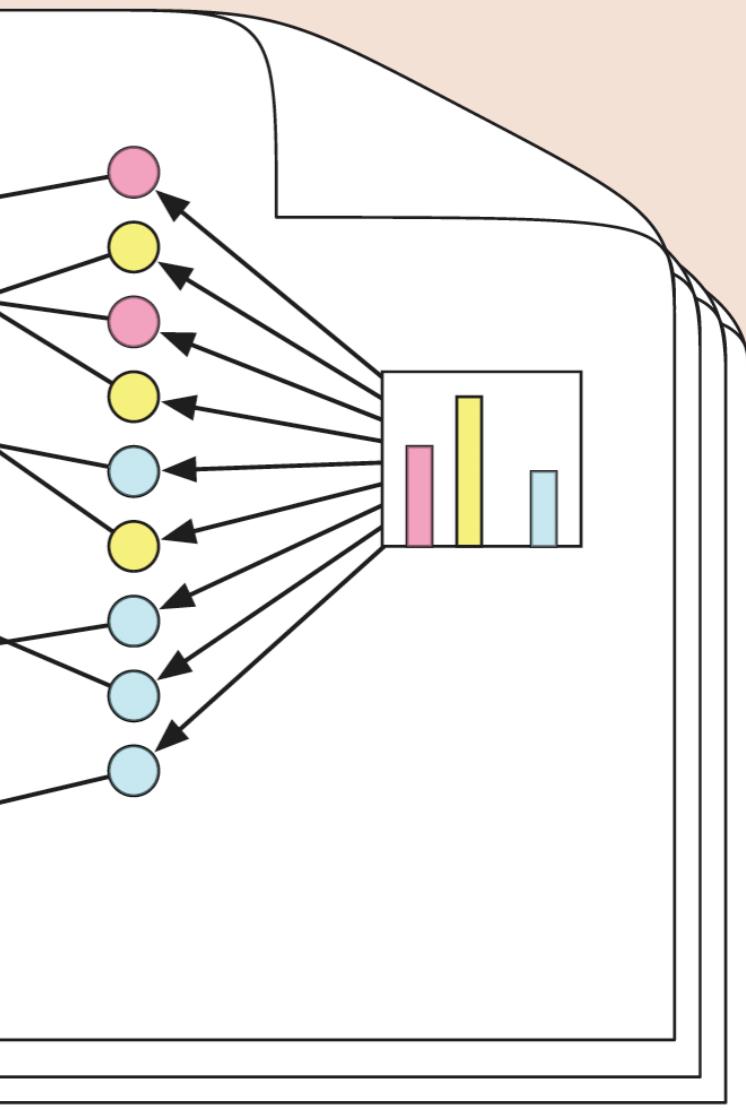
"are not all that far apart," especially in comparison to the 75,000 genes in the human genome, notes Siv Andersson of Uppsala University in Sweden, who arrived at the 800 number. But coming up with a consensus answer may be more than just a genetic numbers game, particularly as more and more genomes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



