

Information Retrieval & Data Mining [COMP0084]

Topic models and vector semantics

Vasileios Lampos

Computer Science, UCL



lampos.net

@lampos



- ▶ In these lectures:
 - Introduction to topic models
 - Introduction to vector semantics
- ▶ Useful additional material
 - “*Speech and language processing*” (Jurafsky, Martin), web.stanford.edu/~jurafsky/slp3/
 - pLSA (Hofmann), iro.umontreal.ca/~nie/IFT6255/Hofmann-UAI99.pdf
 - LDA (Blei, Ng, Jordan), jmlr.org/papers/volume3/blei03a/blei03a.pdf
 - word2vec (Mikolov et al.), arxiv.org/abs/1301.3781
 - Blei on LDA, youtube.com/watch?v=DDq3OVp9dNA
 - Boyd-Graber on topic models, youtube.com/watch?v=yK7nN3FcgUs
 - Manning on word2vec, youtube.com/watch?v=ERibwqs9p38
- ▶ Some slides adapted from WSDM ’14 tutorial on “Multilingual Probabilistic Topic Modelling” — liir.cs.kuleuven.be/tutorial/WSDM2014Tutorial.pdf

What is a topic model?

What is a topic model?

- ▶ **Informally:** groupings (or clusters) of words (terms, n -grams) that are somehow related

What is a topic model?

- ▶ **Informally:** groupings (or clusters) of words (terms, n -grams) that are somehow related
- ▶ **Still informally:** method for automatically organising, understanding, searching, and summarising large (digitised) document collections
 - uncovers hidden (latent) topical patterns (**topics!**) in the collection
 - can annotate, and then organise or summarise, the documents based on these topics

What is a topic model?

- ▶ **Informally:** groupings (or clusters) of words (terms, n -grams) that are somehow related
- ▶ **Still informally:** method for automatically organising, understanding, searching, and summarising large (digitised) document collections
 - uncovers hidden (latent) topical patterns (**topics!**) in the collection
 - can annotate, and then organise or summarise, the documents based on these topics
- ▶ As we will see, it is often defined as a **probabilistic structure** expressing a certain set of assumptions about how the documents in our collection were generated

What is a topic model?

- ▶ **Informally:** groupings (or clusters) of words (terms, n -grams) that are somehow related
- ▶ **Still informally:** method for automatically organising, understanding, searching, and summarising large (digitised) document collections
 - uncovers hidden (latent) topical patterns (**topics!**) in the collection
 - can annotate, and then organise or summarise, the documents based on these topics
- ▶ As we will see, it is often defined as a **probabilistic structure** expressing a certain set of assumptions about how the documents in our collection were generated
- ▶ **Note:** we can also learn topic models (word clusters) using clustering techniques with no explicit probabilistic structure

Why do need topics?



Why do need topics?

- ▶ Too many documents and we can't read them all!



Why do need topics?

- ▶ Too many documents and we can't read them all!
- ▶ Topic models can automatically categorise large document collections, so that we can browse through them much more efficiently
- ▶ Applicable on various corpus collections attracting inter-disciplinary interest
e.g. newspapers, books, social media, health reports



Why do need topics?

- ▶ Too many **documents** and we can't read them all!
- ▶ Topic models can automatically **categorise** large document collections, so that we can browse through them much more efficiently
- ▶ Applicable on **various corpus collections** attracting inter-disciplinary interest
 - e.g. newspapers, books, social media, health reports
- ▶ Can **improve natural language processing tasks**
 - e.g. machine translation, word sense disambiguation
- ▶ Can **improve downstream tasks** in text mining
- ▶ Let's see a few **examples**



Topics in news articles

Blei, Ng & Jordan. JMLR, 2003
[jmlr.org/papers/volume3/
blei03a/blei03a.pdf](http://jmlr.org/papers/volume3/blei03a/blei03a.pdf)

“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.

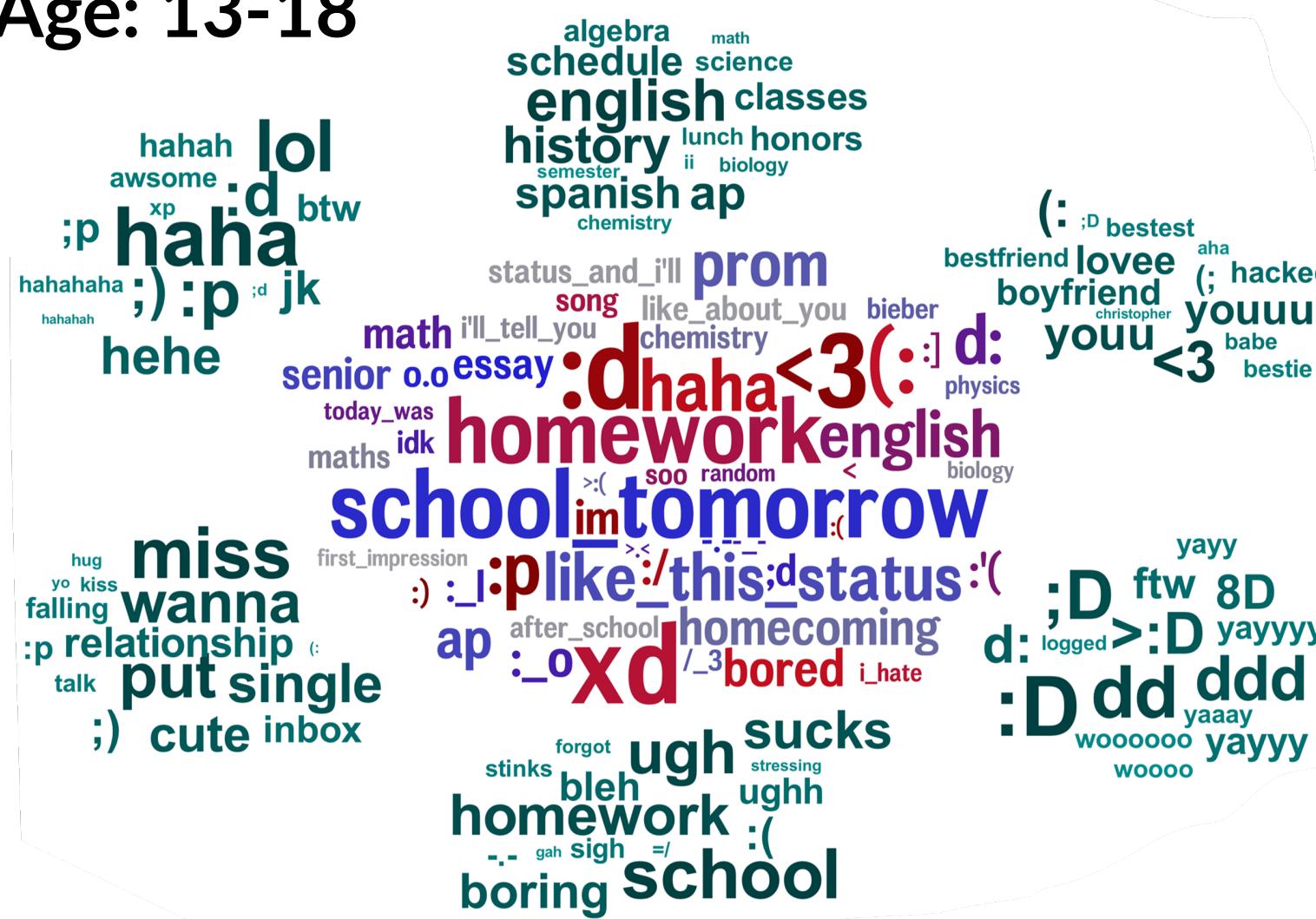
17,000 articles from the journal “Science”

“Genetics”	“Evolution”	“Disease”	“Computers”
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

Blei. CACM, 2012
doi.org/10.1145/2133806.2133826

Age-group topics on Facebook

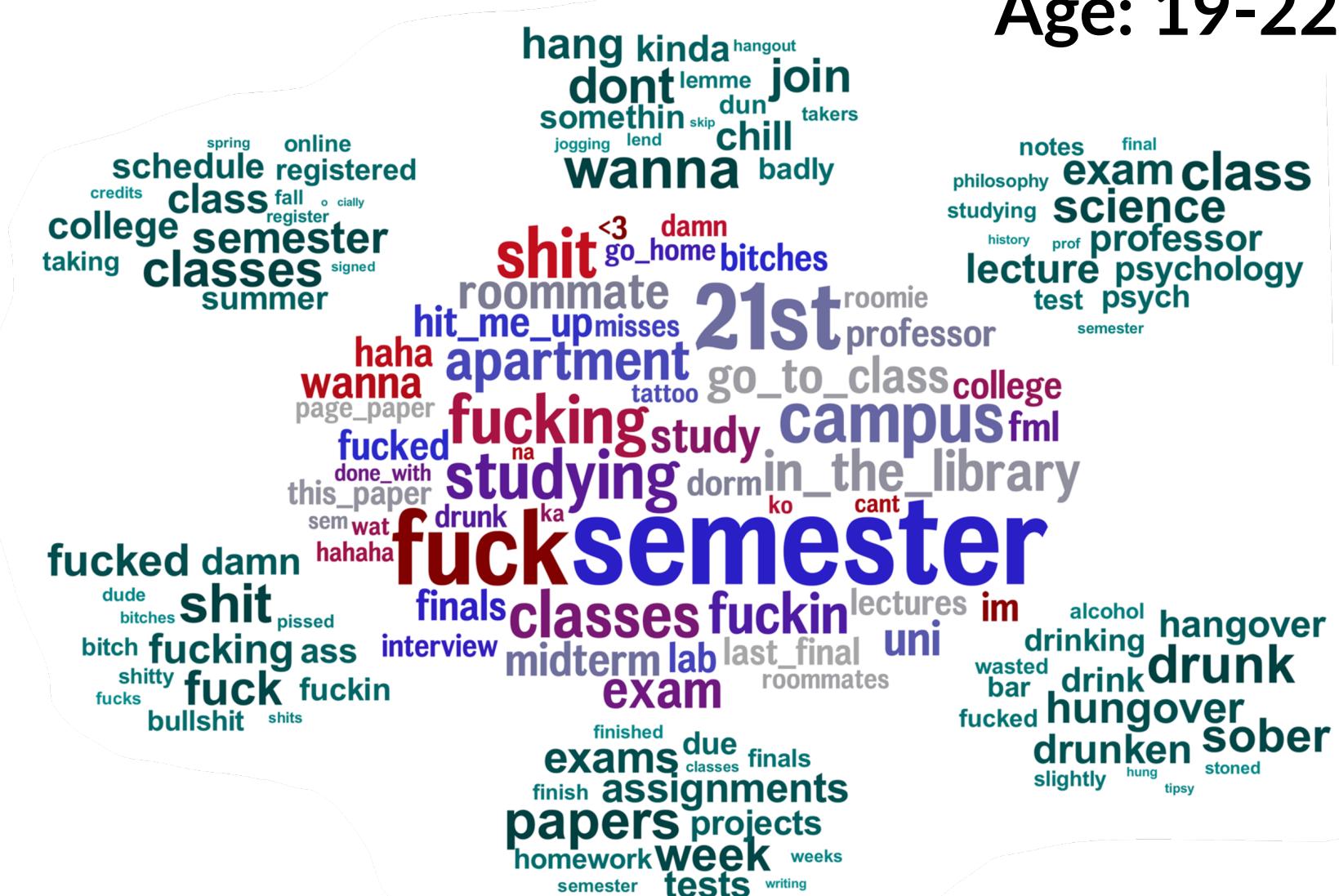
Age: 13-18



Age: 23-29

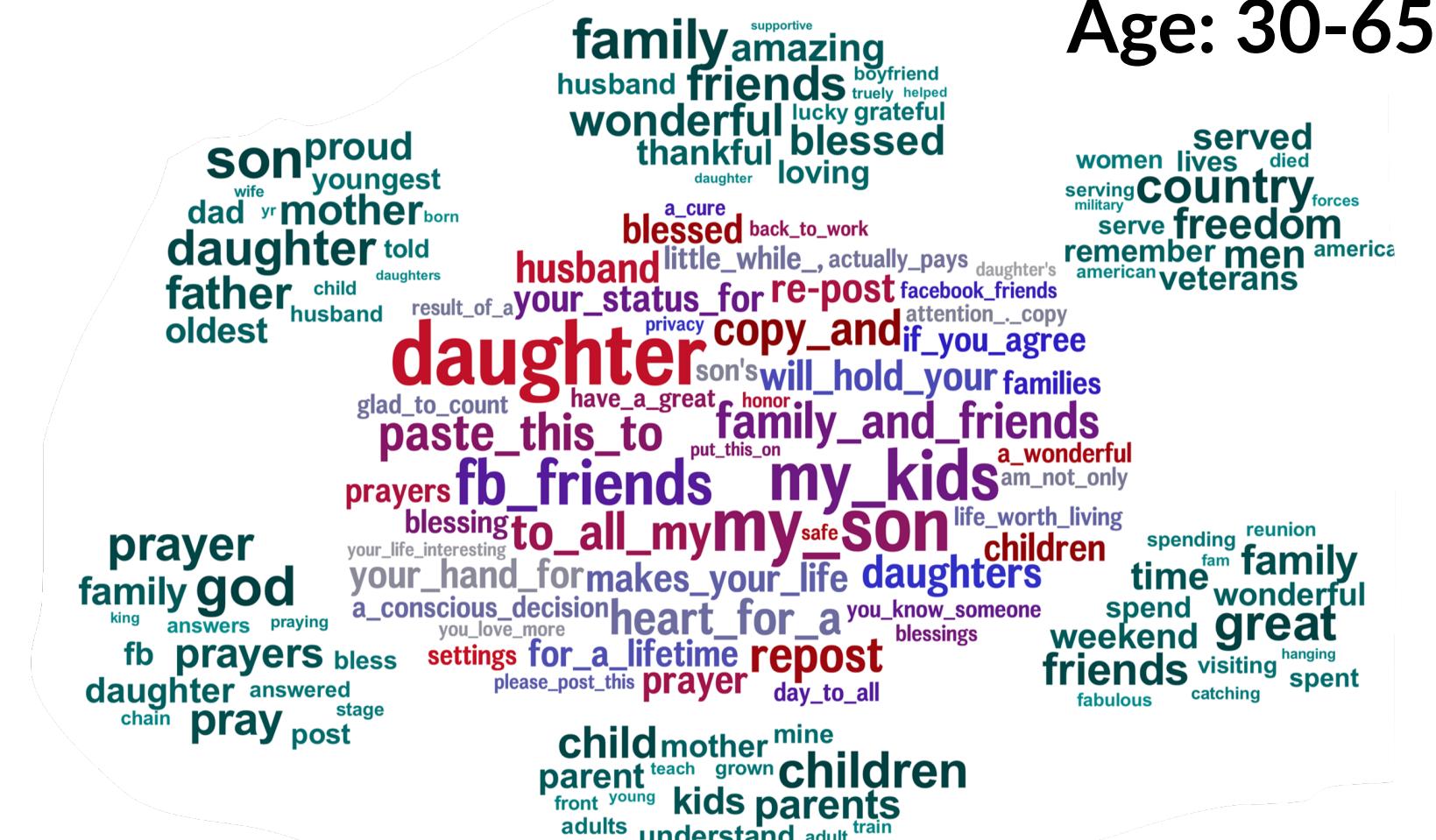


Age: 19-22



Schwartz et al. PLOS ONE, 2013
doi.org/10.1371/journal.pone.0073791

Age: 30-65

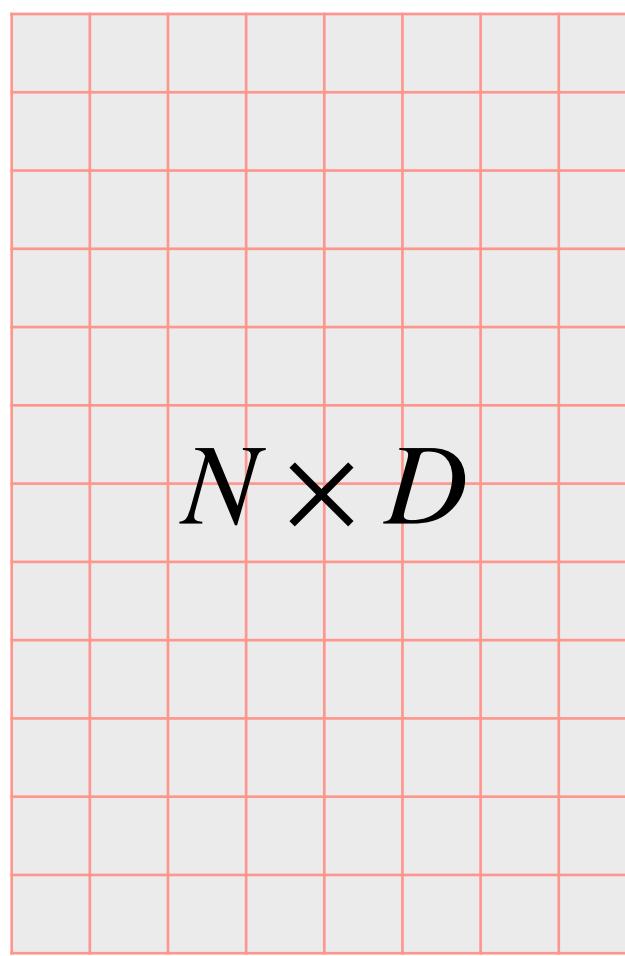


Predicting judicial decisions

Label	Words	w
	Violation of Article 3 that prohibits inhuman treatment	
Positive State Obligations	Top-5 Violation injury, protection, ordered, damage, civil, caused, failed, claim, course, connection, region, effective, quashed, claimed, suffered, suspended, carry, compensation, pecuniary, ukraine	13.50
Detention conditions	prison, detainee, visit, well, regard, cpt, access, food, situation, problem, remained, living, support, visited, establishment, standard, admissibility merit, overcrowding, contact, good	11.70
Treatment by state officials	police, officer, treatment, police officer, July, ill, force, evidence, ill treatment, arrest, allegation, police station, subjected, arrested, brought, subsequently, allegedly, ten, treated, beaten	10.20
	Top-5 No Violation	
Prior Violation of Article 2	june, statement, three, dated, car, area, jurisdiction, gendarmerie, perpetrator, scene, June applicant, killing, prepared, bullet, wall, weapon, kidnapping, dated June, report dated, stopped	-12.40
Issues of Proof	witness, asked, told, incident, brother, heard, submission, arrived, identity, hand, killed, called, involved, started, entered, find, policeman, returned, father, explained	-15.20
Sentencing	sentence, year, life, circumstance, imprisonment, release, set, president, administration, sentenced, term, constitutional, federal, appealed, twenty, convicted, continued, regime, subject, responsible	-17.40