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

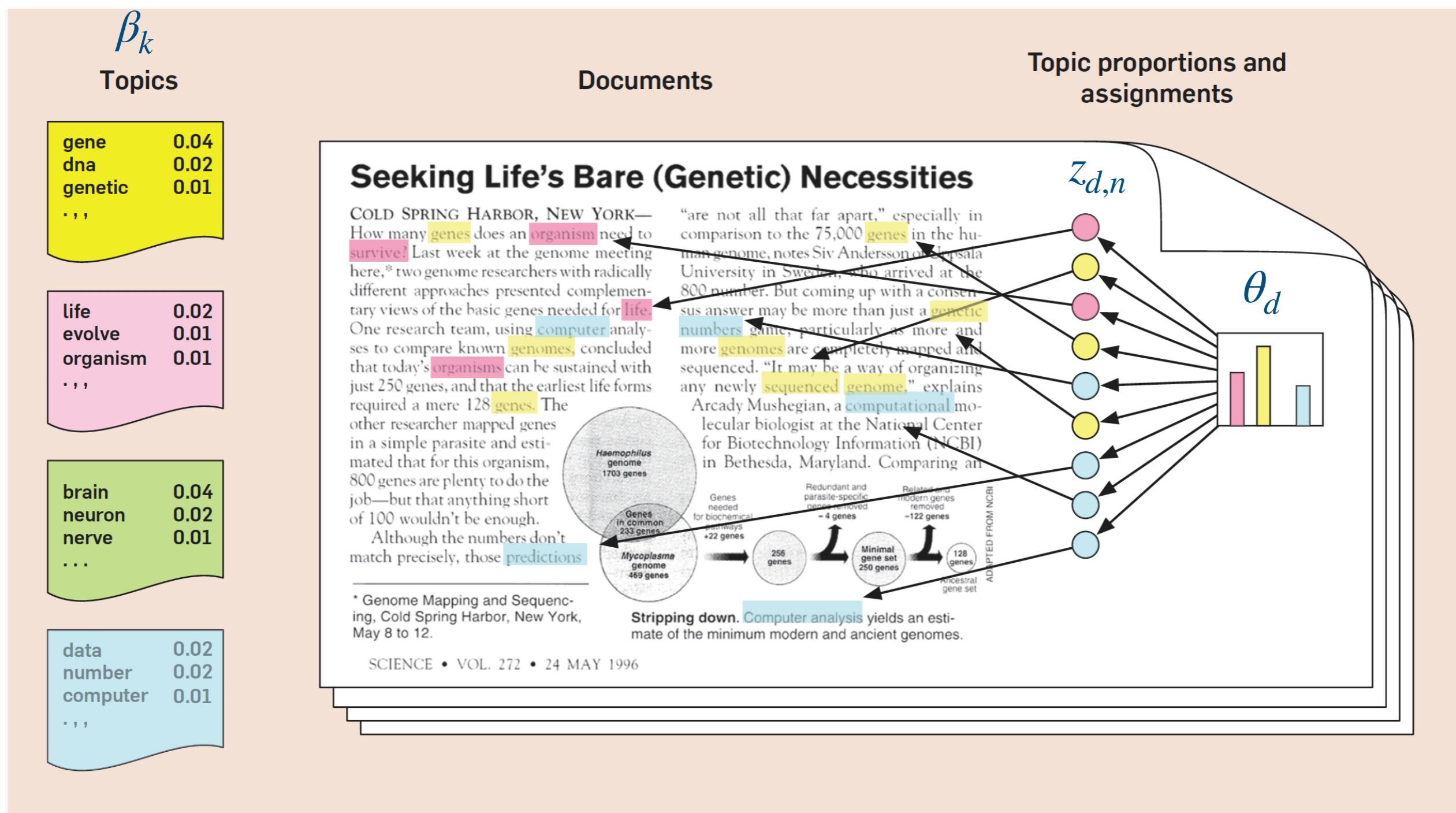
SCIENCE • VOL. 272 • 24 MAY 1996

Topic proportions and assignments

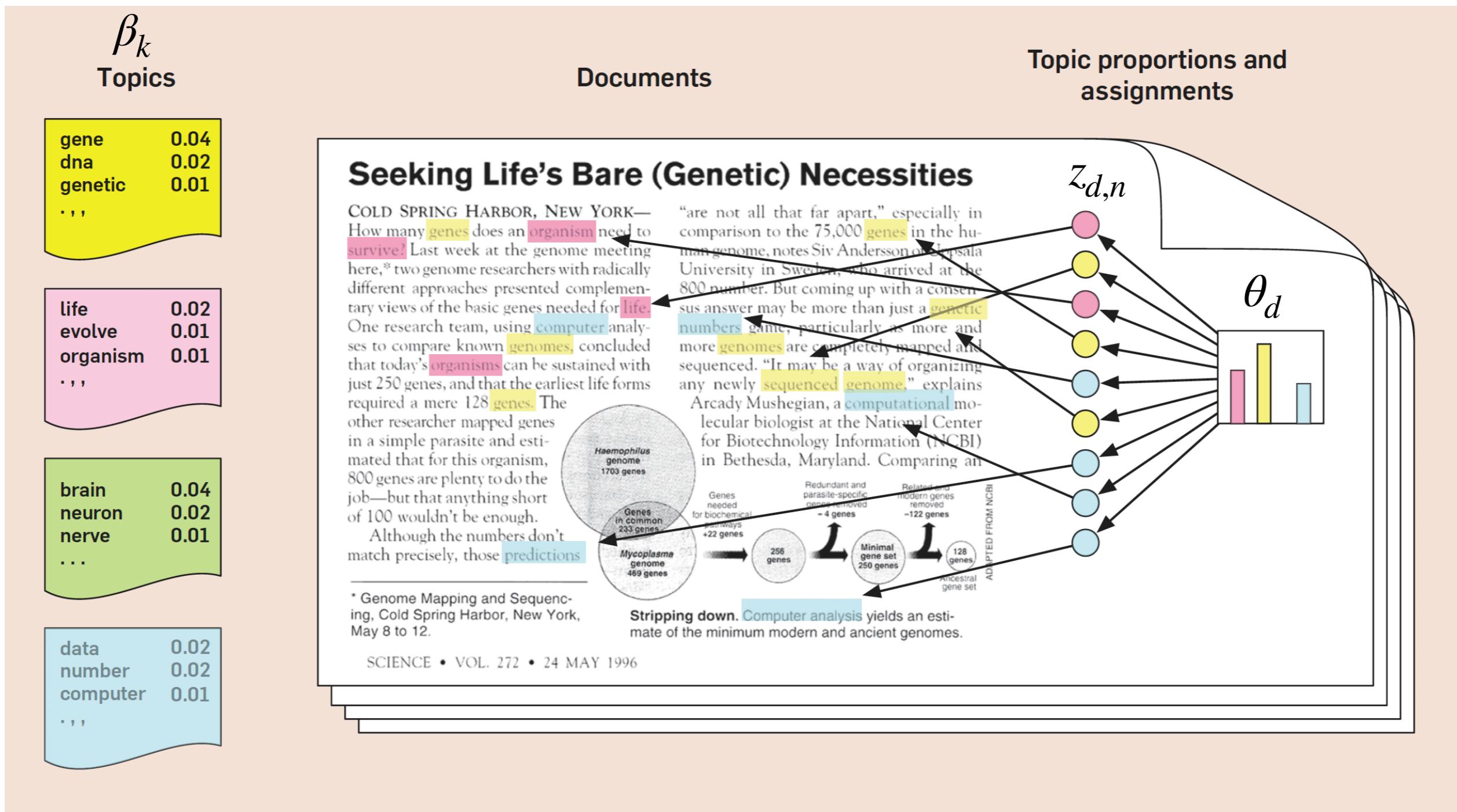


- Our goal is to infer or estimate the hidden variables, i.e. computing their distribution conditioned on the documents.

$$P(\text{topics, proportions, assignments} \mid \text{documents}) = P(\beta, z, \theta \mid w)$$

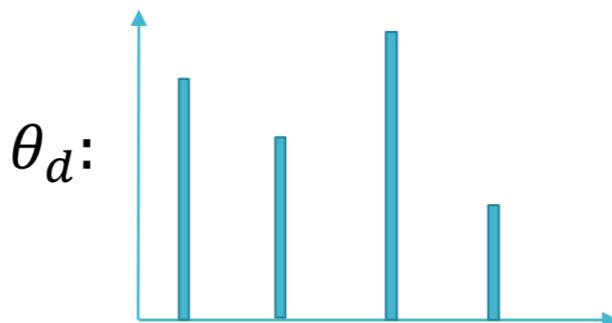


- Use posterior expectations ($E(\beta | w)$ for the corpus, $E(\theta | w)$ for each document) to perform the tasks: information retrieval, document similarity, exploration.

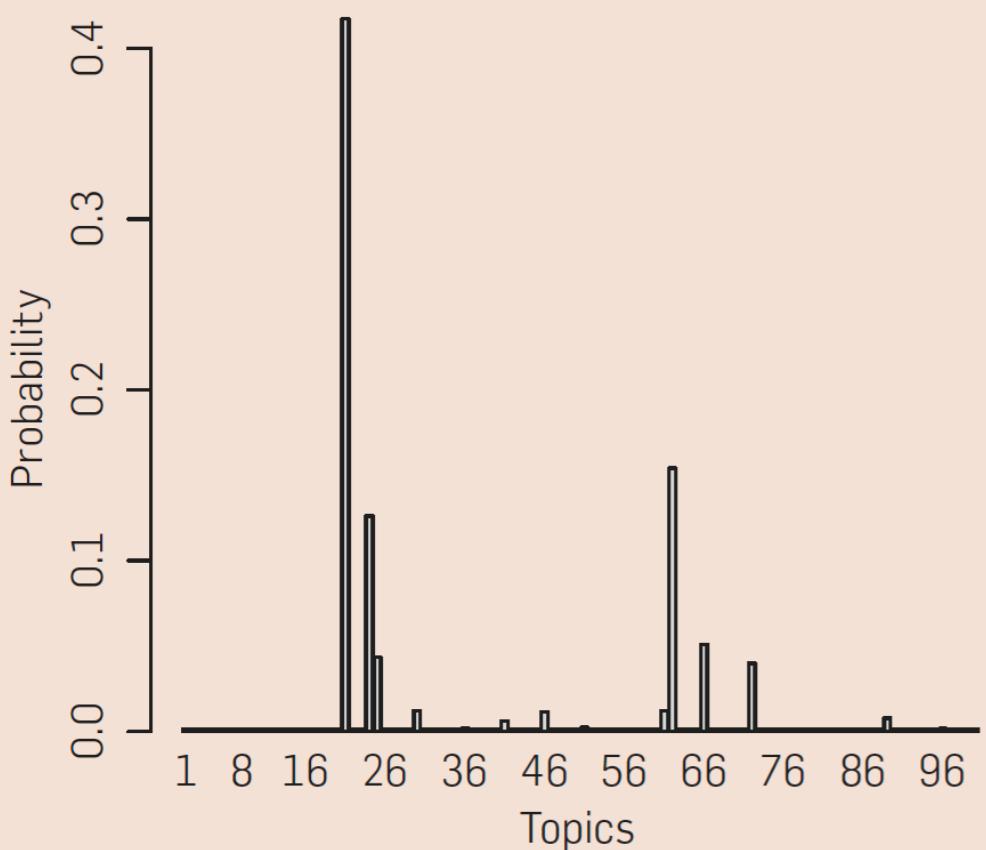


$p(\theta \mid w)$ Per document topic distributions

$p(\beta \mid w)$ Topic word distributions



Topics	Word probabilities for each topic		
	Topic 1	Topic 2	Topic 3
Topic 1			
Topic 2			
Topic 3			



“Genetics”

human
genome
dna
genetic
genes
sequence
gene
molecular
sequencing
map
information
genetics
mapping
project
sequences

“Evolution”

evolution
evolutionary
species
organisms
life
origin
biology
groups
phylogenetic
living
diversity
group
new
two
common

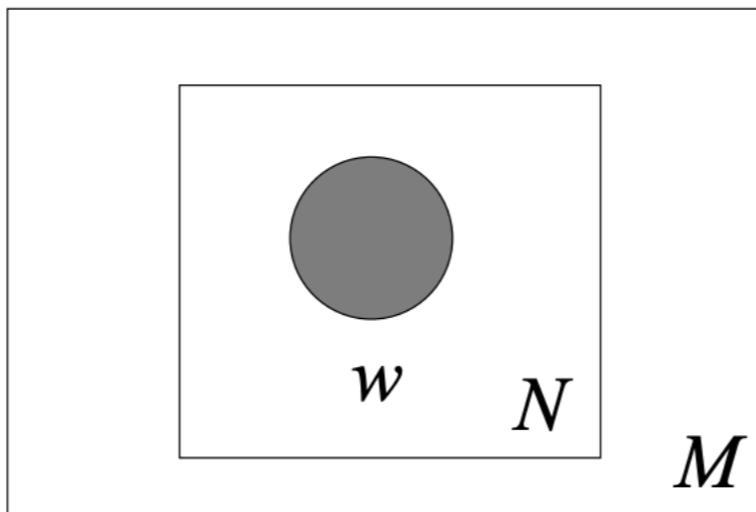
“Disease”

disease
host
bacteria
diseases
resistance
bacterial
new
strains
control
infectious
malaria
parasite
parasites
united
tuberculosis

“Computers”

computer
models
information
data
computers
system
network
systems
model
parallel
methods
networks
software
new
simulations