Aletras, Tsarapatsanis, Preotiuc, Lampos. PeerJ Computer Science, 2016
doi.org/10.7717/peerj-cs.93

Latent Semantic Analysis (or Indexing) – LSA



X

Singular Value Decomposition (SVD; truncated) on the term-document matrix **X** representing the frequency of N terms (or n -grams) in D documents

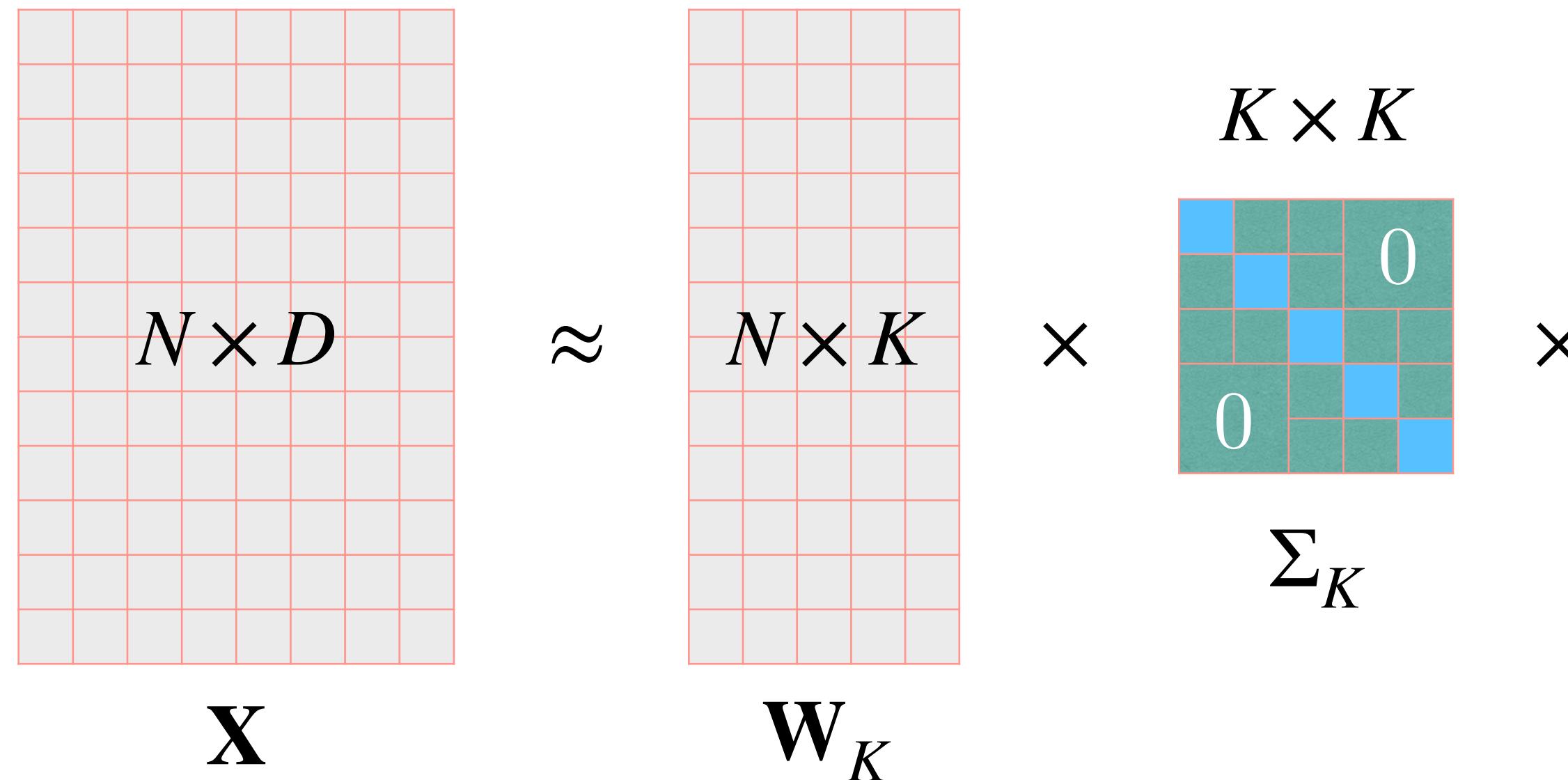
Latent Semantic Analysis (or Indexing) – LSA

The diagram illustrates a mathematical operation involving three components: X , \approx , and W_K . Each component is represented by a light gray square grid with a red rectangular border. The grid for X has 8 horizontal and 8 vertical red lines, creating a 7x7 internal area. The grid for W_K also has 8 horizontal and 8 vertical red lines, creating a 7x7 internal area. The symbol \approx is positioned between the first grid and the second grid. To the right of the second grid, there is a partial view of a third grid, which appears to have 8 horizontal and 8 vertical red lines, creating a 7x7 internal area. This visual representation suggests that the product of X and W_K is approximately equal to the product of two matrices, one of size $N \times D$ and the other of size $N \times K$.

Singular Value Decomposition (SVD; truncated) on the term-document matrix \mathbf{X}
representing the frequency of N terms (or n -grams) in D documents

\mathbf{W}_K : each topic's (K topics) distribution over N terms

Latent Semantic Analysis (or Indexing) – LSA

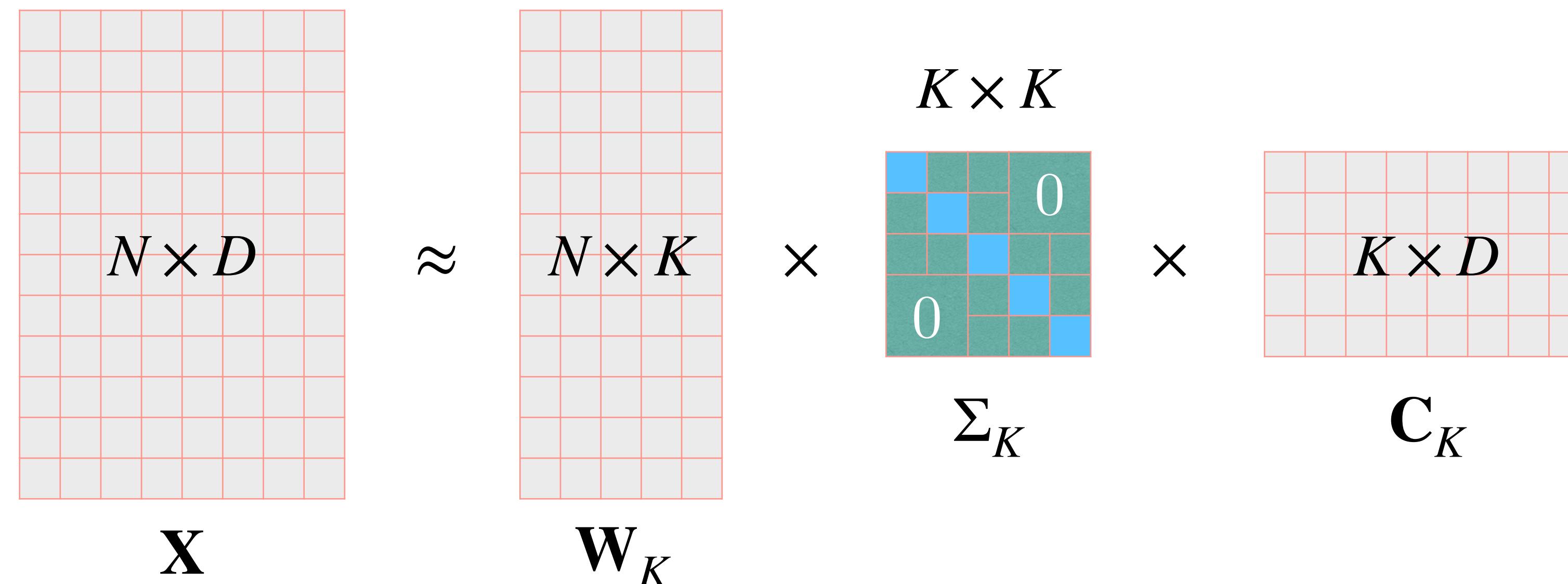


Singular Value Decomposition (SVD; truncated) on the term-document matrix \mathbf{X}
representing the frequency of N terms (or n -grams) in D documents

\mathbf{W}_K : each topic's (K topics) distribution over N terms

Σ_K : topic importance

Latent Semantic Analysis (or Indexing) – LSA



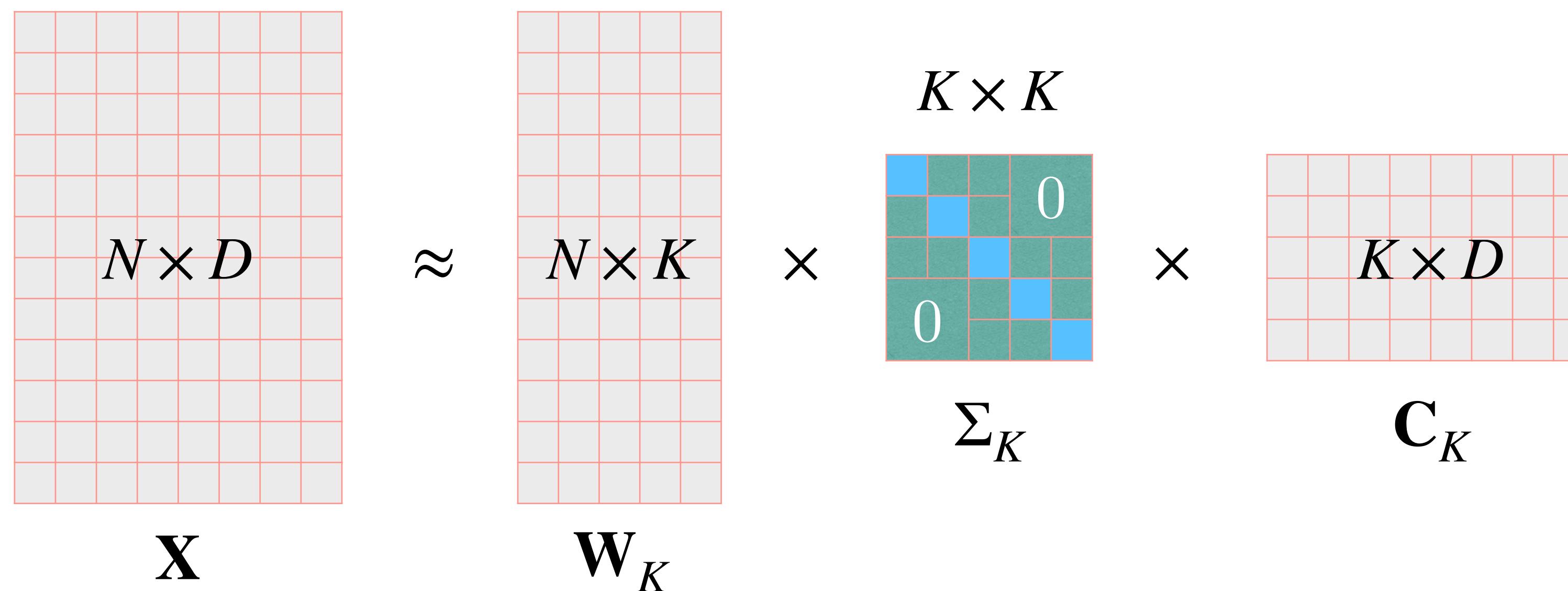
Singular Value Decomposition (SVD; truncated) on the term-document matrix \mathbf{X} representing the frequency of N terms (or n -grams) in D documents

\mathbf{W}_K : each topic's (K topics) distribution over N terms

Σ_K : topic importance

\mathbf{C}_K : each document's (D documents) distribution over K topics

Latent Semantic Analysis (or Indexing) – LSA



Disadvantages

- SVD has a significant computational cost $\approx \mathcal{O}(NDK^2)$
- No intuition about the origin of the topics (*brute force method*)

Probabilistic LSA – pLSA

d_1

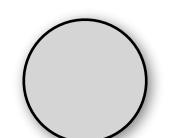
d_2

:

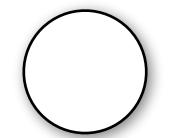
d_D

For all j documents (1 to D):

- Select a document d_j with probability $p(d_j)$
- Choose a mixture of K topics Θ_j for document d_j
- For each word position i (1 to N) in the document d_j :
 - Choose a topic z_k with probability $p(z_k | d_j)$
 - Choose a term/word w_i with probability $p(w_i | z_k)$



observed



latent/hidden

Probabilistic LSA – pLSA

d_1

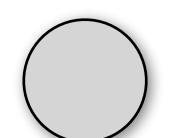
d_2

:

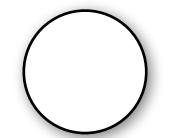
d_D

For all j documents (1 to D):

- Select a document d_j with probability $p(d_j)$
- Choose a mixture of K topics Θ_j for document d_j
- For each word position i (1 to N) in the document d_j :
 - Choose a topic z_k with probability $p(z_k | d_j)$
 - Choose a term/word w_i with probability $p(w_i | z_k)$