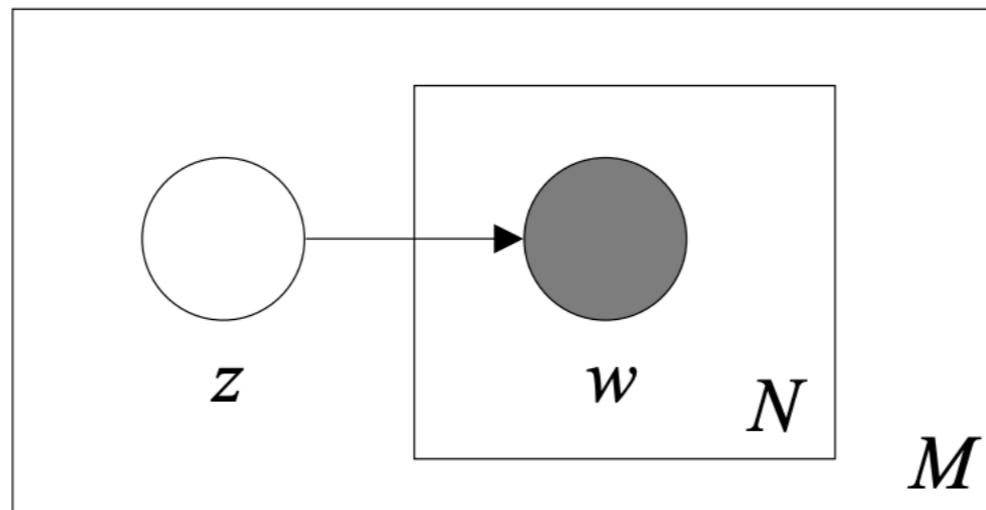
Unigram Model



- the words of every document are drawn independently from a single multinomial distribution

$$p(\mathbf{w}) = \prod_{n=1}^N p(w_n).$$

Mixture of Unigrams Model



- Argument the unigram model with a discrete random topic variables z and obtain a mixture of unigrams model.
- Each document is generated by first choosing a topic z and then generating words independently from the conditional multinomial

$$p(\mathbf{w}) = \sum_z p(z) \prod_{n=1}^N p(w_n | z).$$

$$p(\mathbf{w}) = \sum_z p(z) \prod_{n=1}^N p(w_n | z).$$

Topic word distributions

$$p(w | z_1)$$

$$z_1$$

$$p(w | z_2)$$

$$z_2$$

$$p(w | z_3)$$

$$z_3$$

$$p(w | z_4)$$

$$z_4$$

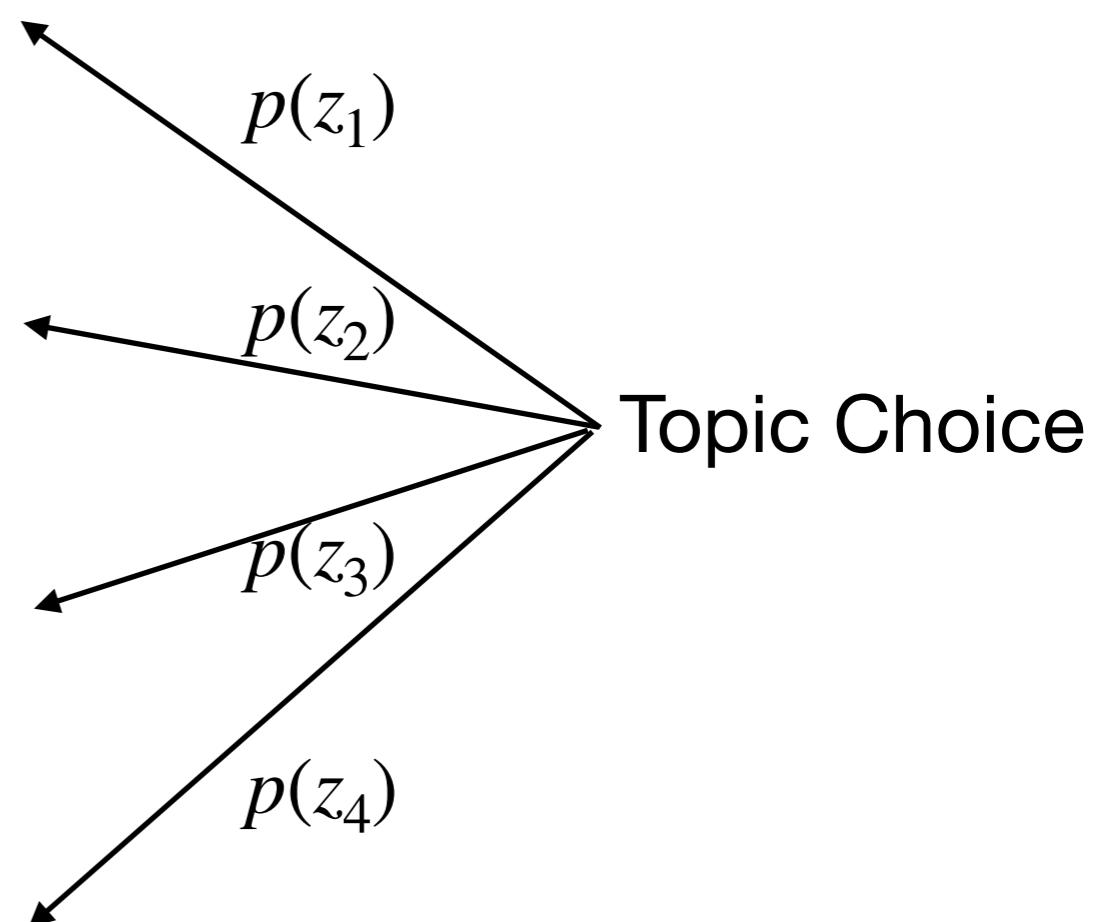
$$p(z_1)$$

$$p(z_2)$$

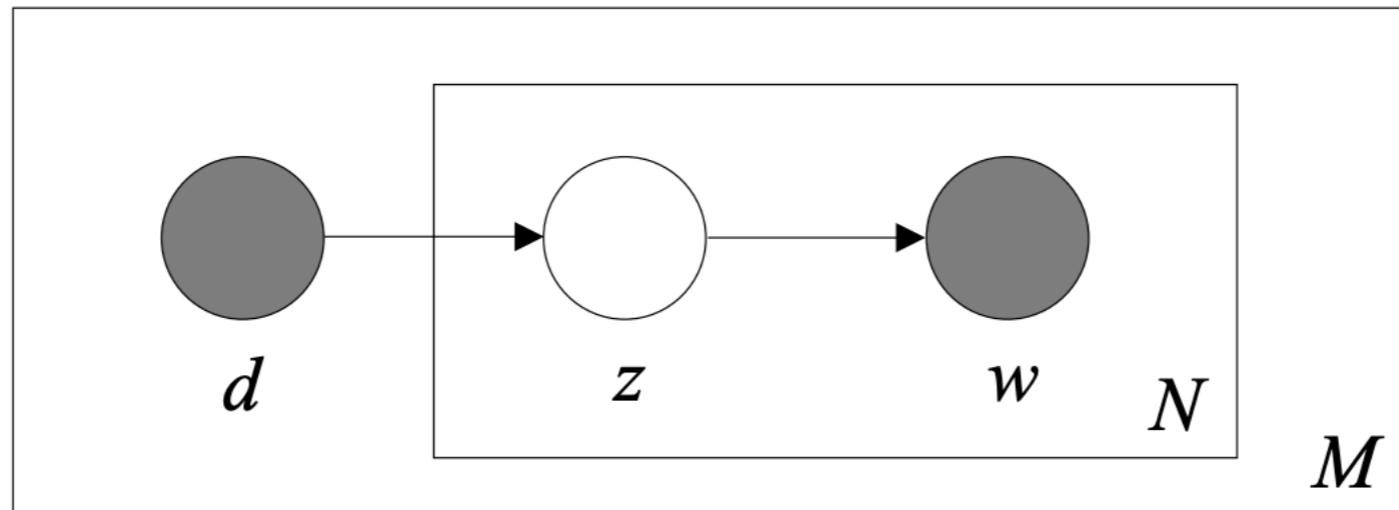
$$p(z_3)$$

$$p(z_4)$$

Topic Choice



pLSI

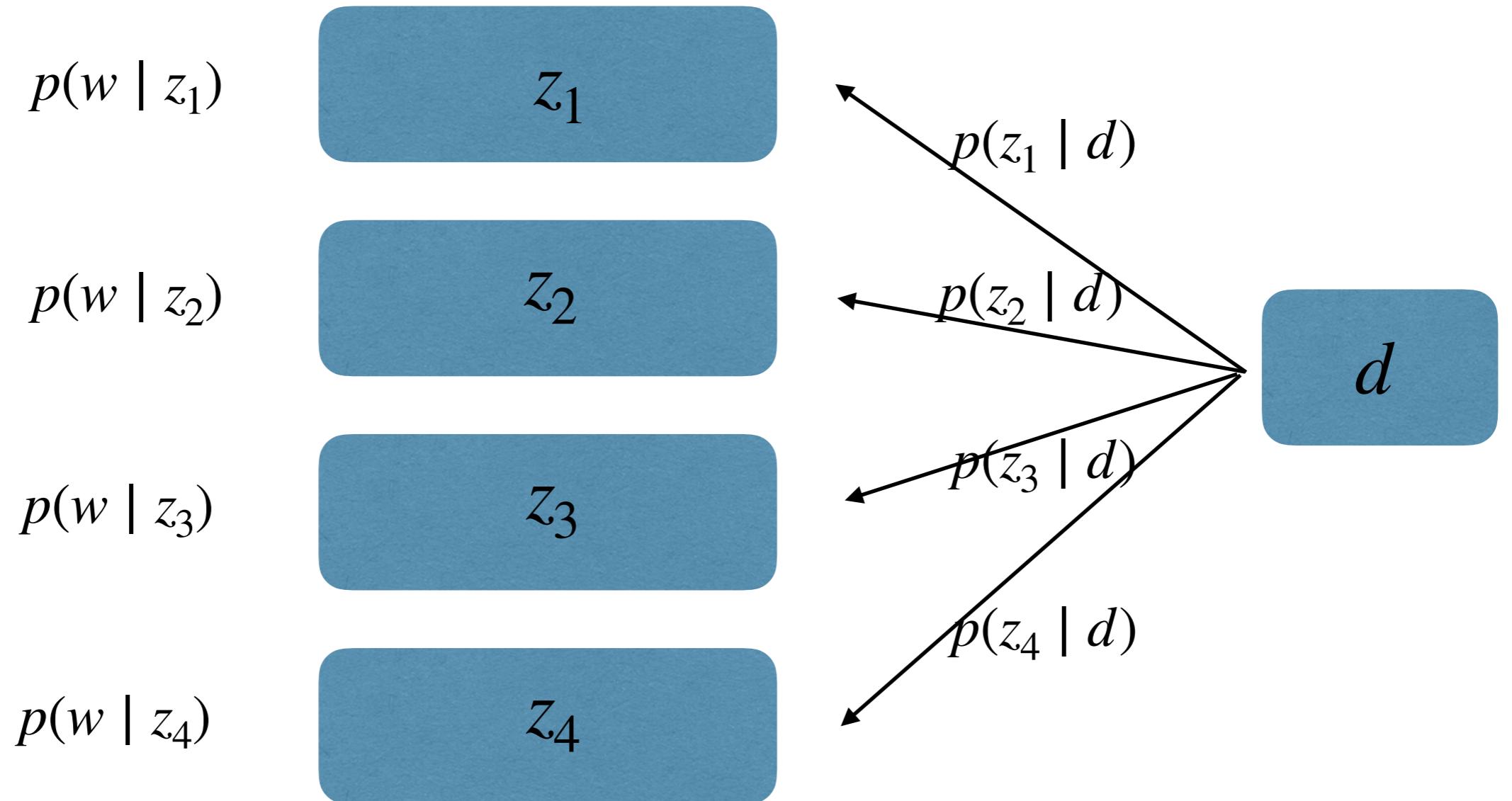


- The pLSI model attempts to relax the simplifying assumption made in the mixture of unigrams model that each document is generated from only one topic.
- Given all parameters, we want to infer the distribution z a word is from