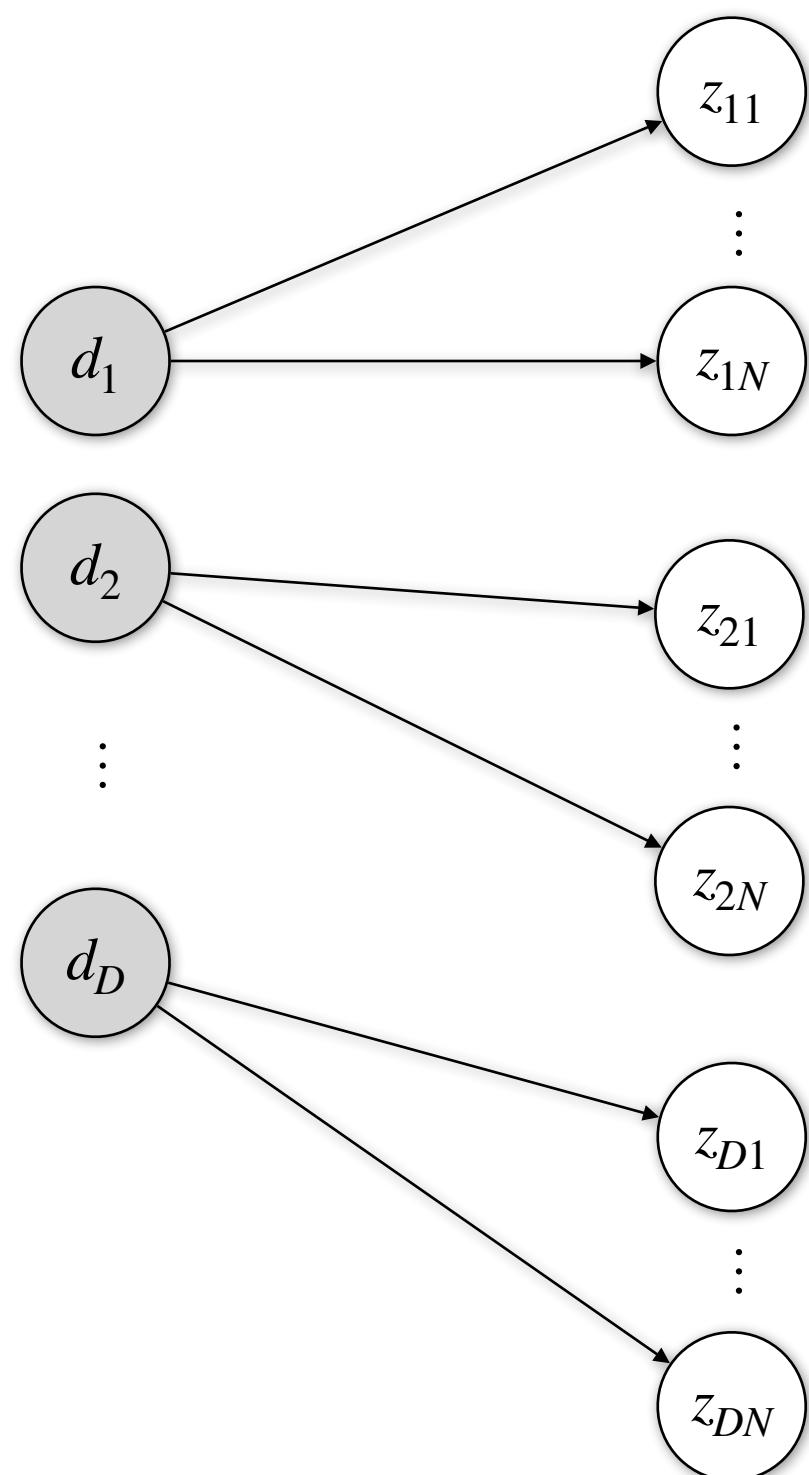


observed



latent/hidden

Probabilistic LSA – pLSA



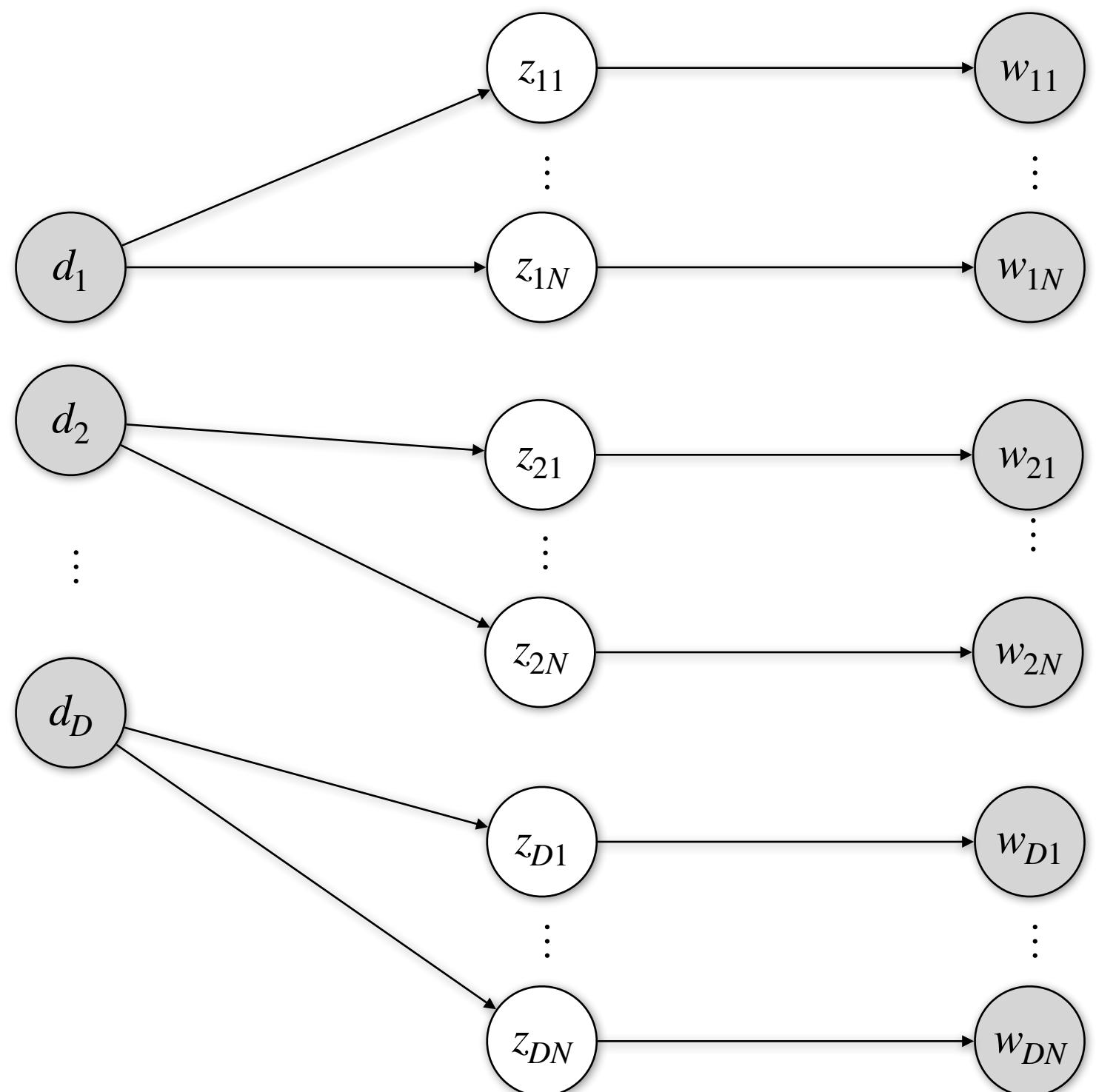
For all j documents (1 to D):

- Select a document d_j with probability $p(d_j)$
- Choose a mixture of K topics Θ_j for document d_j
- For each word position i (1 to N) in the document d_j :
 - Choose a topic z_k with probability $p(z_k | d_j)$
 - Choose a term/word w_i with probability $p(w_i | z_k)$

observed

latent/hidden

Probabilistic LSA – pLSA



For all j documents (1 to D):

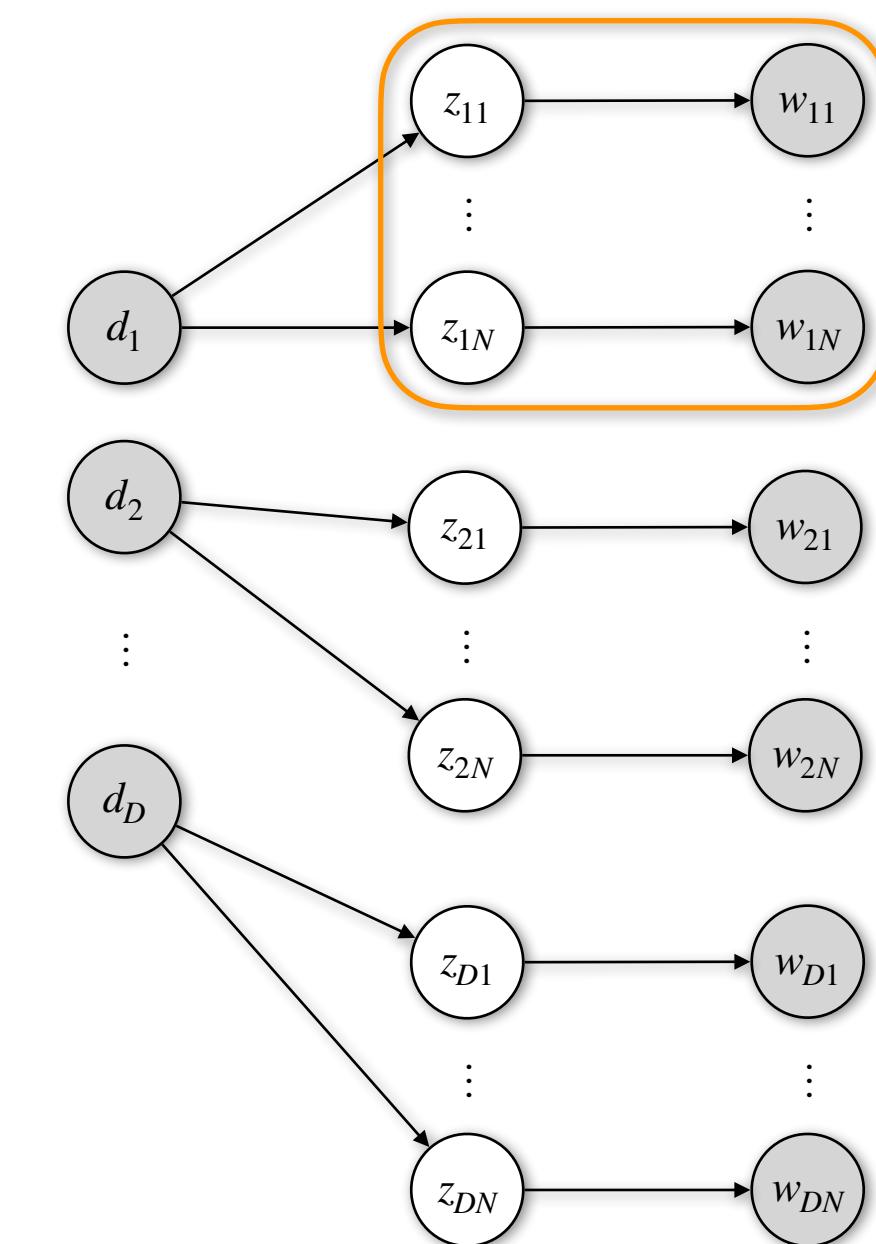
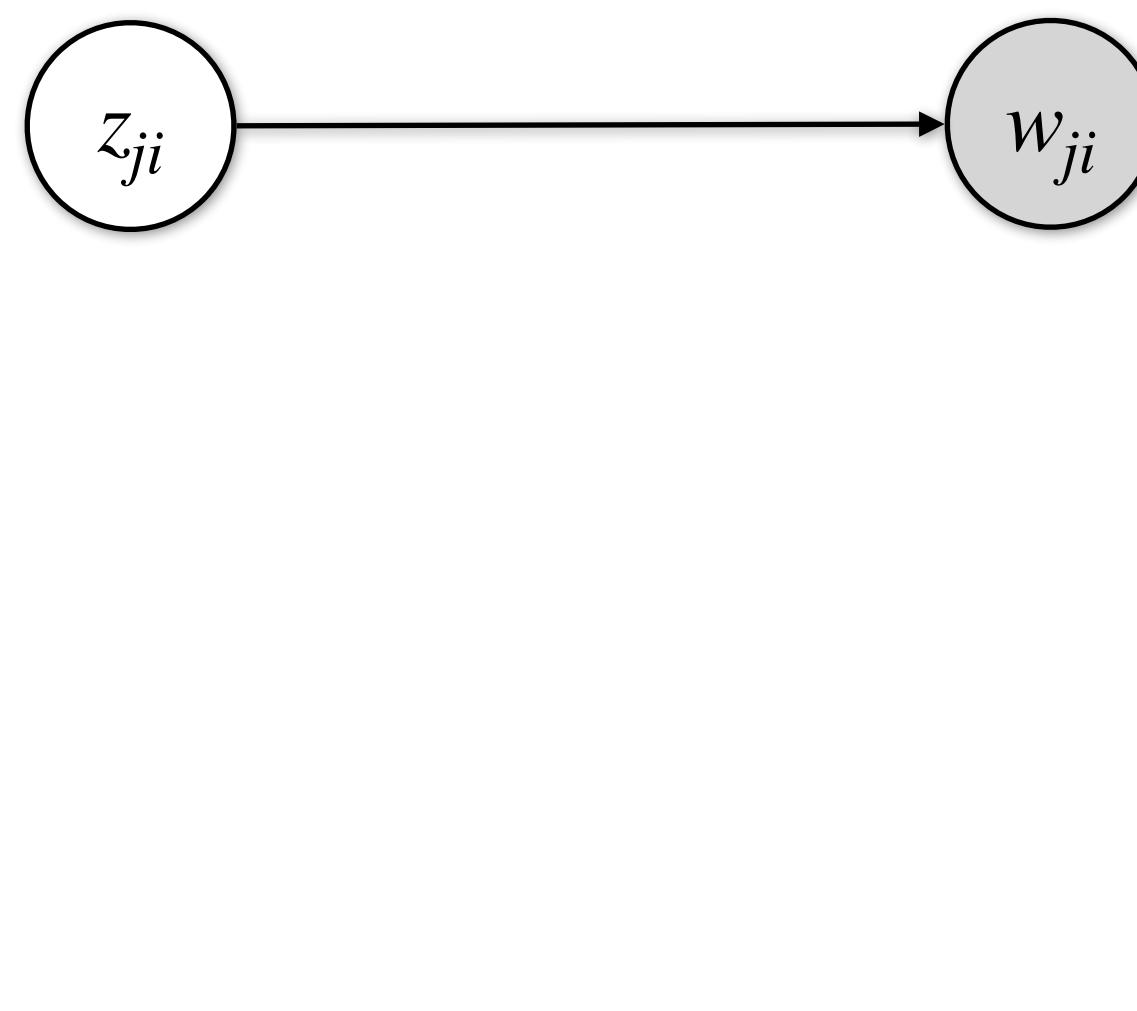
- Select a document d_j with probability $p(d_j)$
- Choose a mixture of K topics Θ_j for document d_j
- For each word position i (1 to N) in the document d_j :
 - Choose a topic z_k with probability $p(z_k | d_j)$
 - Choose a term/word w_i with probability $p(w_i | z_k)$

- observed
- latent/hidden

Probabilistic LSA – pLSA

Plate notation

observed
latent/hidden



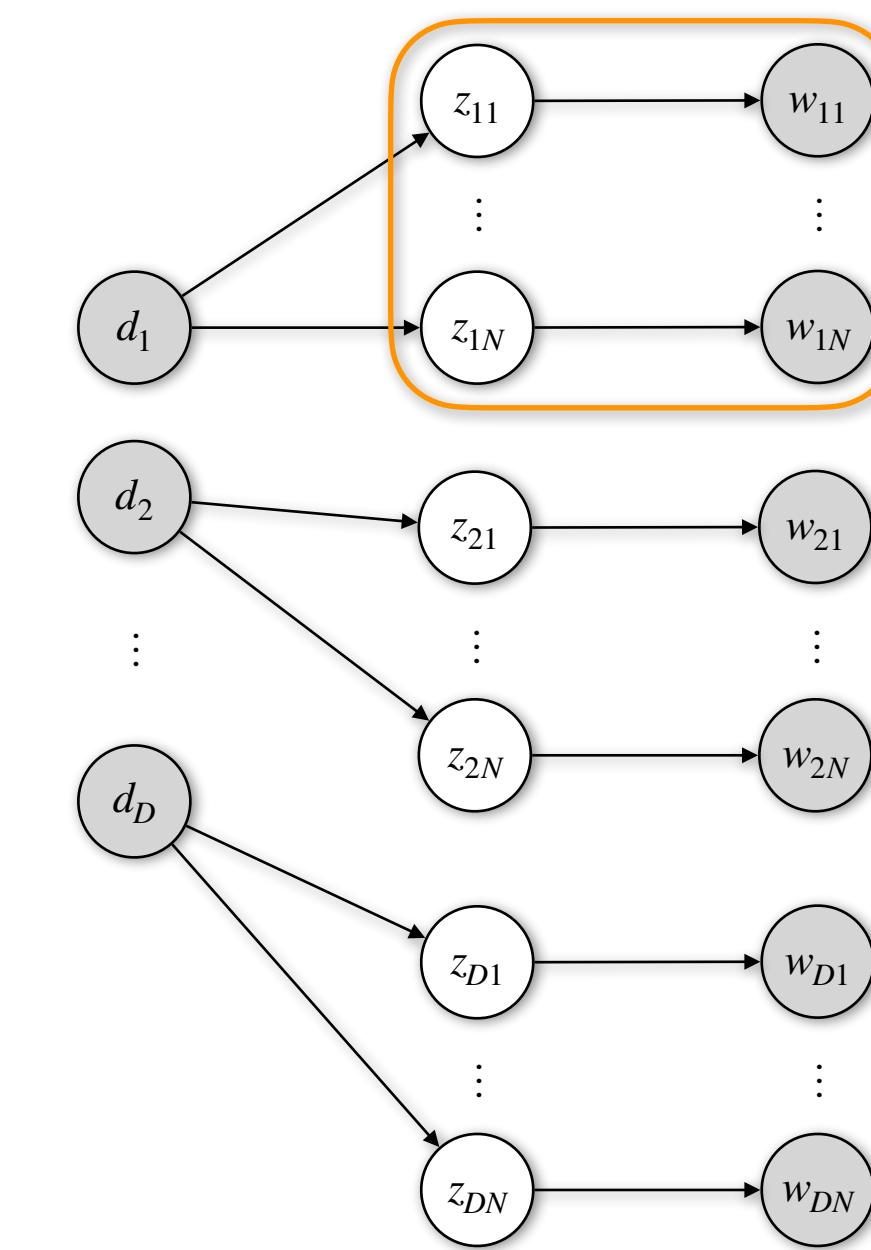
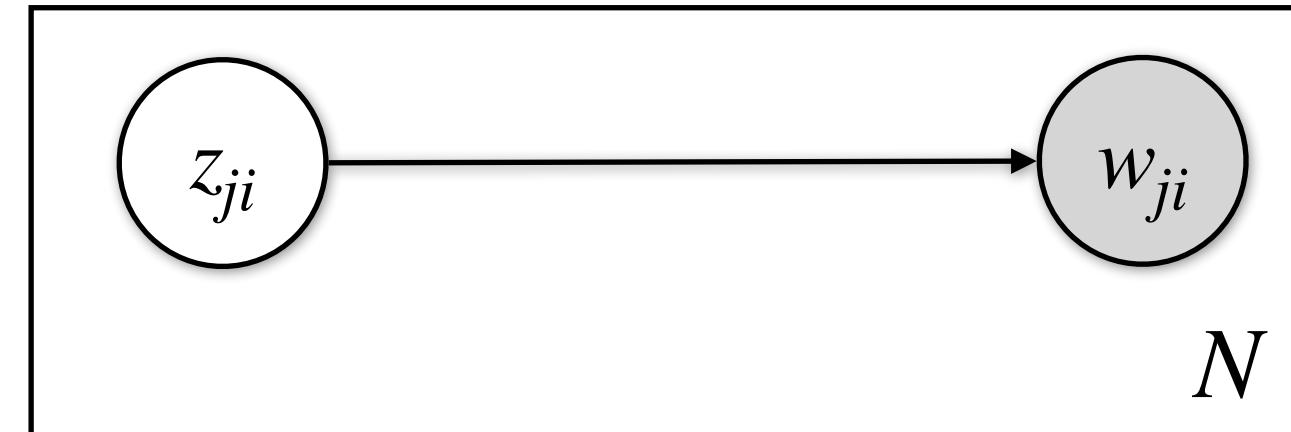
For all j documents (1 to D):

- Select a document d_j with probability $p(d_j)$
- Choose a mixture of K topics θ_j for document d_j
- For each word position i (1 to N) in the document d_j :
 - Choose a topic z_k with probability $p(z_k | d_j)$
 - Choose a term/word w_i with probability $p(w_i | z_k)$