$$p(d, w_n) = p(d) \sum_z p(w_n | z) p(z | d).$$

$$p(d, w_n) = p(d) \sum_z p(w_n | z) p(z | d).$$

Topic word distributions



pLSI Deficiency

K * V

Topic word distributions

$p(w | z_1)$

z_1

$p(w | z_2)$

z_2

$p(w | z_3)$

z_3

$p(w | z_4)$

z_4

K * M

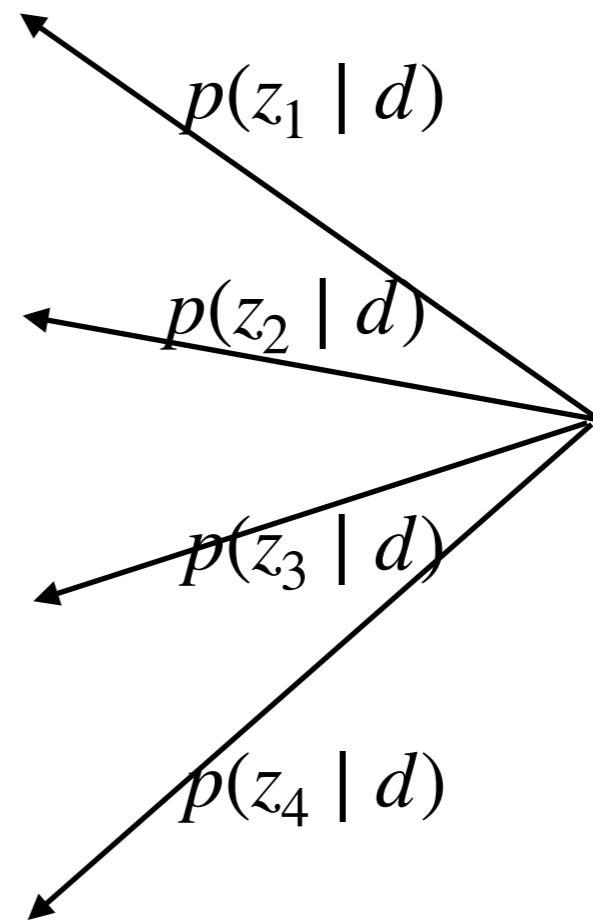
$p(z_1 | d)$

$p(z_2 | d)$

$p(z_3 | d)$

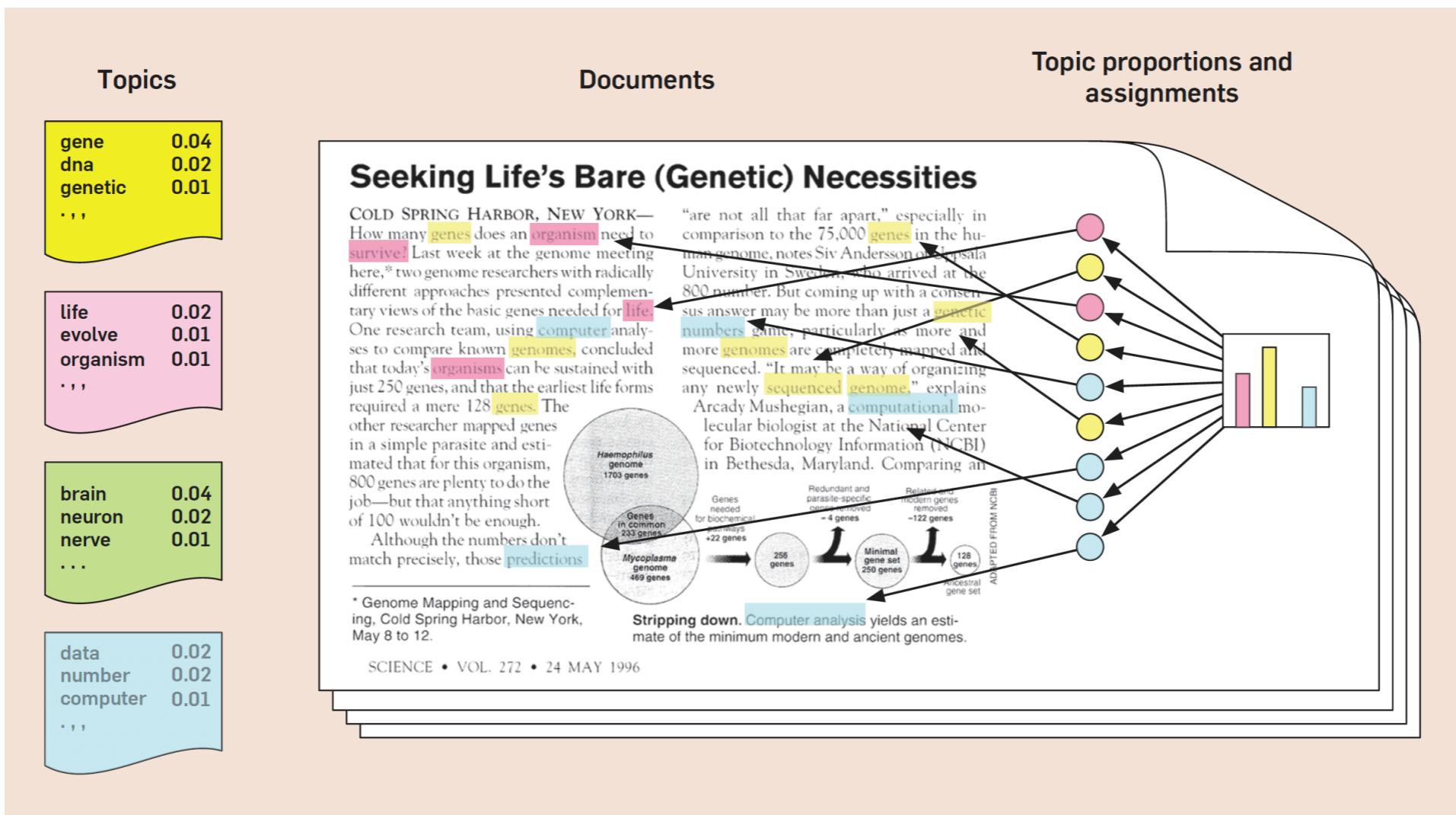
$p(z_4 | d)$

d

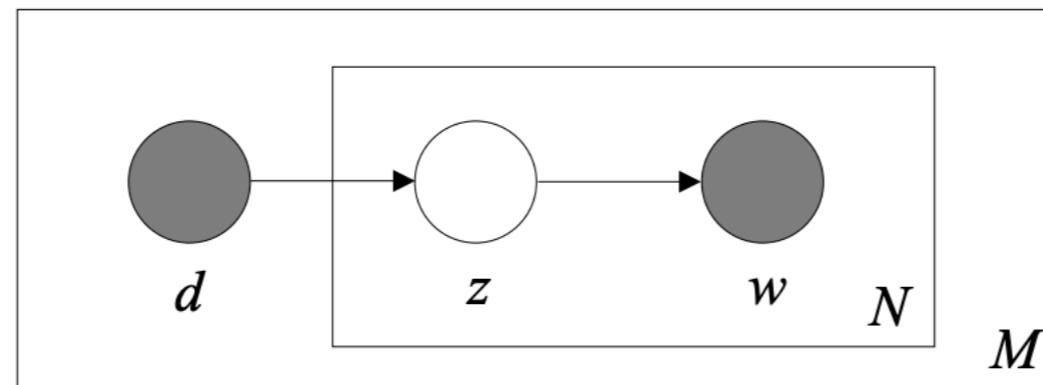


Exchangeability

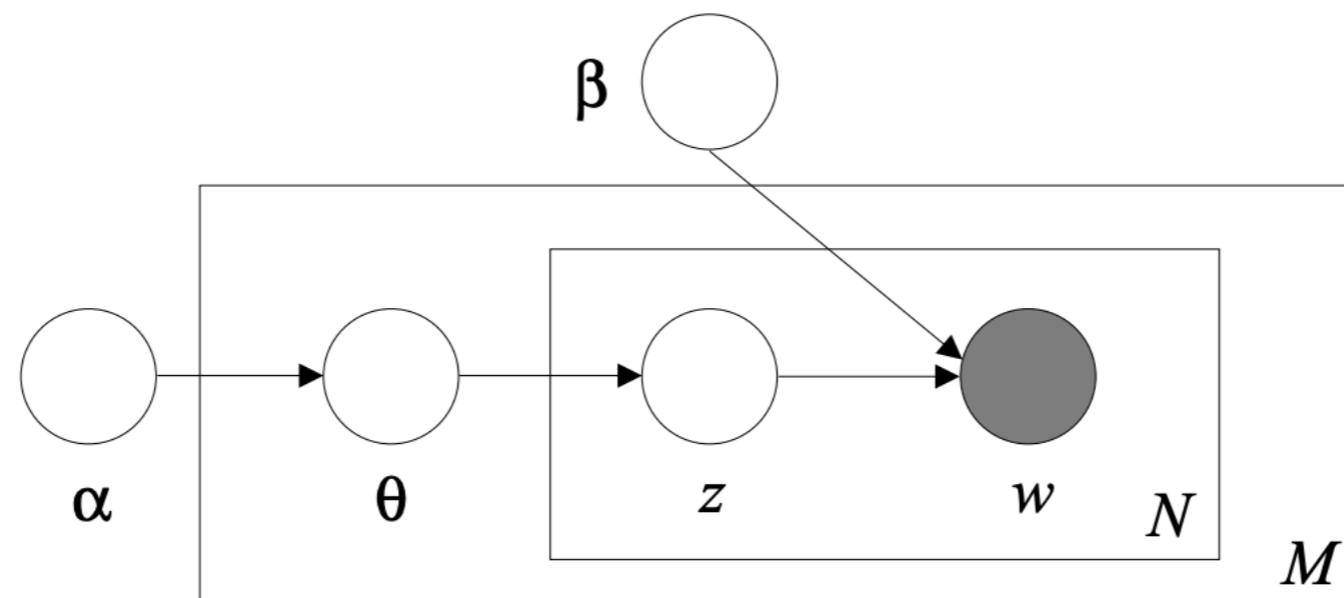
- De Finetti's theorem states that any set of exchangeable random variables has a representation as a mixture distribution.



LDA



- LDA improves pLSI to be a generative model by imposing Dirichlet priors on the model parameters.

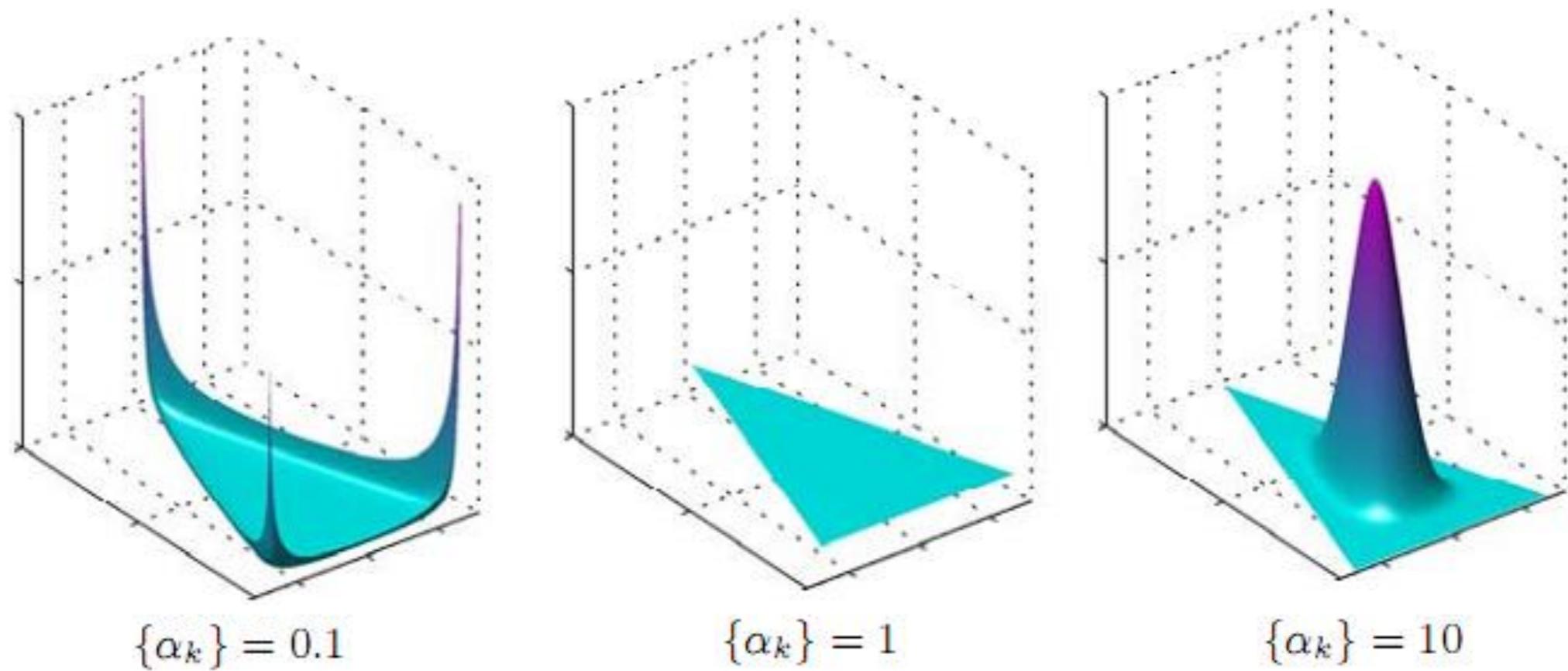


$$p(\mathbf{w}) = \int_{\theta} \left(\prod_{n=1}^N \sum_{z_n=1}^k p(w_n|z_n; \beta) p(z_n|\theta) \right) p(\theta; \alpha) d\theta,$$

Dirichlet Distribution

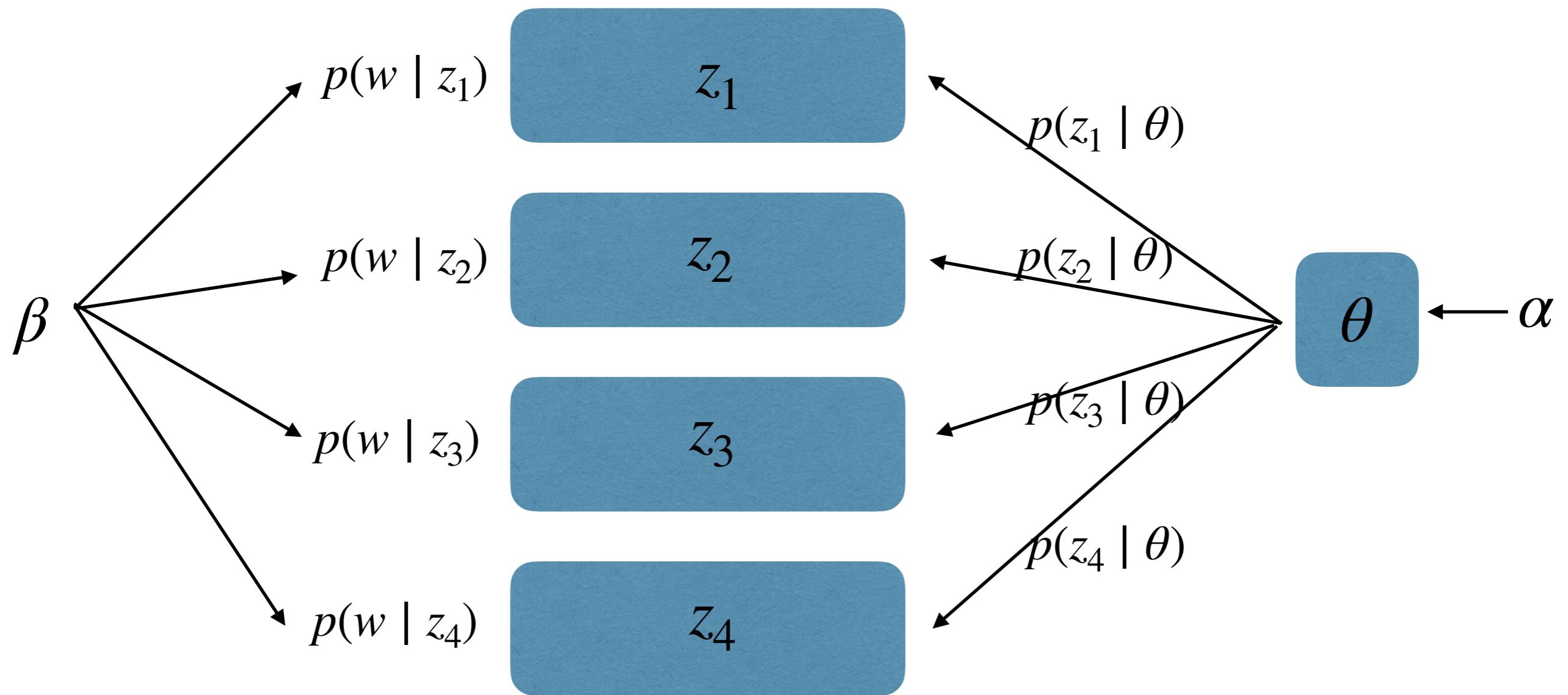
- The topic mixture proportions for each document are drawn from some distribution
- The space of all of these multinomials has a nice geometric interpretation as $(k-1)$ -simplex
- The distribution is defined over a $(k-1)$ -simplex. It takes k non-negative arguments which sum to one. A natural distribution to use over multinomial distribution
- Explain probability distribution of probability distribution

$$p(\theta | \alpha) = \frac{\Gamma\left(\sum_{i=1}^k \alpha_i\right)}{\prod_{i=1}^k \Gamma(\alpha_i)} \theta_1^{\alpha_1-1} \cdots \theta_k^{\alpha_k-1},$$



$$p(\mathbf{w}) = \int_{\theta} \left(\prod_{n=1}^N \sum_{z_n=1}^k p(w_n | z_n; \beta) p(z_n | \theta) \right) p(\theta; \alpha) d\theta,$$

Topic word distributions



Continuous mixture of unigram

Mixture of unigram

$$p(\mathbf{w} | \alpha, \beta) = \int p(\theta | \alpha) \left(\prod_{n=1}^N \sum_{z_n} p(z_n | \theta) p(w_n | z_n, \beta) \right) d\theta.$$

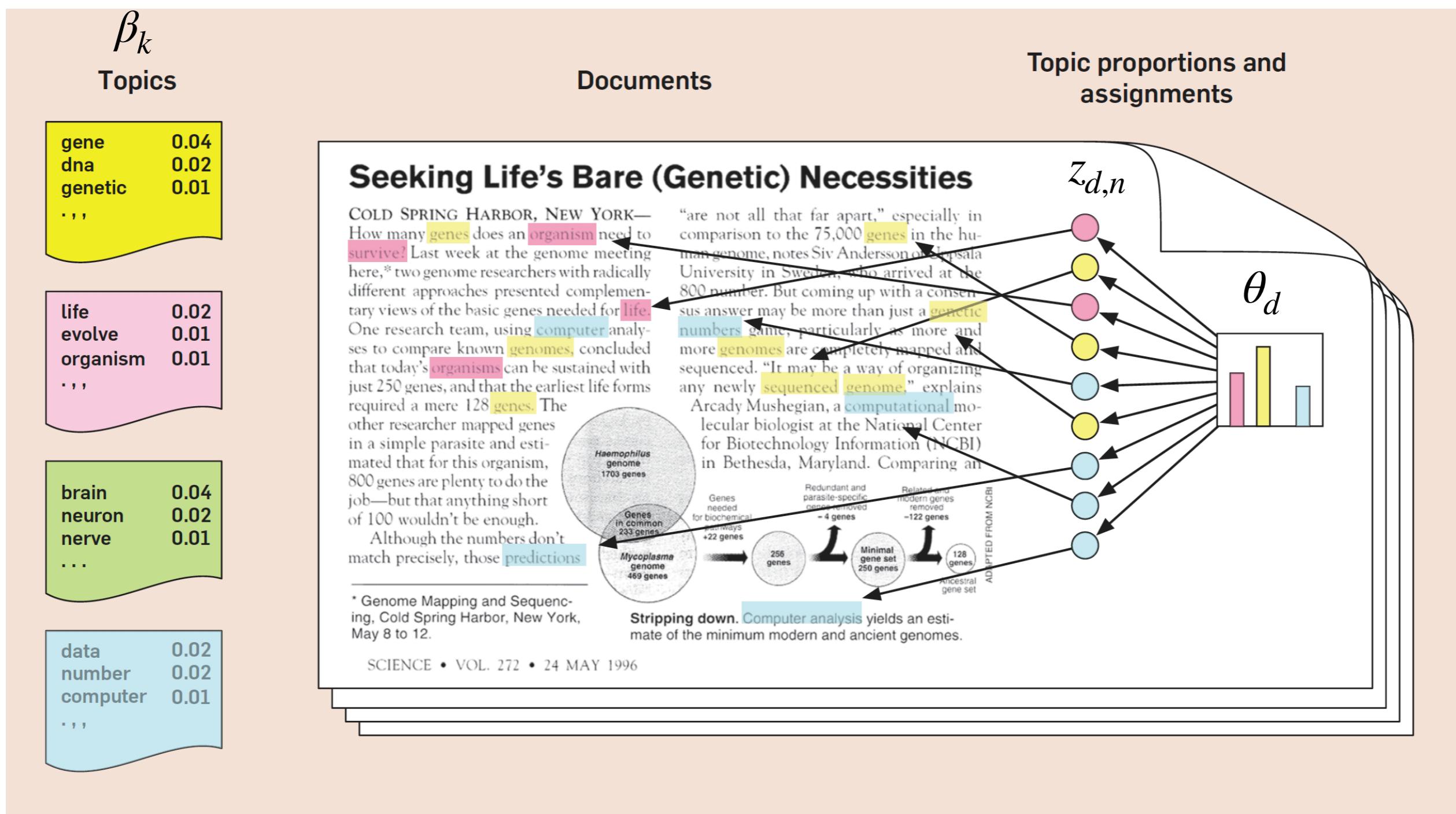
- Integral over all possibilities because document topic distribution is not fixed anymore.