Probabilistic LSA – pLSA

 observed
 latent/hidden

Plate notation



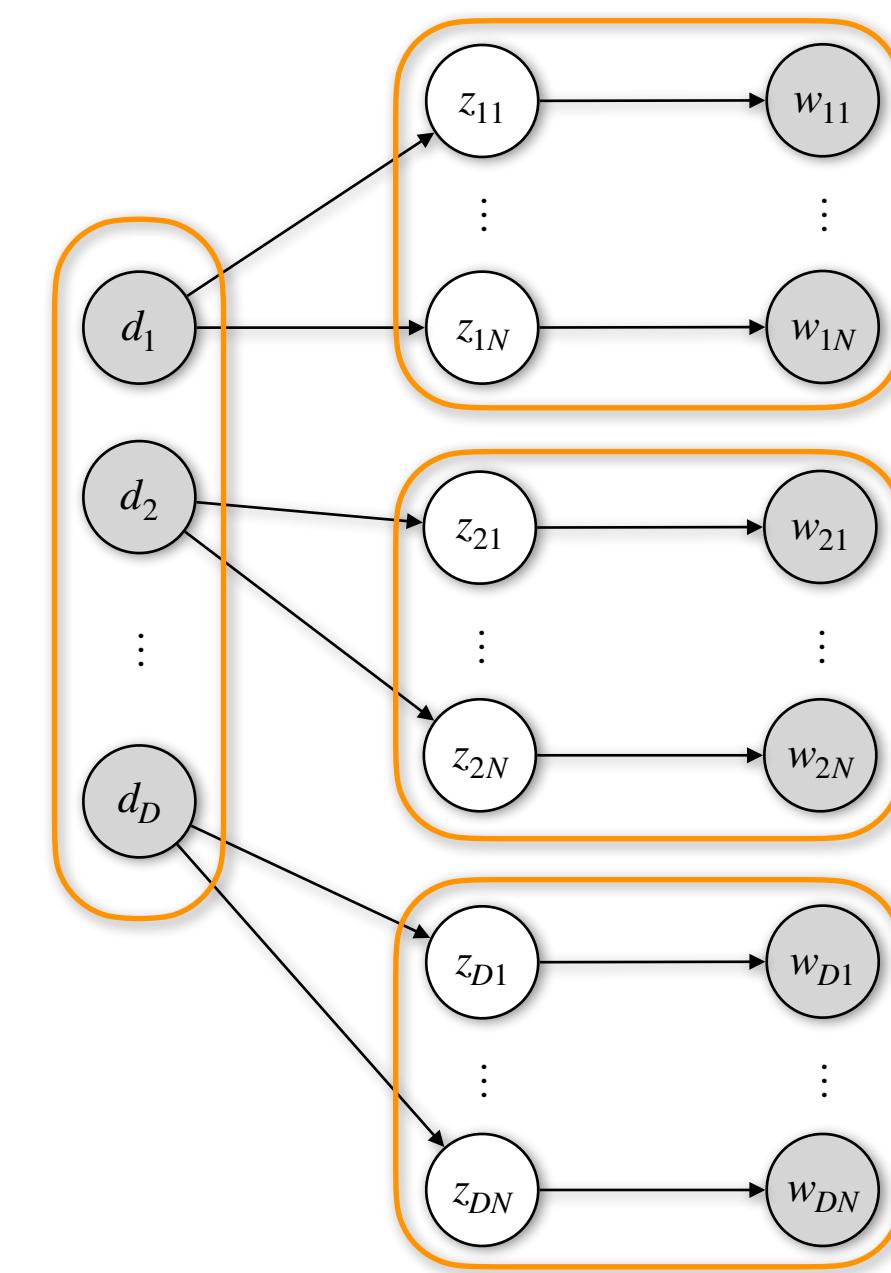
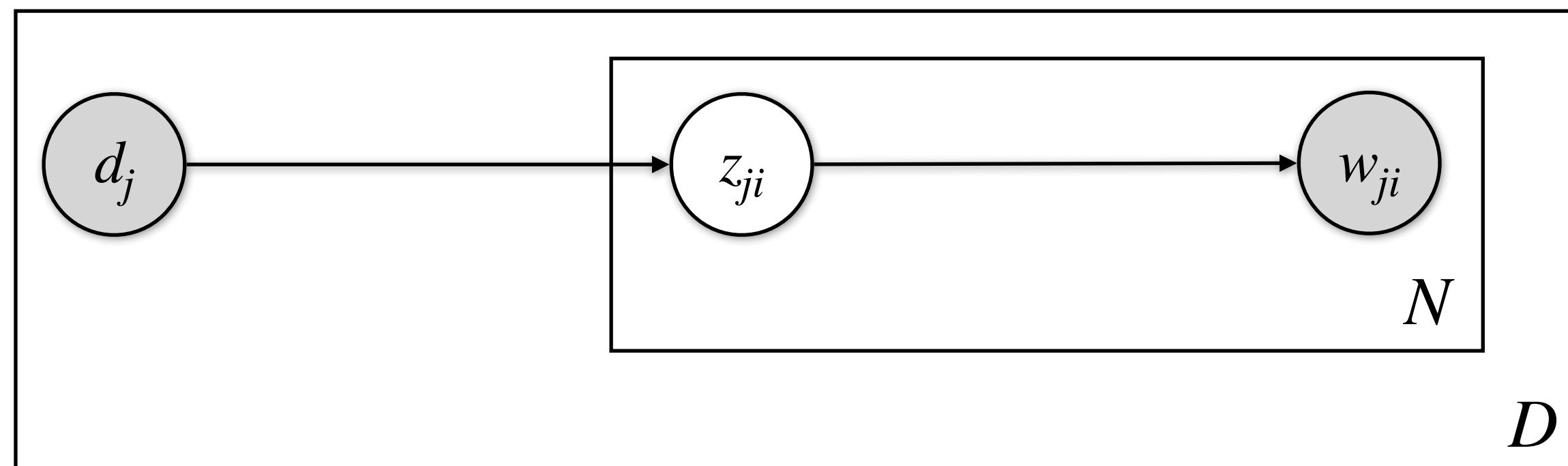
For all j documents (1 to D):

- Select a document d_j with probability $p(d_j)$
- Choose a mixture of K topics θ_j for document d_j
- For each word position i (1 to N) in the document d_j :
 - Choose a topic z_k with probability $p(z_k | d_j)$
 - Choose a term/word w_i with probability $p(w_i | z_k)$

Probabilistic LSA – pLSA

Plate notation

observed
latent/hidden

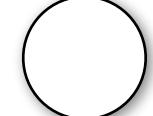


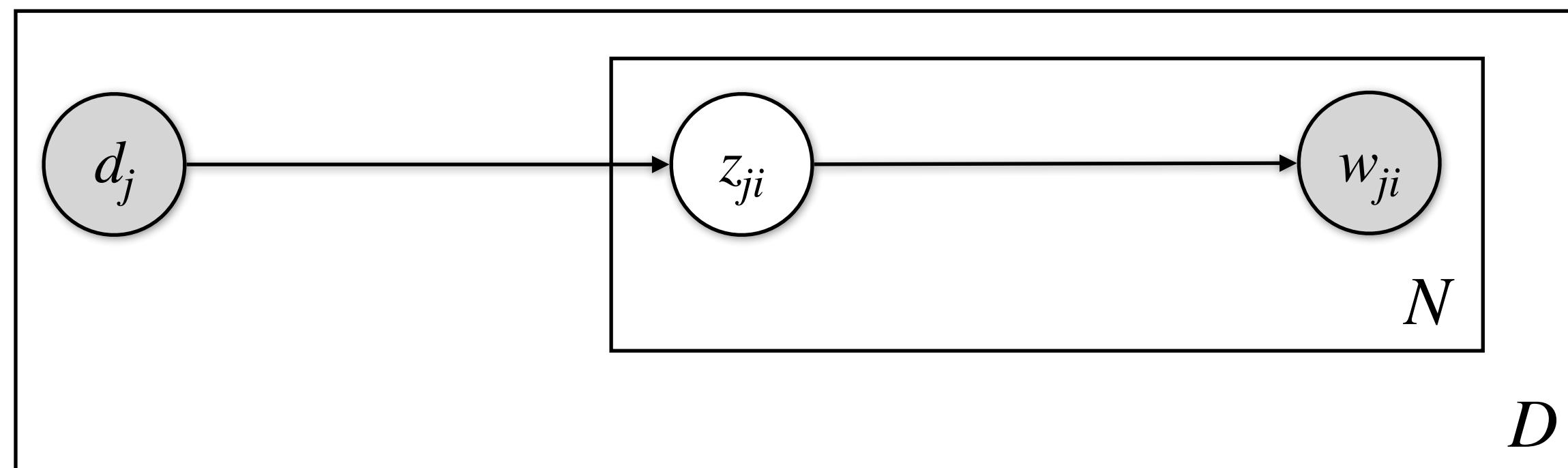
For all j documents (1 to D):

- Select a document d_j with probability $p(d_j)$
- Choose a mixture of K topics θ_j for document d_j
- For each word position i (1 to N) in the document d_j :
 - Choose a topic z_k with probability $p(z_k | d_j)$
 - Choose a term/word w_i with probability $p(w_i | z_k)$

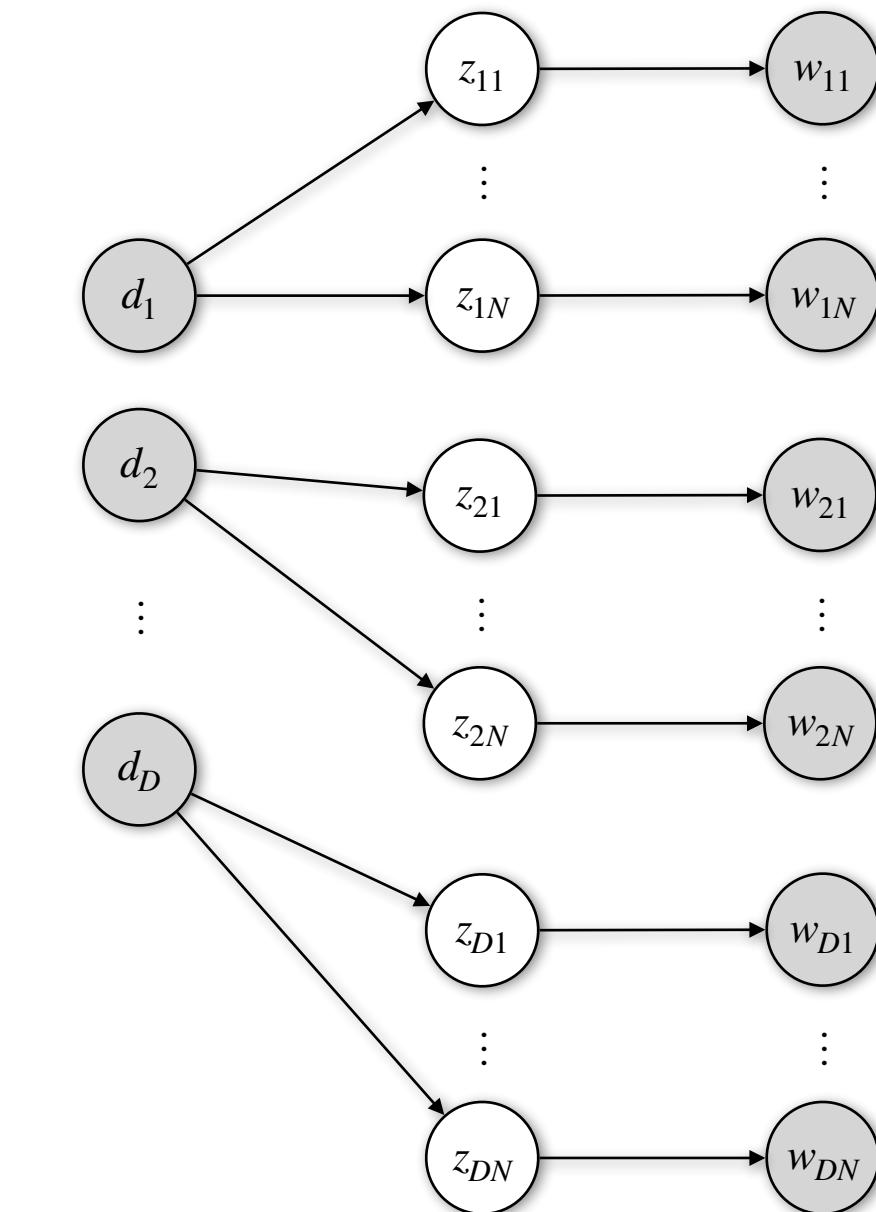
Probabilistic LSA – pLSA

Plate notation

 observed
 latent/hidden



$$p(\mathbf{d}, \mathbf{W}) = \prod_{j=1}^D p(d_j) \prod_{i=1}^N \sum_{k=1}^K p(z_{ji} = k | d_j) p(w_{ji} | z_{ji} = k)$$

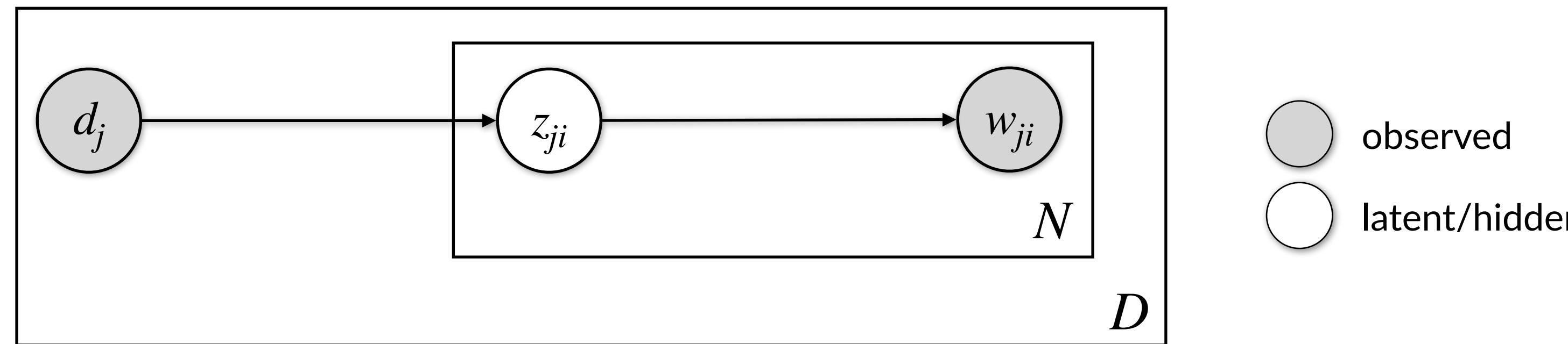


For all j documents (1 to D):

- Select a document d_j with probability $p(d_j)$
- Choose a mixture of K topics θ_j for document d_j
- For each word position i (1 to N) in the document d_j :
 - Choose a topic z_k with probability $p(z_k | d_j)$
 - Choose a term/word w_i with probability $p(w_i | z_k)$

Probabilistic LSA – pLSA

Plate notation



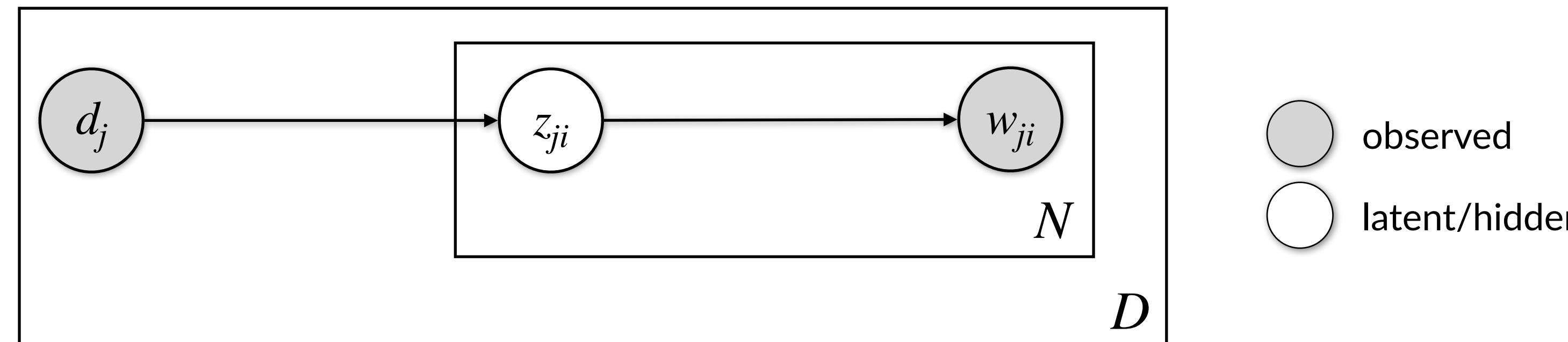
Assumptions: In a document d_j , a word w_{ji} is generated from a single topic z_{ji} from the K assumed ones, and given that topic, the word is independent of all of the other words in that document.

$$p(d_j, w_i) = p(d_j) p(w_i | d_j) = p(d_j) \sum_{k=1}^K p(z = k | d_j) p(w_i | z = k)$$

Joint probability distribution for d_j and w_i

Probabilistic LSA – pLSA

Plate notation



Assumptions: In a document d_j , a word w_{ji} is generated from a single topic z_{ji} from the K assumed ones, and given that topic, the word is independent of all of the other words in that document.

$$p(d_j, w_i) = p(d_j) p(w_i | d_j) = p(d_j) \sum_{k=1}^K p(z = k | d_j) p(w_i | z = k)$$

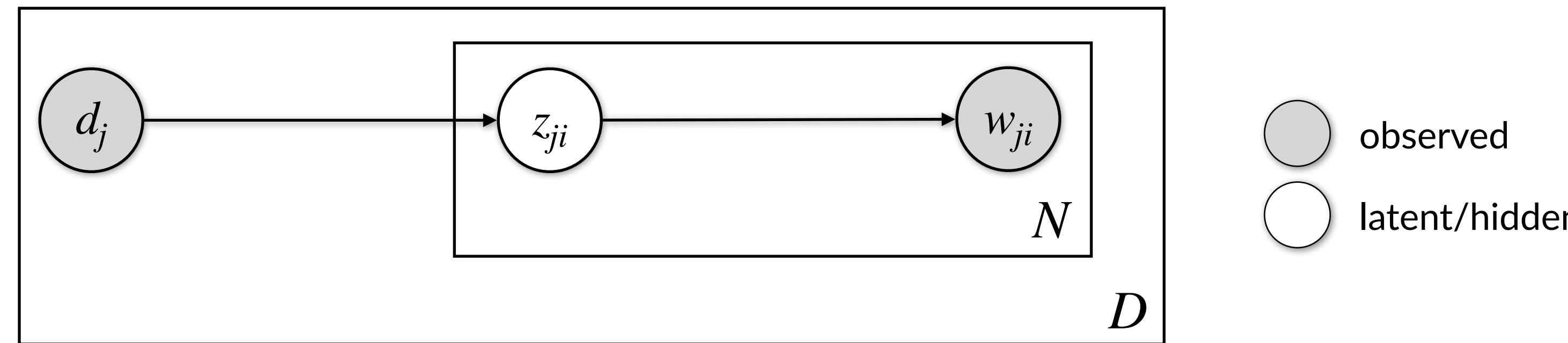
Joint probability distribution for d_j and w_i

$$p(d_j, \mathbf{w}_j) = p(d_j) \prod_{i=1}^N \sum_{k=1}^K p(z_i = k | d_j) p(w_{ji} | z_i = k)$$