Marginal likelihood

$$p(\mathcal{D} | \alpha, \beta) = \prod_{d=1}^M \int p(\theta_d | \alpha) \left(\prod_{n=1}^{N_d} \sum_{z_{dn}} p(z_{dn} | \theta_d) p(w_{dn} | z_{dn}, \beta) \right) d\theta_d.$$

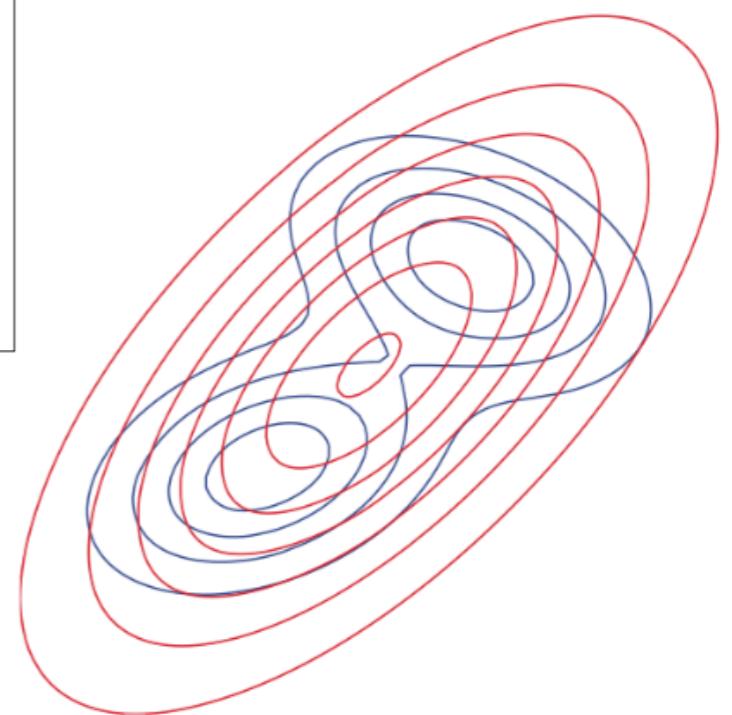
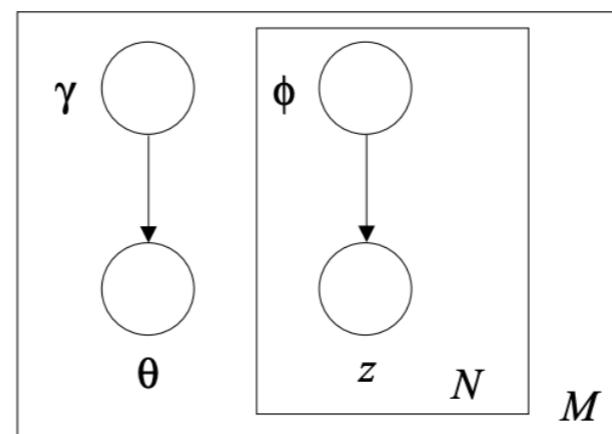
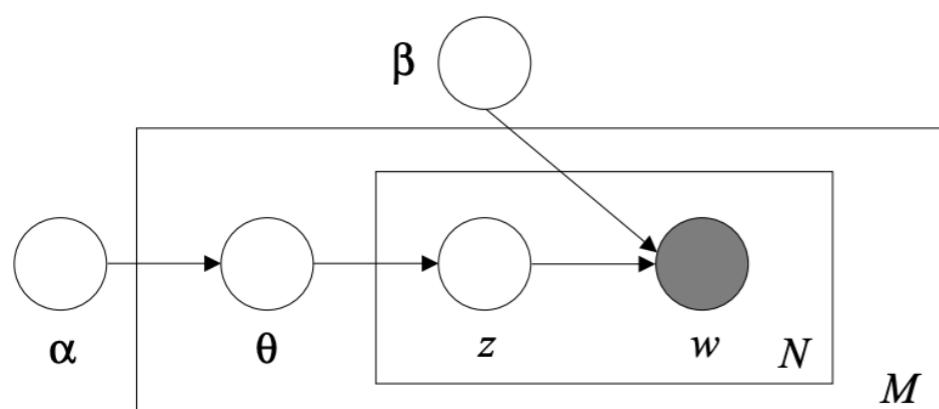
- Our goal is to infer or estimate the hidden variables, i.e. computing their distribution conditioned on the documents.

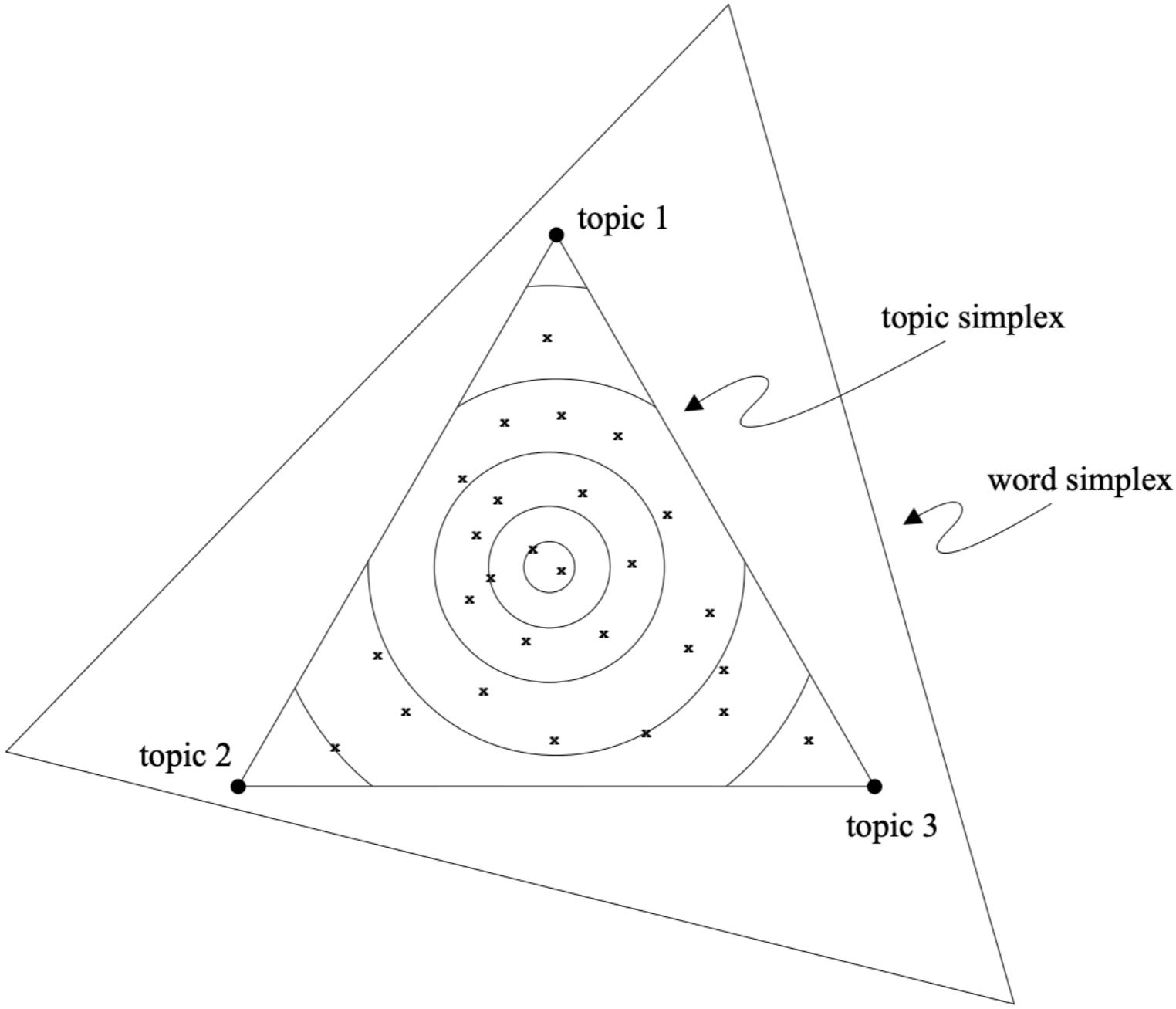
$$P(\text{topics, proportions, assignments} \mid \text{documents}) = P(\beta, z, \theta \mid w)$$



Variation Inference

- The key inferential problem that we need to solve in order to use LDA is that of computing the posterior distribution of the hidden variables given a document
- Unfortunately, this distribution is intractable to compute in general