Joint probability distribution for d_j and \mathbf{w}_j

Probabilistic LSA – pLSA

Plate notation



Assumptions: In a document d_j , a word w_{ji} is generated from a single topic z_{ji} from the K assumed ones, and given that topic, the word is independent of all of the other words in that document.

$$p(d_j, w_i) = p(d_j) p(w_i | d_j) = p(d_j) \sum_{k=1}^K p(z = k | d_j) p(w_i | z = k)$$

Joint probability distribution for d_j and w_i

$$p(d_j, \mathbf{w}_j) = p(d_j) \prod_{i=1}^N \sum_{k=1}^K p(z_i = k | d_j) p(w_{ji} | z_i = k)$$

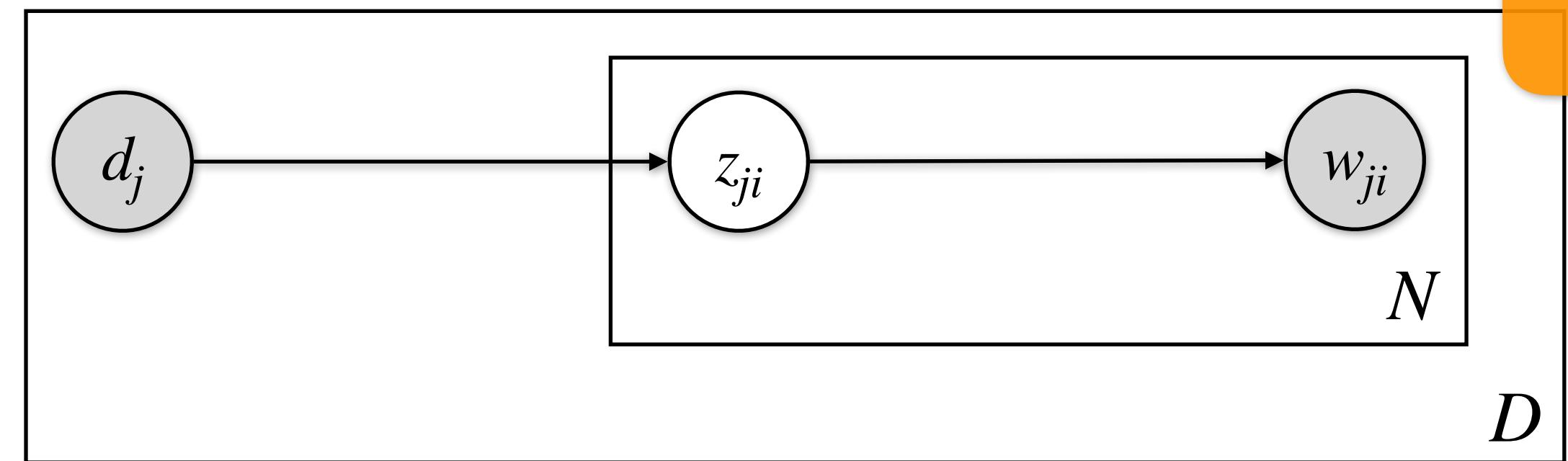
Joint probability distribution for d_j and \mathbf{w}_j

$$p(\mathbf{d}, \mathbf{W}) = \prod_{j=1}^D p(d_j) \prod_{i=1}^N \sum_{k=1}^K p(z_{ji} = k | d_j) p(w_{ji} | z_{ji} = k)$$

Joint probability distribution

Probabilistic LSA – pLSA

Plate notation



Find a **minor mistake** in this slide and previous ones

- Shaded circle: observed
- White circle: latent/hidden

Assumptions: In a document d_j , a word w_{ji} is generated from a single topic z_{ji} from the K assumed ones, and given that topic, the word is independent of all of the other words in that document.

$$p(d_j, w_i) = p(d_j) p(w_i | d_j) = p(d_j) \sum_{k=1}^K p(z = k | d_j) p(w_i | z = k)$$

Joint probability distribution for d_j and w_i

$$p(d_j, \mathbf{w}_j) = p(d_j) \prod_{i=1}^N \sum_{k=1}^K p(z_i = k | d_j) p(w_{ji} | z_i = k)$$

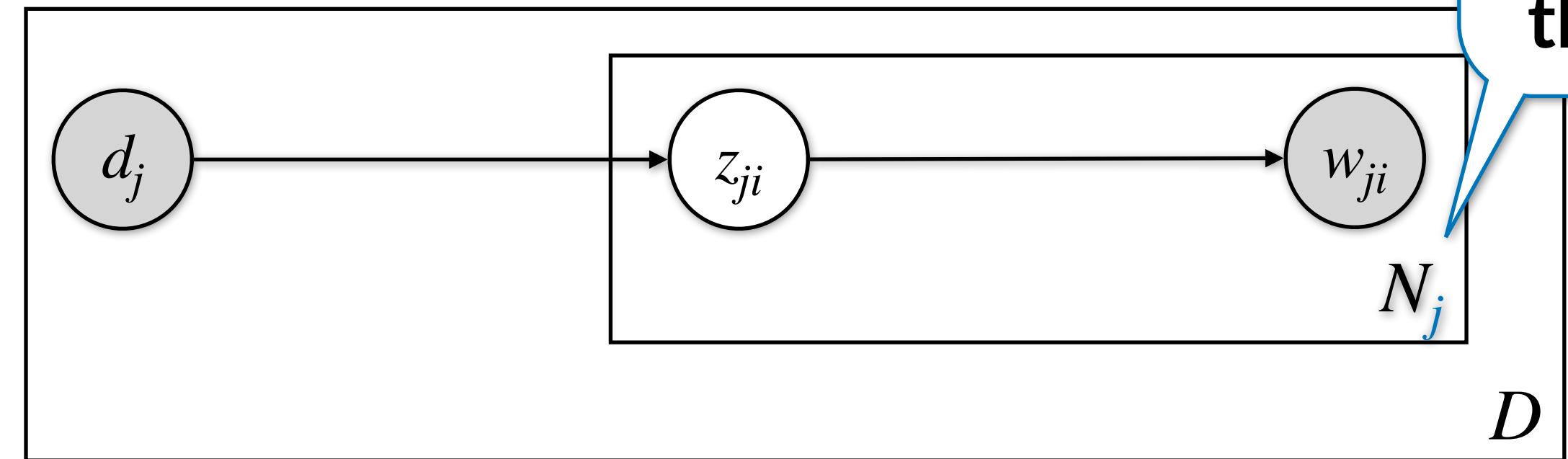
Joint probability distribution for d_j and \mathbf{w}_j

$$p(\mathbf{d}, \mathbf{W}) = \prod_{j=1}^D p(d_j) \prod_{i=1}^N \sum_{k=1}^K p(z_{ji} = k | d_j) p(w_{ji} | z_{ji} = k)$$

Joint probability distribution

Probabilistic LSA – pLSA

Plate notation



The number of words may not be the same for all the documents!

- observed
- latent/hidden

Assumptions: In a document d_j , a word w_{ji} is generated from a single topic z_{ji} from the K assumed ones, and given that topic, the word is independent of all of the other words in that document.

$$p(d_j, w_i) = p(d_j) p(w_i | d_j) = p(d_j) \sum_{k=1}^K p(z = k | d_j) p(w_i | z = k)$$

Joint probability distribution for d_j and w_i

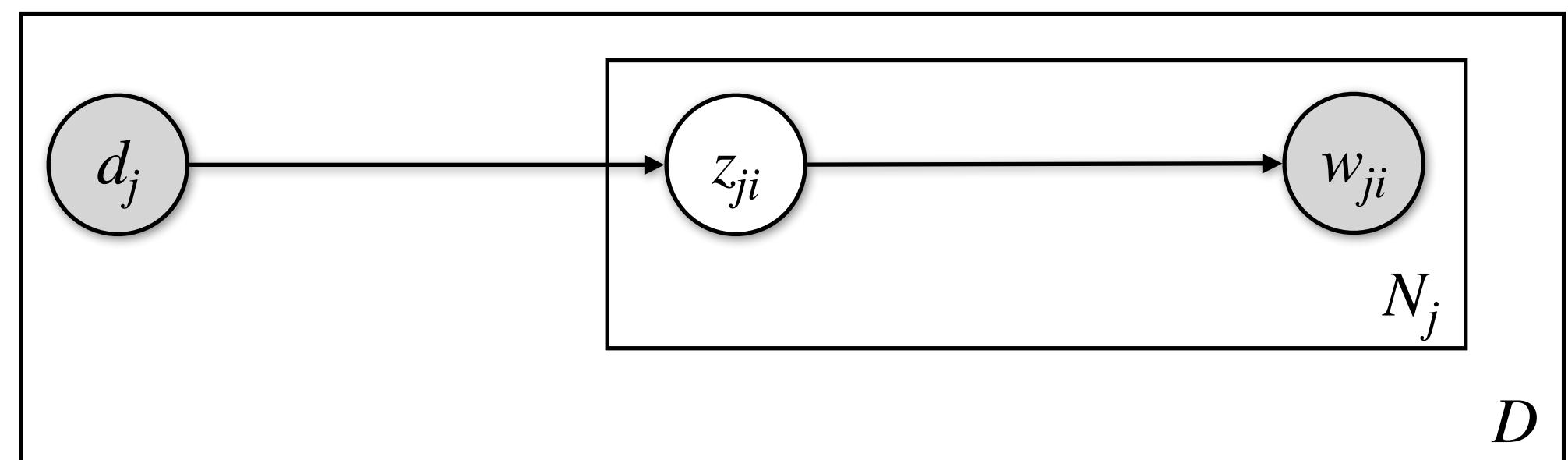
$$p(d_j, \mathbf{w}_j) = p(d_j) \prod_{i=1}^{N_j} \sum_{k=1}^K p(z_i = k | d_j) p(w_{ji} | z_i = k)$$

Joint probability distribution for d_j and \mathbf{w}_j

$$p(\mathbf{d}, \mathbf{W}) = \prod_{j=1}^D p(d_j) \prod_{i=1}^{N_j} \sum_{k=1}^K p(z_{ji} = k | d_j) p(w_{ji} | z_{ji} = k)$$

Joint probability distribution

$$p(\mathbf{d}, \mathbf{W}) = \prod_{j=1}^D p(d_j) \prod_{i=1}^{N_j} \sum_{k=1}^K p(z_{ji} = k | d_j) p(w_{ji} | z_{ji} = k)$$

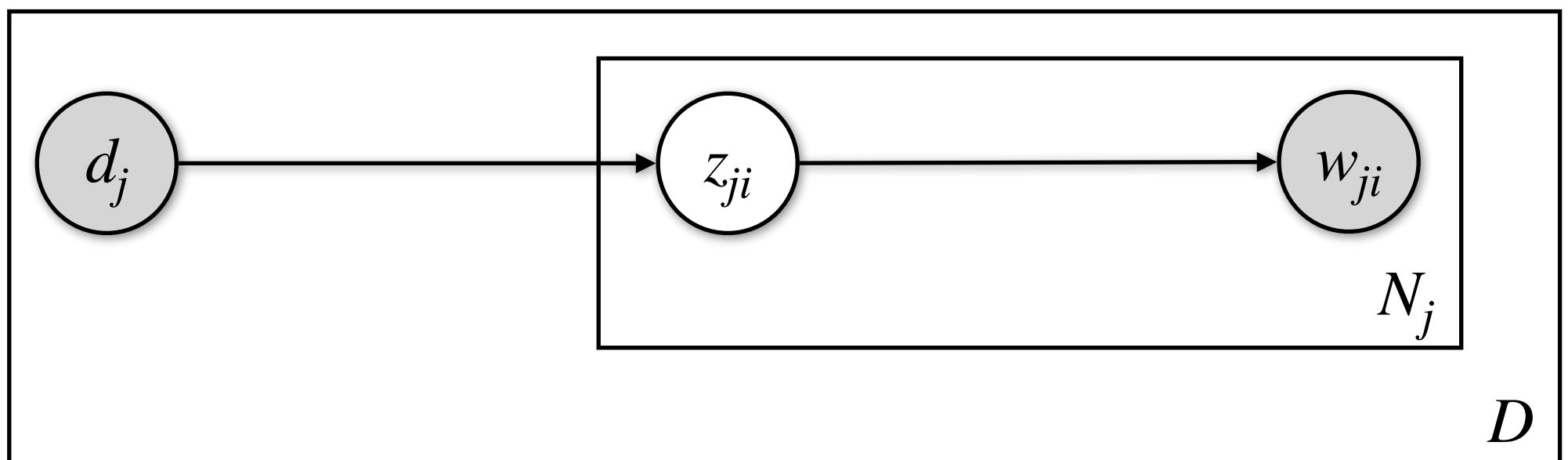


Expectation Maximisation (EM)

- ▶ **E-step:** Compute expected values of the variables, given the current parametrisation of the model. In the very beginning, start with a *random* or *uniform* parametrisation
- ▶ **M-step:** Then, using the above values, update the model parameters
- ▶ Go back to the E-step; **repeat until convergence**

pLSA – Inference

$$p(\mathbf{d}, \mathbf{W}) = \prod_{j=1}^D p(d_j) \prod_{i=1}^{N_j} \sum_{k=1}^K p(z_{ji} = k | d_j) p(w_{ji} | z_{ji} = k)$$

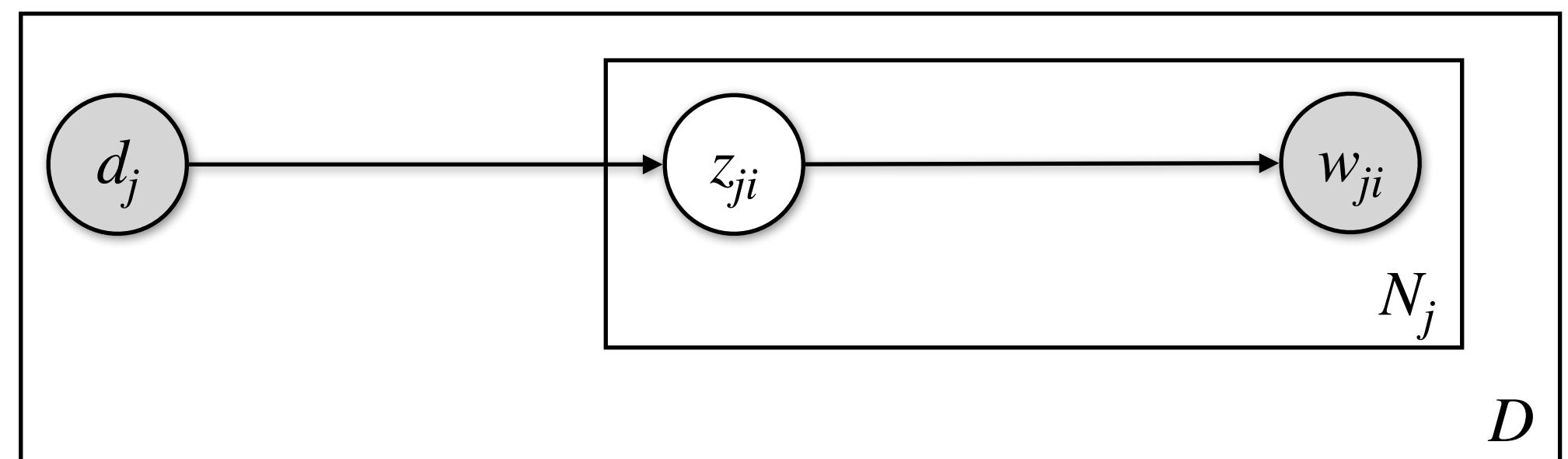


- ▶ Initialise $p(z_k | d_j)$ and $p(w_i | z_k)$ to positive quantities
- ▶ **E-step:** Estimate the probability of each topic given the words in each document

$$p(z_k | d_j, w_i) = \frac{p(z_k | d_j) p(w_i | z_k)}{\sum_{k'=1}^K p(z_{k'} | d_j) p(w_i | z_{k'})}$$

pLSA – Inference

$$p(\mathbf{d}, \mathbf{W}) = \prod_{j=1}^D p(d_j) \prod_{i=1}^{N_j} \sum_{k=1}^K p(z_{ji} = k | d_j) p(w_{ji} | z_{ji} = k)$$



- ▶ Initialise $p(z_k | d_j)$ and $p(w_i | z_k)$ to positive quantities

- ▶ **E-step:** Estimate the probability of each topic given the words in each document

- ▶ **M-step:** Re-estimate $p(z_k | d_j), p(w_i | z_k)$ given the revised $p(z_k | d_j, w_i)$

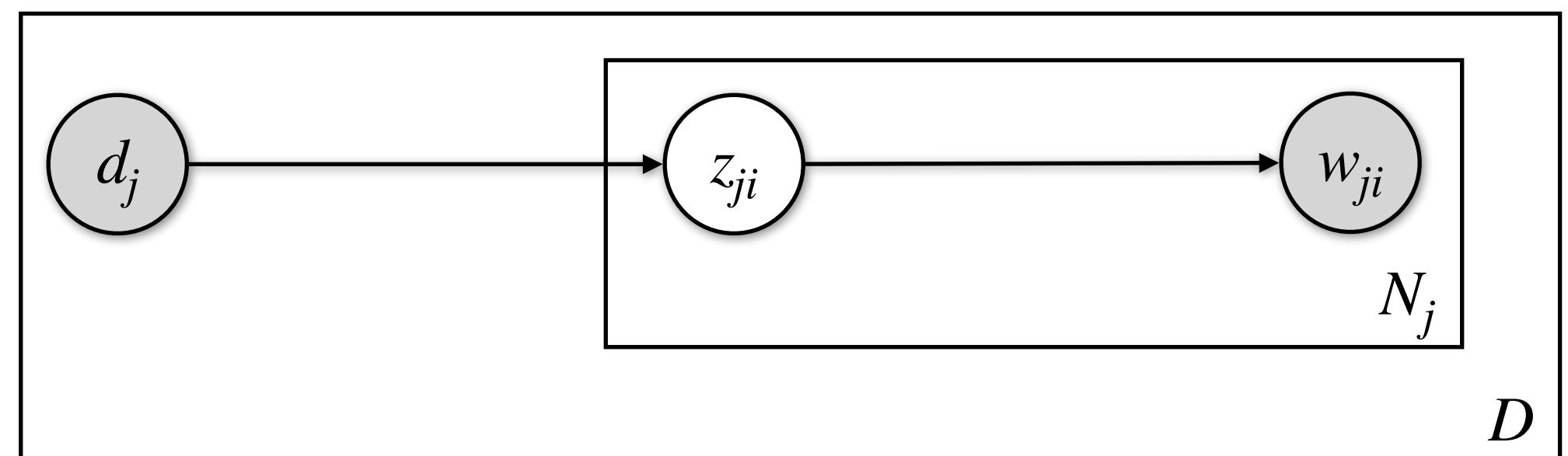
$$p(w_i | z_k) = \frac{\sum_{j=1}^D n(d_j, w_i) p(z_k | d_j, w_i)}{\sum_{j=1}^D \sum_{i'=1}^{N_j} n(d_j, w_{i'}) p(z_k | d_j, w_{i'})}$$

$$p(z_k | d_j, w_i) = \frac{p(z_k | d_j) p(w_i | z_k)}{\sum_{k'=1}^K p(z_{k'} | d_j) p(w_i | z_{k'})}$$

$$p(z_k | d_j) = \frac{\sum_{i=1}^{N_j} n(d_j, w_i) p(z_k | d_j, w_i)}{\sum_{i=1}^{N_j} \sum_{k'=1}^K n(d_j, w_i) p(z_{k'} | d_j, w_i)}$$

pLSA – Inference

$$p(\mathbf{d}, \mathbf{W}) = \prod_{j=1}^D p(d_j) \prod_{i=1}^{N_j} \sum_{k=1}^K p(z_{ji} = k | d_j) p(w_{ji} | z_{ji} = k)$$



- ▶ Initialise $p(z_k | d_j)$ and $p(w_i | z_k)$ to positive quantities
- ▶ **E-step:** Estimate the probability of each topic given the words in each document
- ▶ **M-step:** Re-estimate $p(z_k | d_j), p(w_i | z_k)$ given the revised $p(z_k | d_j, w_i)$

$$p(w_i | z_k) = \frac{\sum_{j=1}^D n(d_j, w_i) p(z_k | d_j, w_i)}{\sum_{j=1}^D \sum_{i'=1}^{N_j} n(d_j, w_{i'}) p(z_k | d_j, w_{i'})}$$

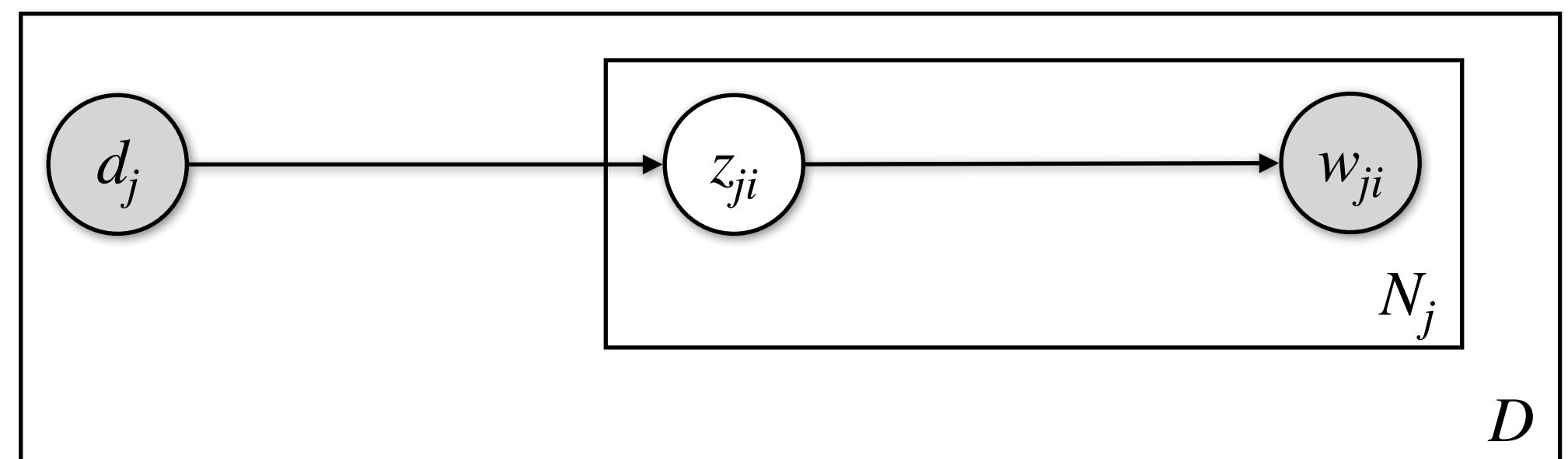
$$p(z_k | d_j, w_i) = \frac{p(z_k | d_j) p(w_i | z_k)}{\sum_{k'=1}^K p(z_{k'} | d_j) p(w_i | z_{k'})}$$

Weighted sums, e.g. $n(d_j, w_i)$ is the number of times word i appears in document j .

$$p(z_k | d_j) = \frac{\sum_{i=1}^{N_j} n(d_j, w_i) p(z_k | d_j, w_i)}{\sum_{i=1}^{N_j} \sum_{k'=1}^K n(d_j, w_i) p(z_{k'} | d_j, w_i)}$$

pLSA – Inference

$$p(\mathbf{d}, \mathbf{W}) = \prod_{j=1}^D p(d_j) \prod_{i=1}^{N_j} \sum_{k=1}^K p(z_{ji} = k | d_j) p(w_{ji} | z_{ji} = k)$$



- Initialise $p(z_k | d_j)$ and $p(w_i | z_k)$ to positive quantities

- E-step:** Estimate the probability of each topic given the words in each document

- M-step:** Re-estimate $p(z_k | d_j), p(w_i | z_k)$ given the revised $p(z_k | d_j, w_i)$

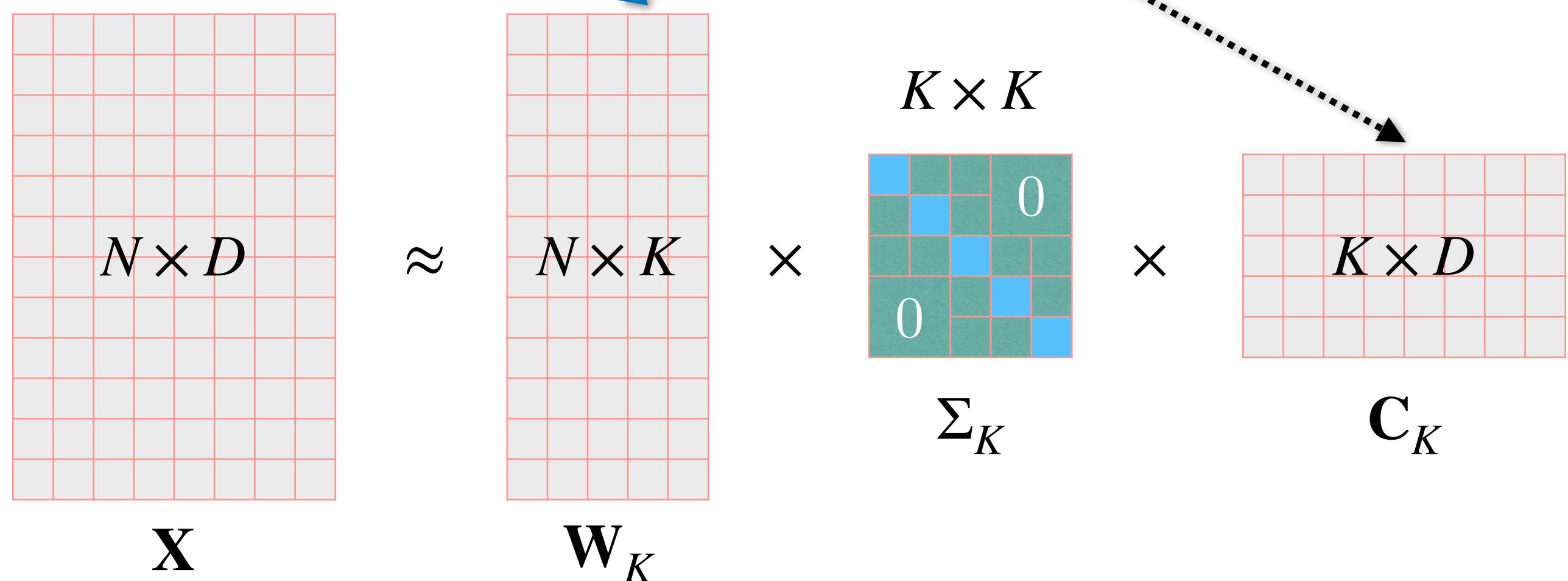
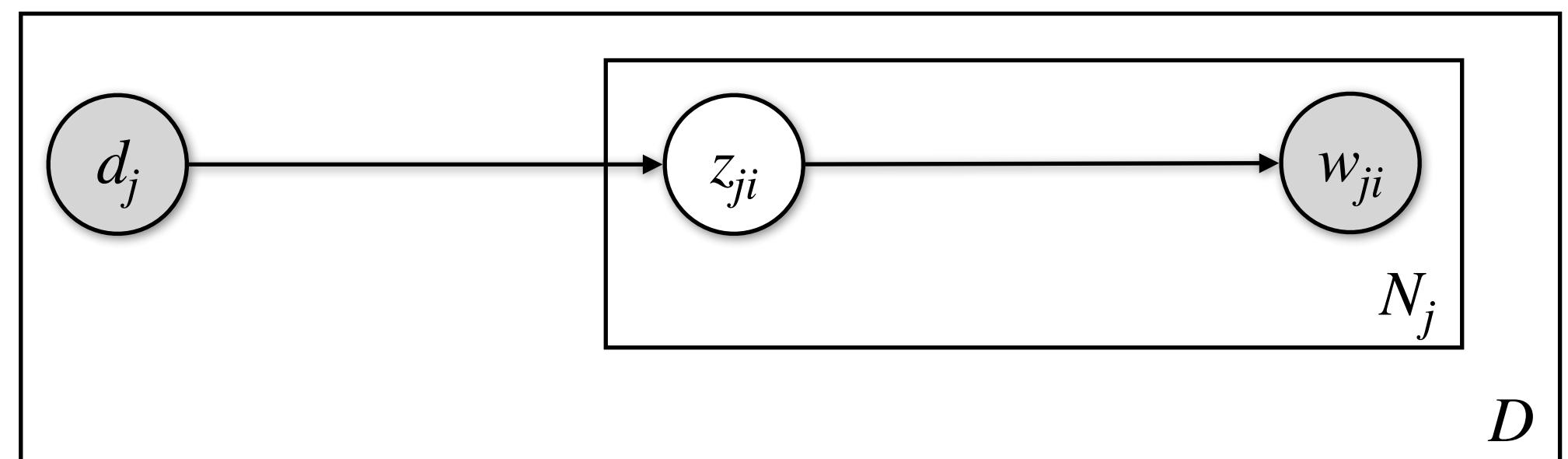
$$p(w_i | z_k) = \frac{\sum_{j=1}^D n(d_j, w_i) p(z_k | d_j, w_i)}{\sum_{j=1}^D \sum_{i'=1}^{N_j} n(d_j, w_{i'}) p(z_k | d_j, w_{i'})}$$

$$p(z_k | d_j, w_i) = \frac{p(z_k | d_j) p(w_i | z_k)}{\sum_{k'=1}^K p(z_{k'} | d_j) p(w_i | z_{k'})}$$

$$p(z_k | d_j) = \frac{\sum_{i=1}^{N_j} n(d_j, w_i) p(z_k | d_j, w_i)}{\sum_{i=1}^{N_j} \sum_{k'=1}^K n(d_j, w_i) p(z_{k'} | d_j, w_i)}$$

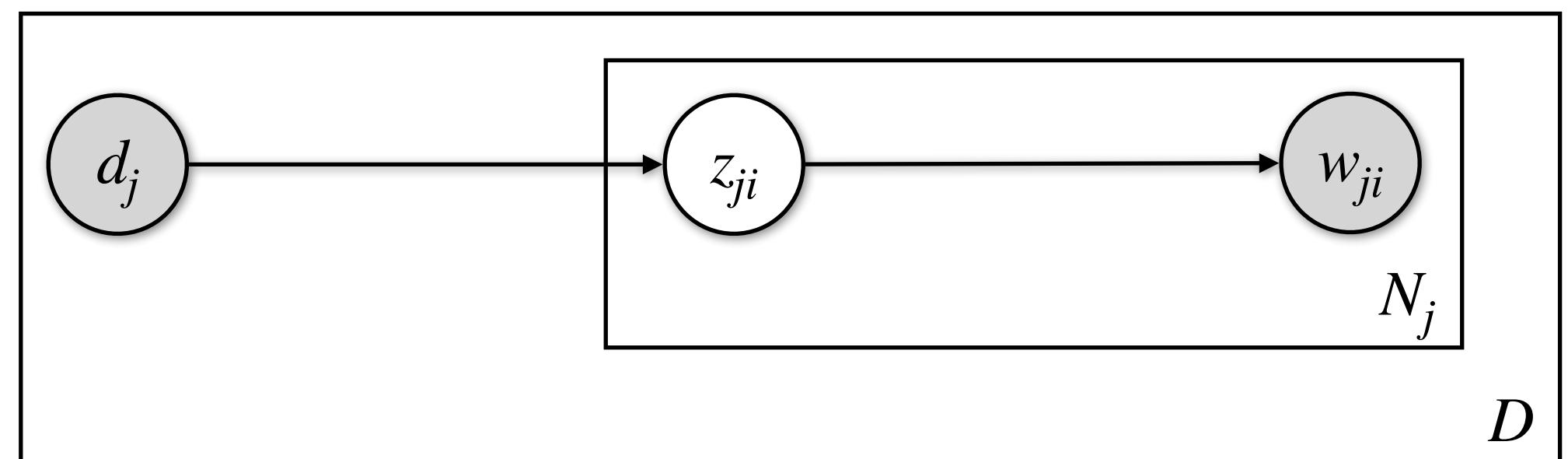
pLSA and LSA

$$p(\mathbf{d}, \mathbf{W}) = \prod_{j=1}^D p(d_j) \prod_{i=1}^{N_j} \sum_{k=1}^K \frac{p(z_{ji} = k | d_j) p(w_{ji} | z_{ji} = k)}{\text{Norm}}$$



pLSA and LSA

$$p(\mathbf{d}, \mathbf{W}) = \prod_{j=1}^D p(d_j) \prod_{i=1}^{N_j} \sum_{k=1}^K \frac{p(z_{ji} = k | d_j) p(w_{ji} | z_{ji} = k)}{\dots}$$



$\mathbf{X} \approx \mathbf{W}_K \times \Sigma_K \times \mathbf{C}_K \times \mathbf{Z}$

Matrix factorization diagram:

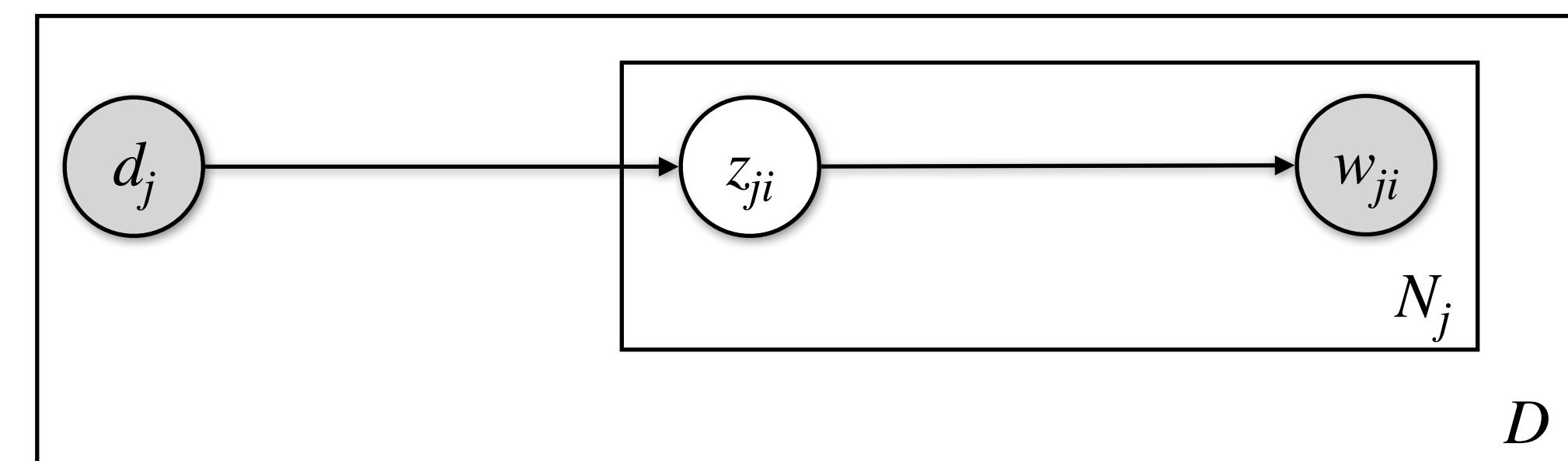
- \mathbf{X} is an $N \times D$ matrix.
- \mathbf{W}_K is an $N \times K$ matrix.
- Σ_K is a $K \times K$ diagonal matrix with zeros on the diagonal.
- \mathbf{C}_K is a $K \times D$ matrix.
- \mathbf{Z} is a latent variable matrix.