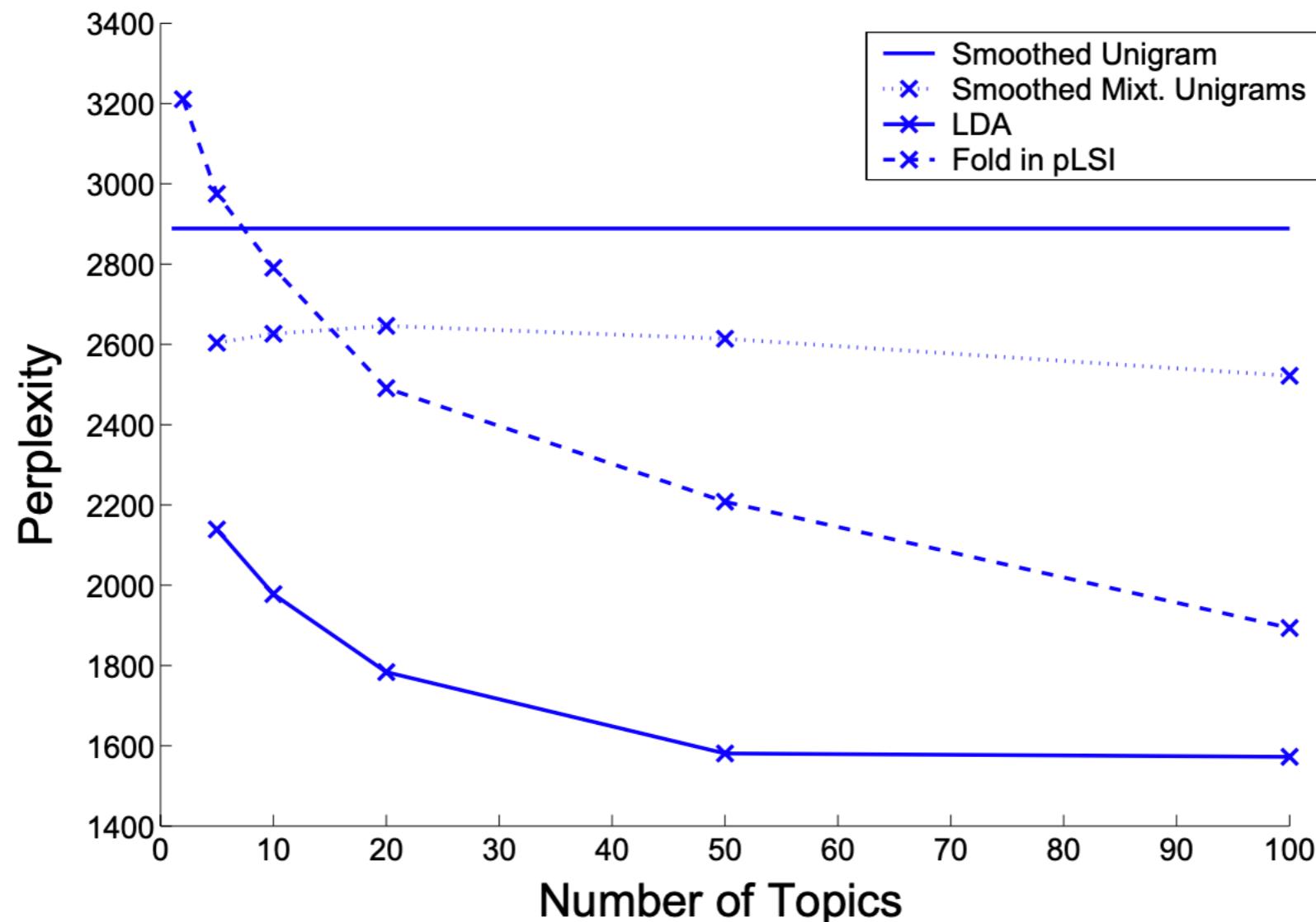


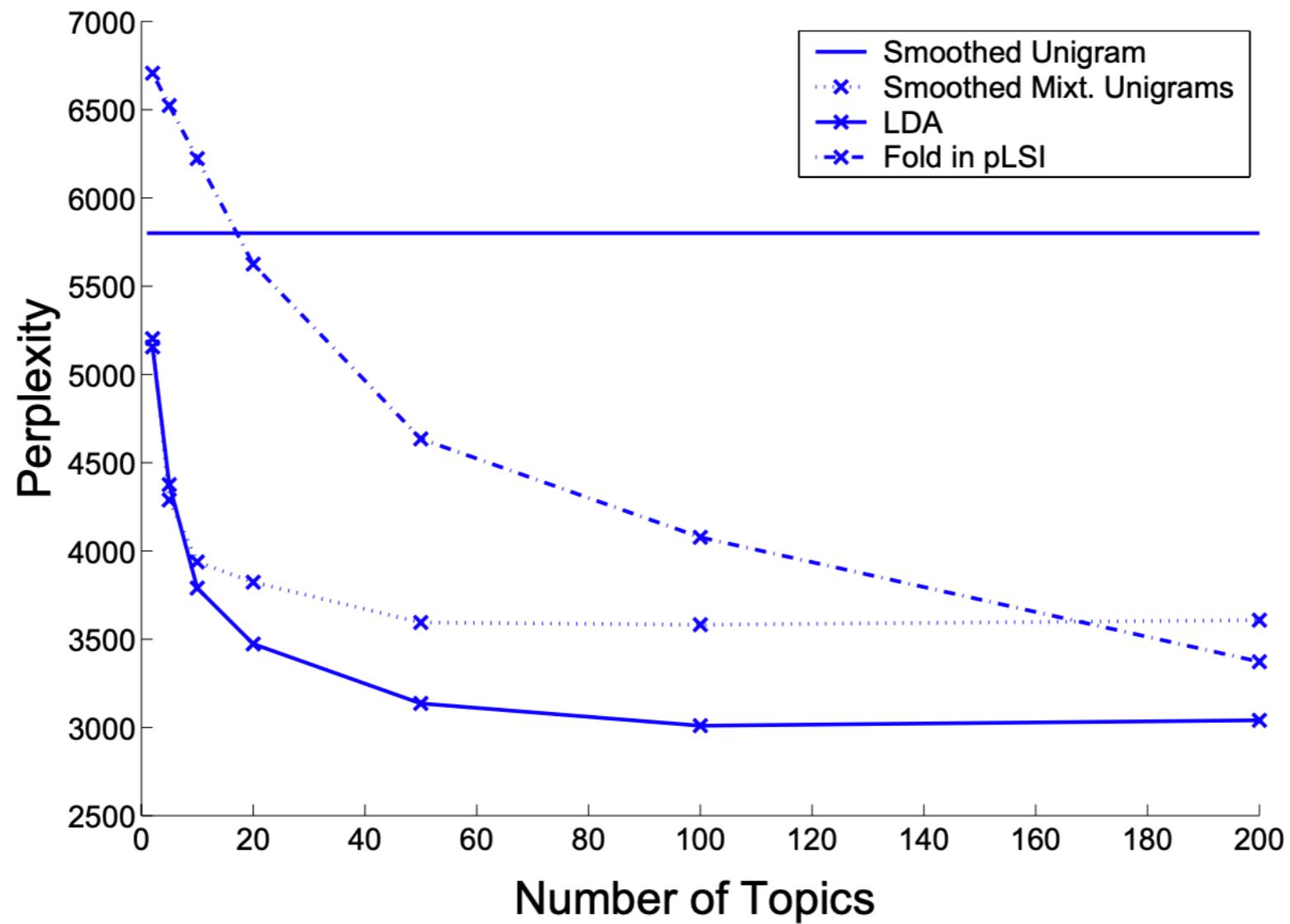
The topic simplex for three topics embedded in the word simplex for three words. The corners of the word simplex correspond to the three distributions where each word (respectively) has probability one. The three points of the topic simplex correspond to three different distributions over words. The mixture of unigrams places each document at one of the corners of the topic simplex. The pLSI model induces an empirical distribution on the topic simplex denoted by x . LDA places a smooth distribution on the topic simplex denoted by the contour lines.

“Arts”	“Budgets”	“Children”	“Education”
NEW	MILLION	CHILDREN	SCHOOL
FILM	TAX	WOMEN	STUDENTS
SHOW	PROGRAM	PEOPLE	SCHOOLS
MUSIC	BUDGET	CHILD	EDUCATION
MOVIE	BILLION	YEARS	TEACHERS
PLAY	FEDERAL	FAMILIES	HIGH
MUSICAL	YEAR	WORK	PUBLIC
BEST	SPENDING	PARENTS	TEACHER
ACTOR	NEW	SAYS	BENNETT
FIRST	STATE	FAMILY	MANIGAT
YORK	PLAN	WELFARE	NAMPHY
OPERA	MONEY	MEN	STATE
THEATER	PROGRAMS	PERCENT	PRESIDENT
ACTRESS	GOVERNMENT	CARE	ELEMENTARY
LOVE	CONGRESS	LIFE	HAITI

The William Randolph Hearst Foundation will give \$1.25 million to Lincoln Center, Metropolitan Opera Co., New York Philharmonic and Juilliard School. “Our board felt that we had a real opportunity to make a mark on the future of the performing arts with these grants an act every bit as important as our traditional areas of support in health, medical research, education and the social services,” Hearst Foundation President Randolph A. Hearst said Monday in announcing the grants. Lincoln Center’s share will be \$200,000 for its new building, which will house young artists and provide new public facilities. The Metropolitan Opera Co. and New York Philharmonic will receive \$400,000 each. The Juilliard School, where music and the performing arts are taught, will get \$250,000. The Hearst Foundation, a leading supporter of the Lincoln Center Consolidated Corporate Fund, will make its usual annual \$100,000 donation, too.

Figure 8: An example article from the AP corpus. Each color codes a different factor from which the word is putatively generated.







Mixture model parameters: $\Lambda = (\{p(w | z)\}, \{p(z)\})$

Likelihood: $p(C | \Lambda)$

ML estimate: $\hat{\Lambda} = \arg \max_{\Lambda} p(C | \Lambda)$

Mixture model parameters: $\Lambda = (\{p(w | z)\}, \{p(z | d)\})$

Likelihood: $p(C | \Lambda)$

ML estimate: $\hat{\Lambda} = \arg \max_{\Lambda} p(C | \Lambda)$