Topic distribution equation:

$$p(z_k) = \frac{\sum_{j=1}^D \sum_{i=1}^{N_j} n(d_j, w_i) p(z_k | d_j, w_i)}{\sum_{j=1}^D \sum_{i=1}^{N_j} n(d_j, w_i)}$$

pLSA and LSA

$$p(\mathbf{d}, \mathbf{W}) = \prod_{j=1}^D p(d_j) \prod_{i=1}^{N_j} \sum_{k=1}^K \frac{p(z_{ji} = k | d_j) p(w_{ji} | z_{ji} = k)}{\text{Norm}}$$



$$\mathbf{X} \approx \mathbf{W}_K \times \Sigma_K \times \mathbf{C}_K$$

where $\Sigma_K = \sum_{j=1}^D \sum_{i=1}^{N_j} n(d_j, w_i) p(z_k | d_j, w_i)$

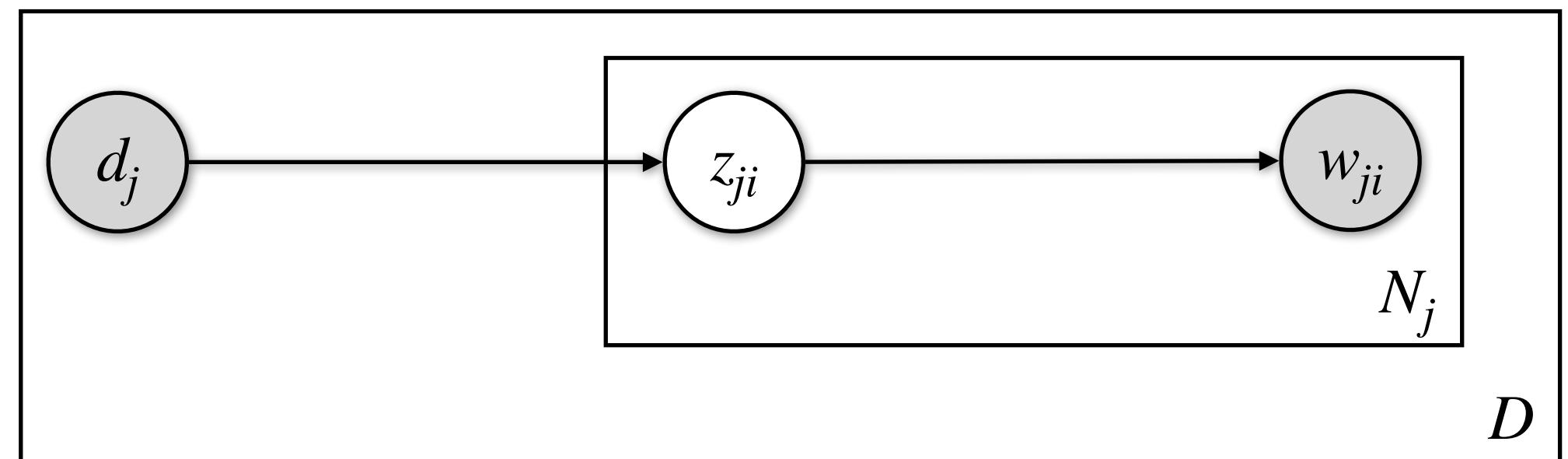
$p(z_k) = \frac{\sum_{j=1}^D \sum_{i=1}^{N_j} n(d_j, w_i) p(z_k | d_j, w_i)}{\sum_{j=1}^D \sum_{i=1}^{N_j} n(d_j, w_i)}$

Main difference

The two techniques have a different objective function – probabilistic vs. deterministic approach

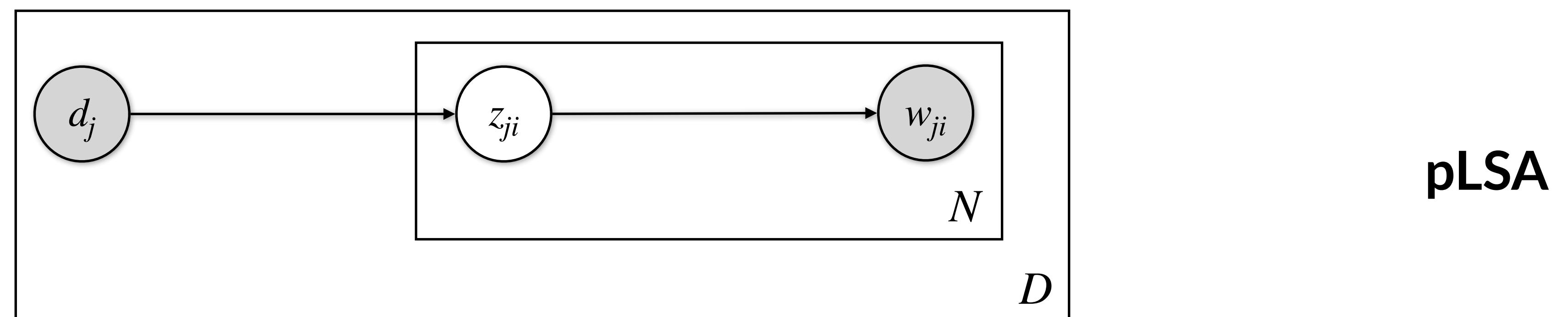
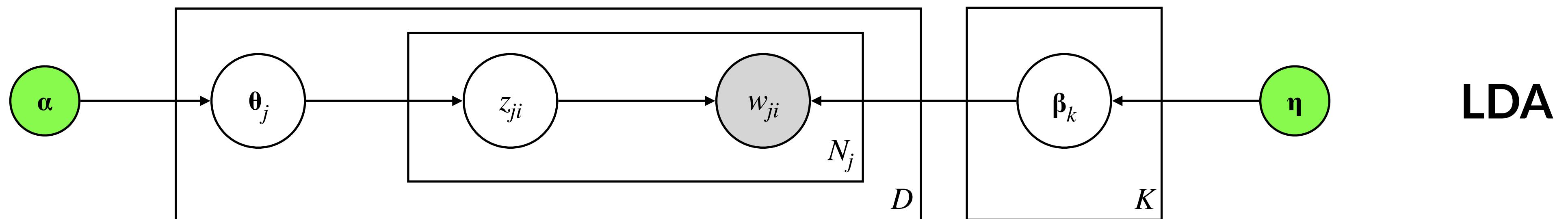
pLSA – Disadvantages

$$p(\mathbf{d}, \mathbf{W}) = \prod_{j=1}^D p(d_j) \prod_{i=1}^{N_j} \sum_{k=1}^K p(z_{ji} = k | d_j) p(w_{ji} | z_{ji} = k)$$

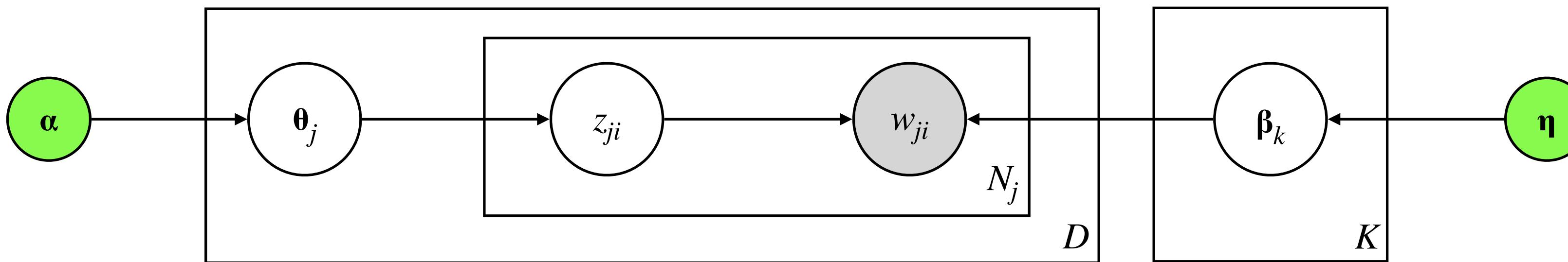


- ▶ The number of parameters that we need to learn during training grows linearly with the number of documents (D), which ultimately leads to overfitting.
- ▶ pLSA learns $p(z_k | d_j)$ only for the documents it sees during the training phase. To deal with a new document, it needs to repeat EM (retrain).

Latent Dirichlet Allocation (LDA)

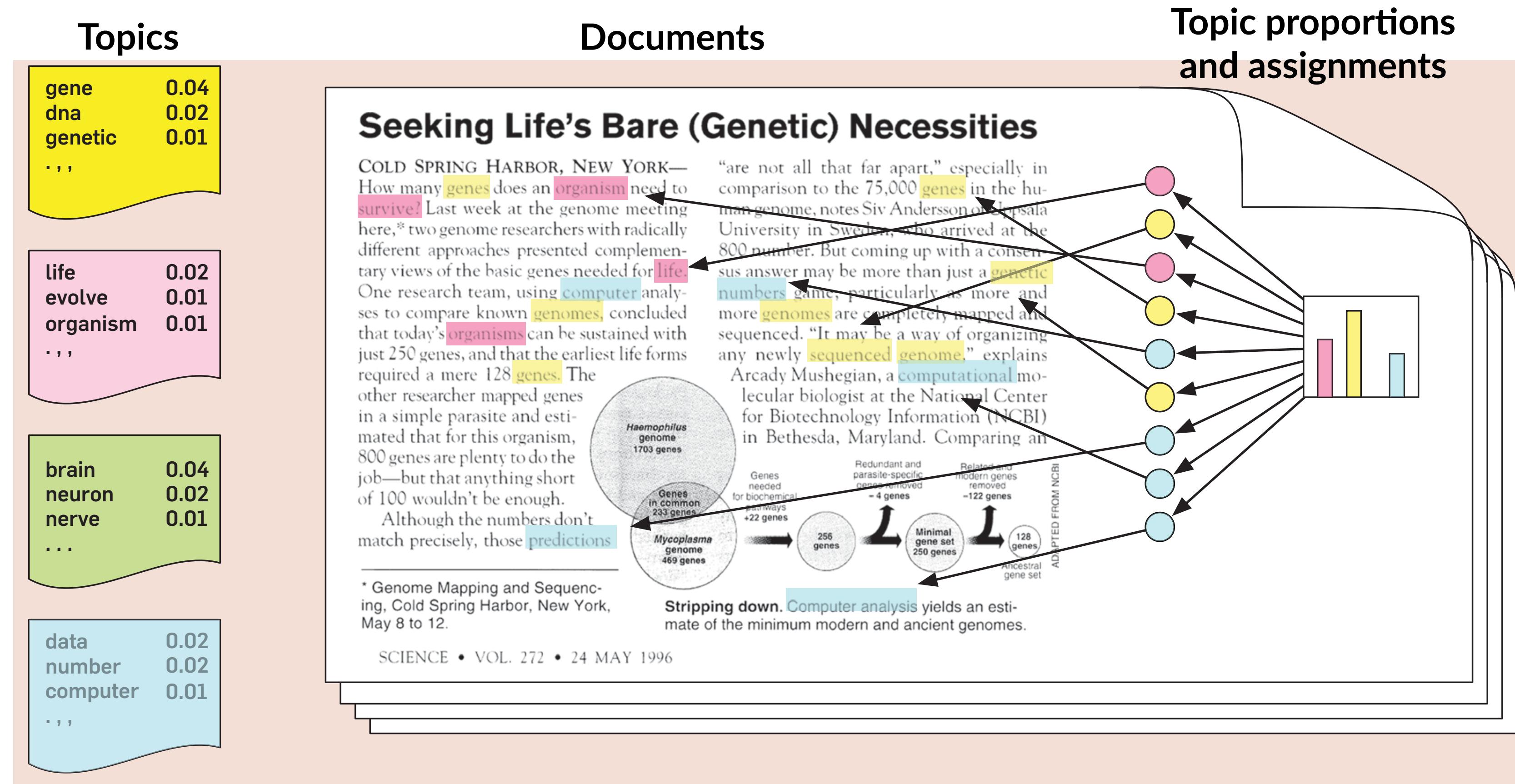


Latent Dirichlet Allocation (LDA)



1. For each of the K topics draw a multinomial distribution β_k from a Dirichlet distribution with parameter η
2. For each of the D documents draw a multinomial distribution θ_j from a Dirichlet distribution with parameter α
3. For each word position i (1 to N_j) in a document j :
 - a. Select a latent topic z_{ji} from the multinomial distribution parametrised by θ_j
 - b. Choose the observed word w_{ji} from the multinomial distribution parametrised by $\beta_{z_{ji}}$

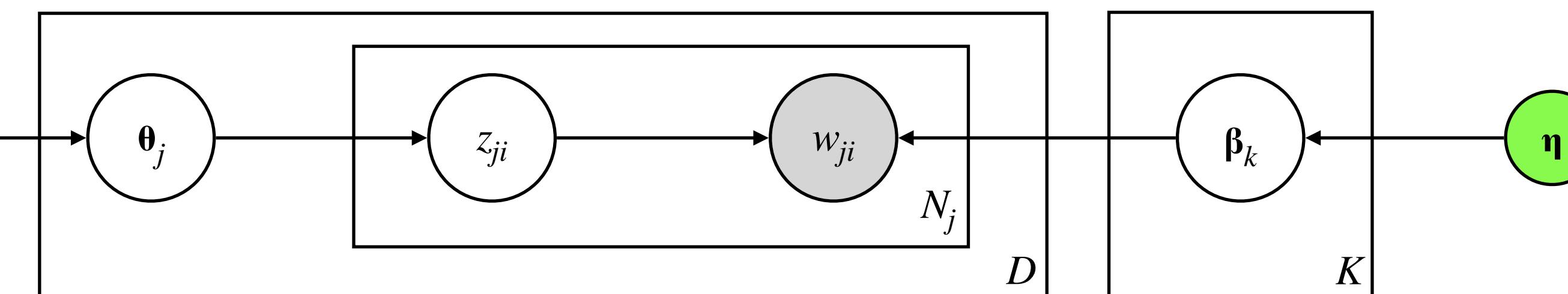
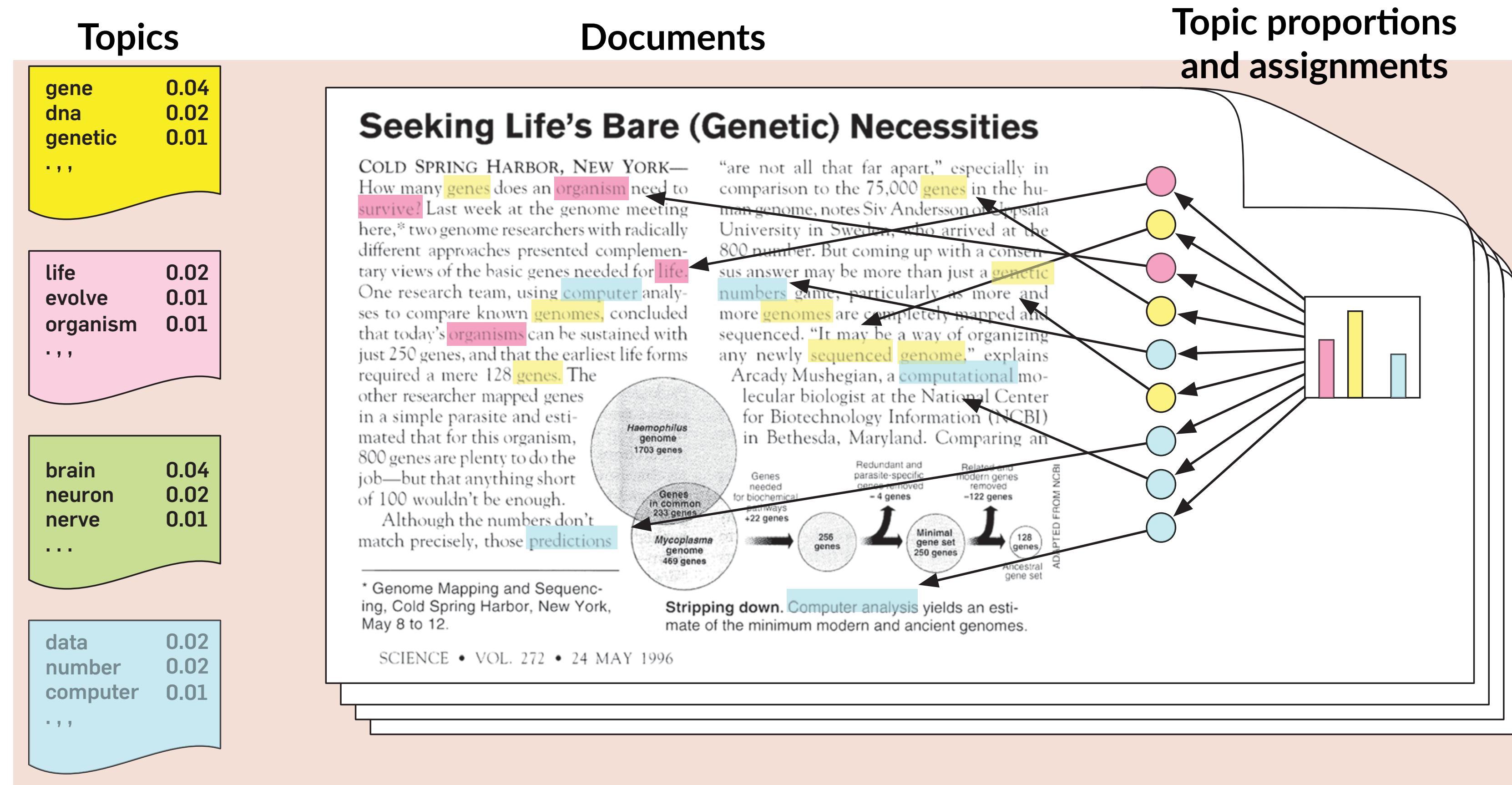
LDA – Generative story



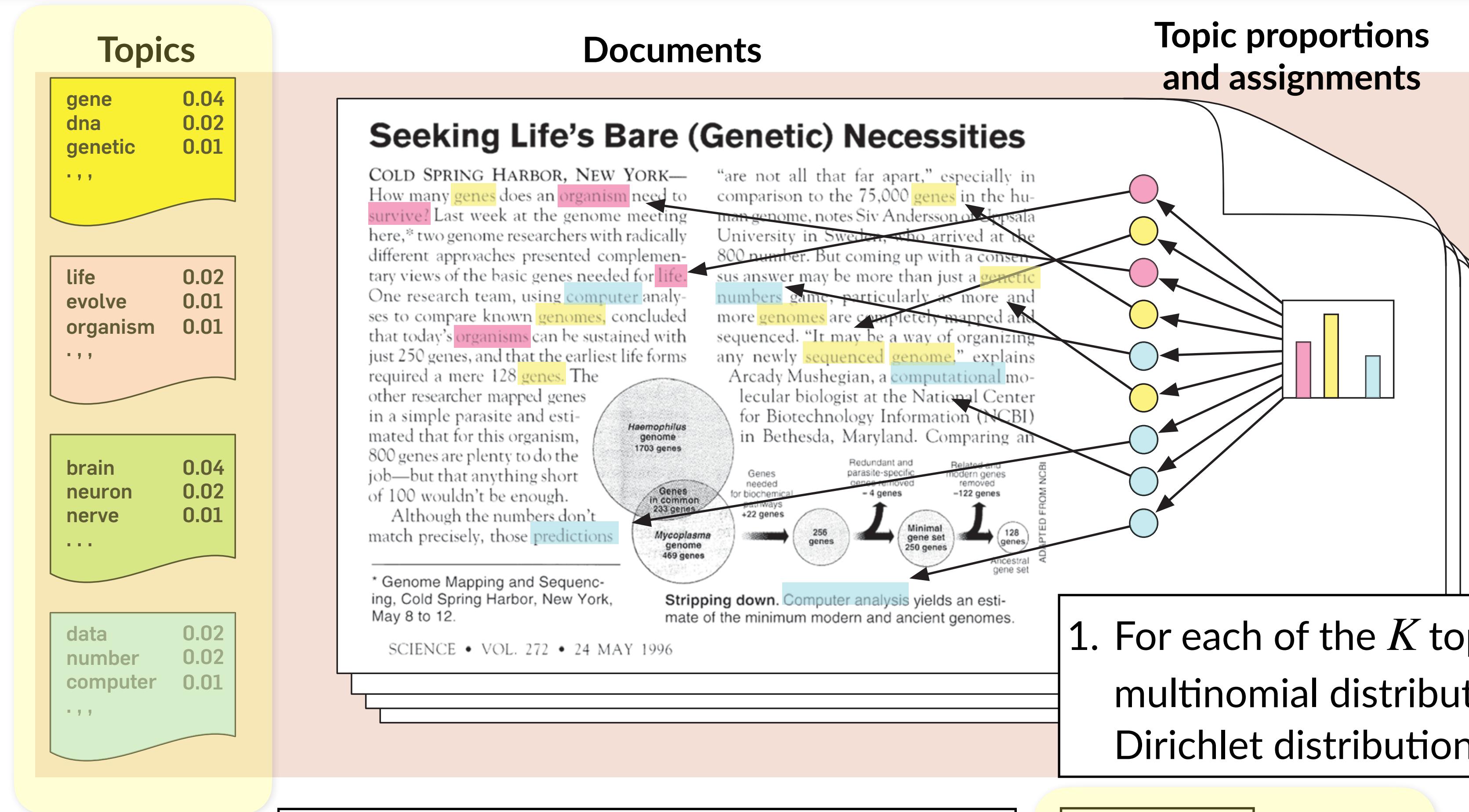
Assume a number of topics, defined as distributions over words (far left). A document is generated by first choosing a distribution over the topics (far right), then for each word position choosing a topic assignment (coloured coins), then choosing a word from the corresponding topic.

Blei. CACM, 2012, doi.org/10.1145/2133806.2133826

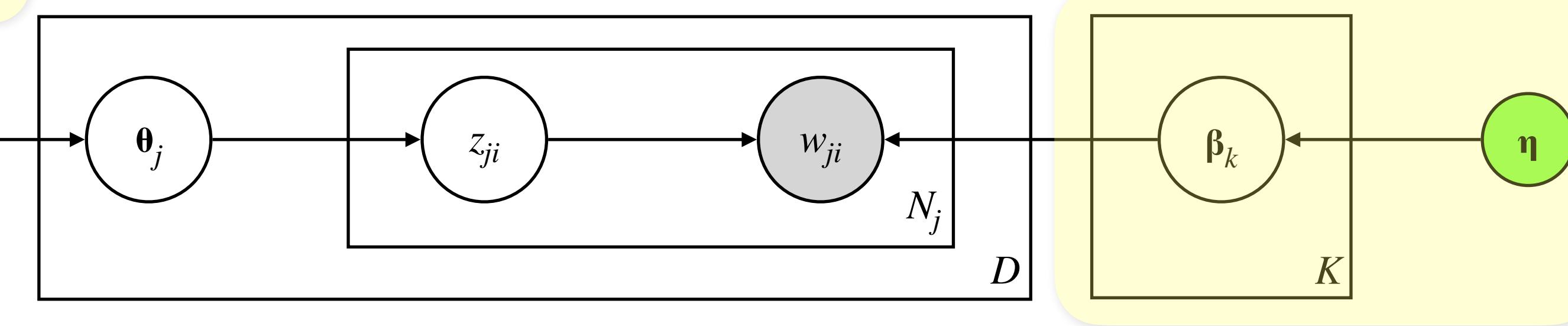
LDA – Generative story



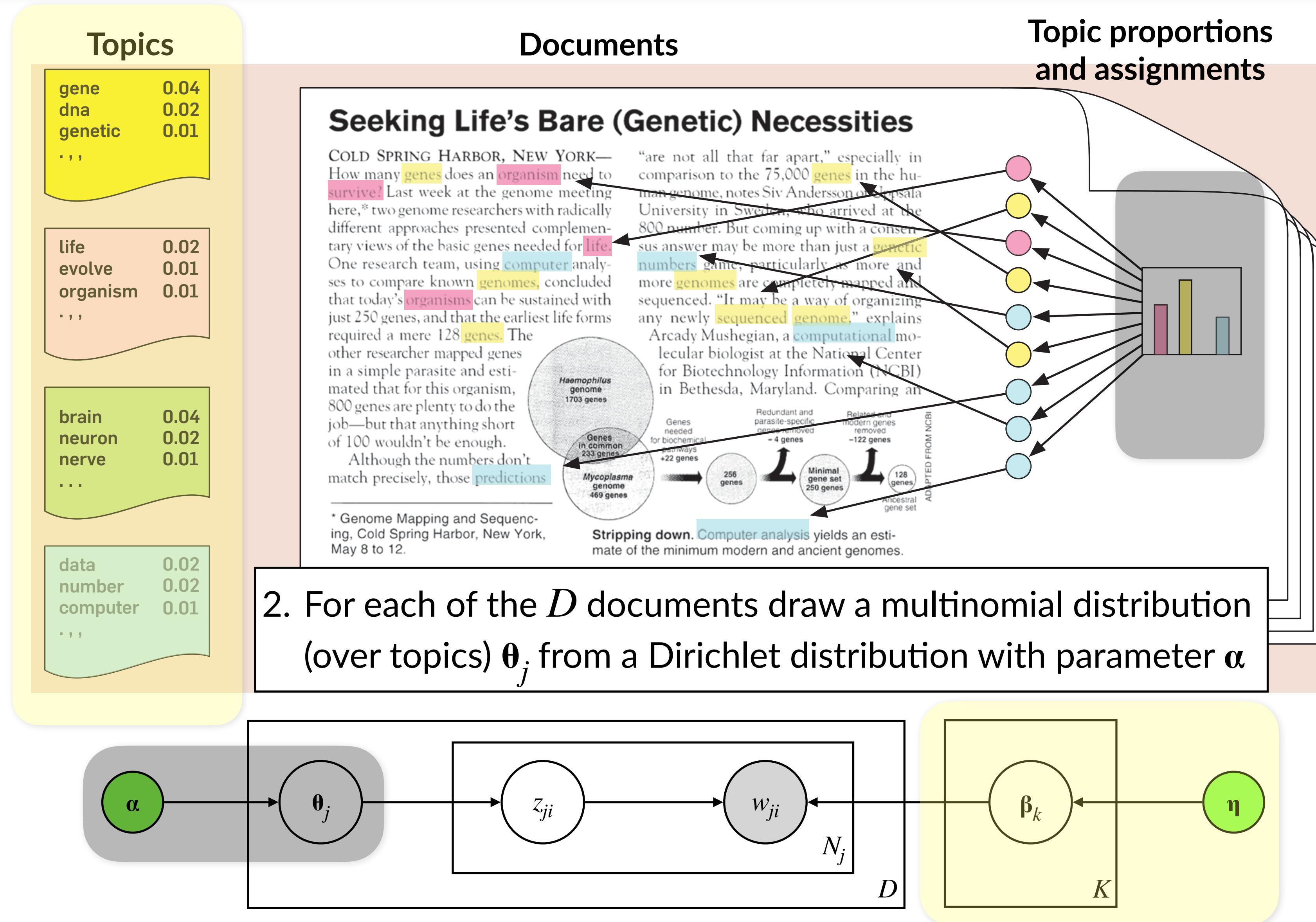
LDA – Generative story



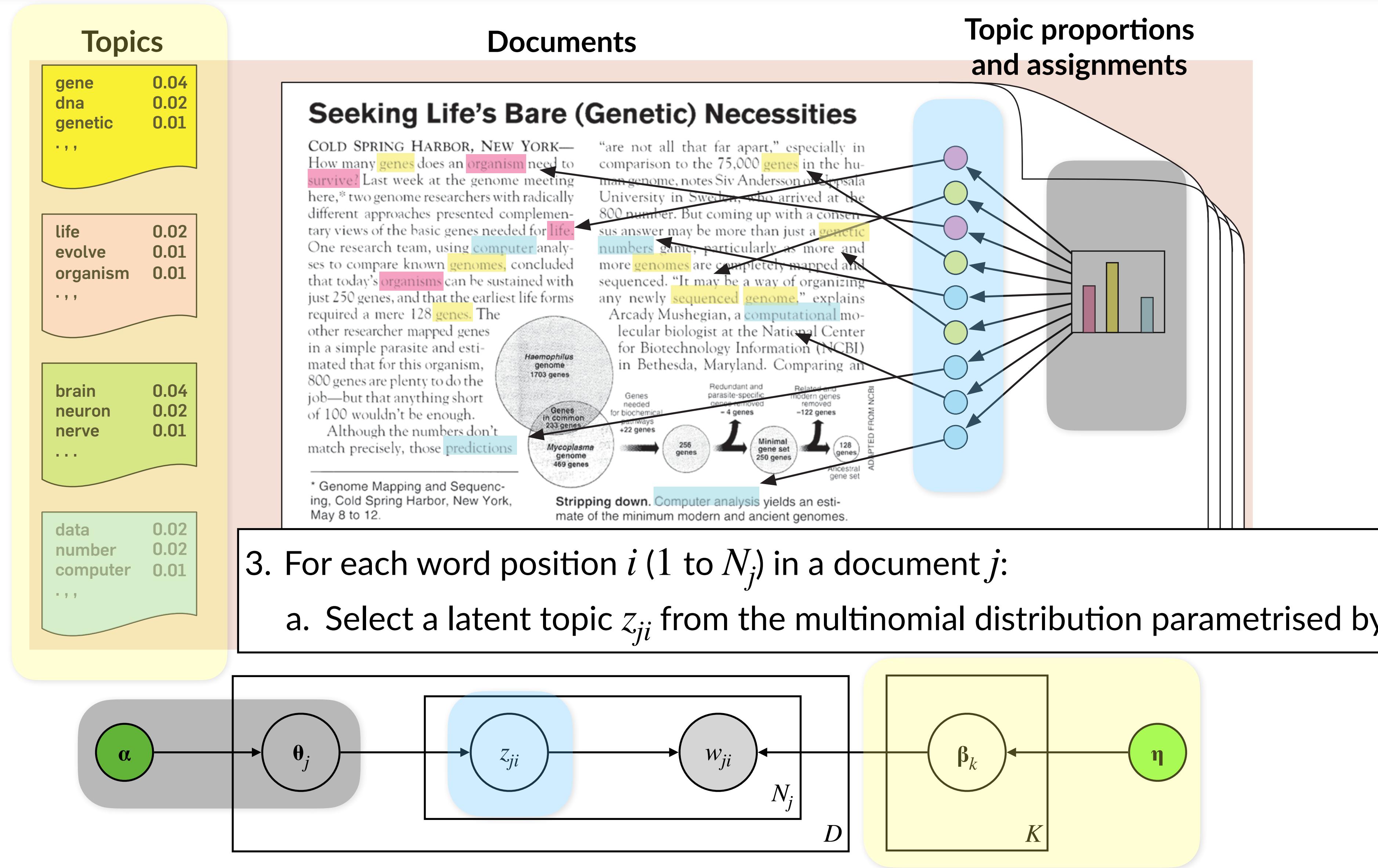
1. For each of the K topics draw a multinomial distribution β_k from a Dirichlet distribution with parameter η



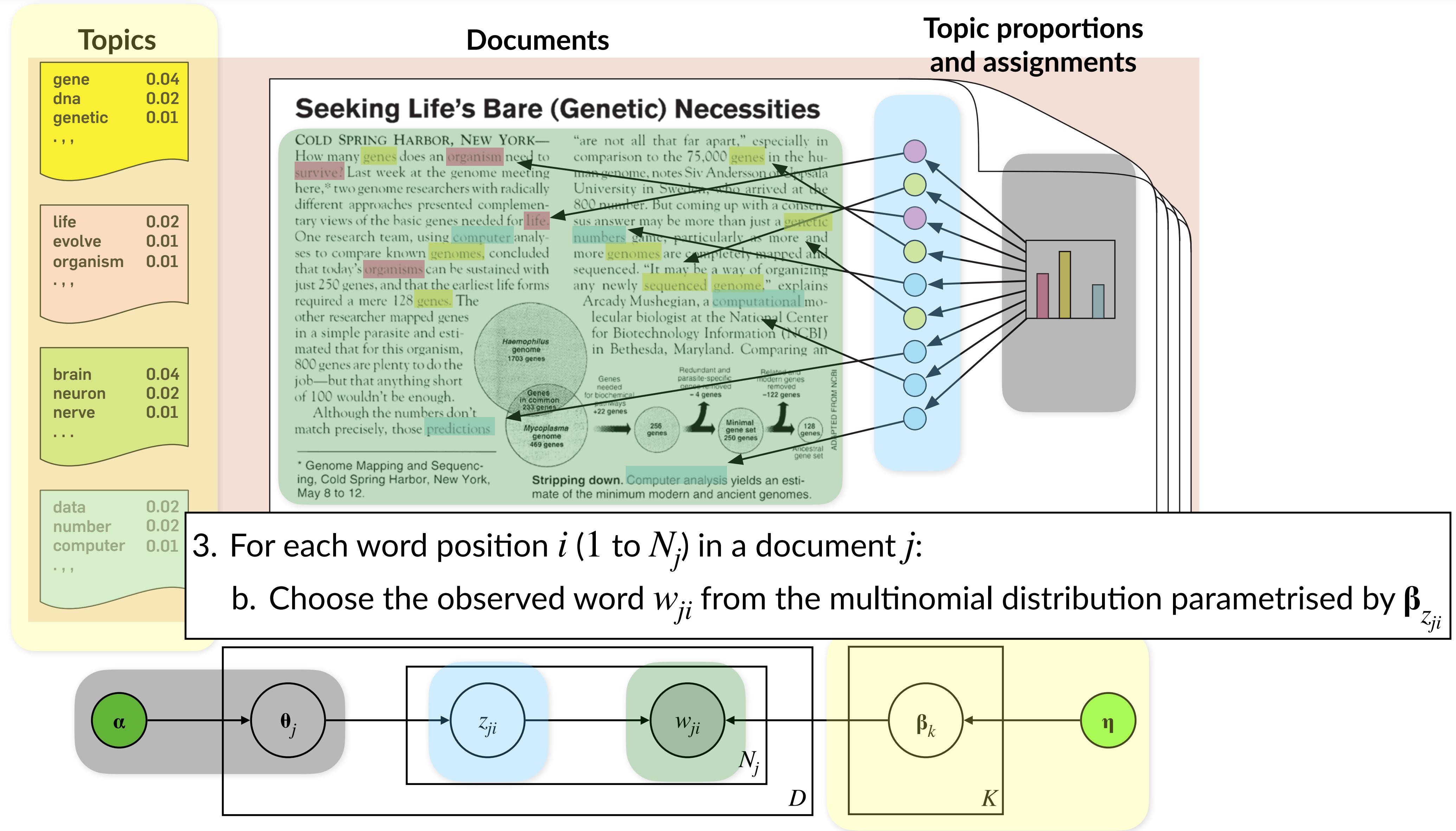
LDA – Generative story



LDA – Generative story

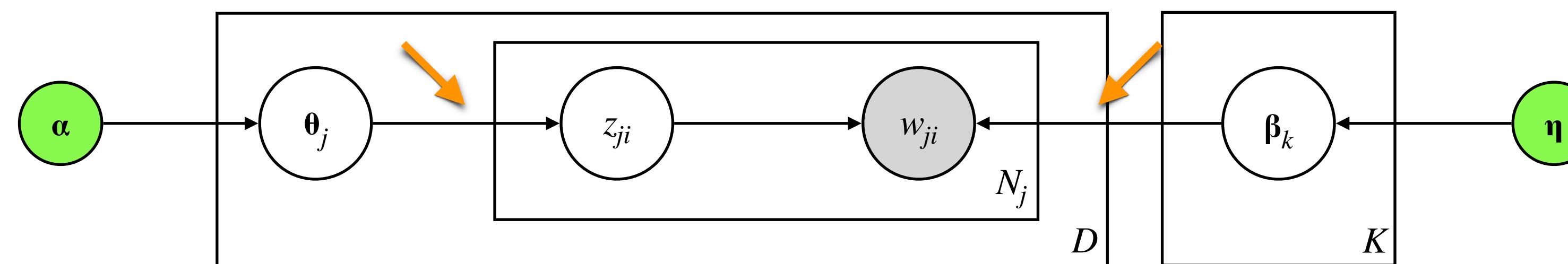


LDA – Generative story



3. For each word position i (1 to N_j) in a document j :
- b. Choose the observed word w_{ji} from the multinomial distribution parametrised by $\beta_{z_{ji}}$

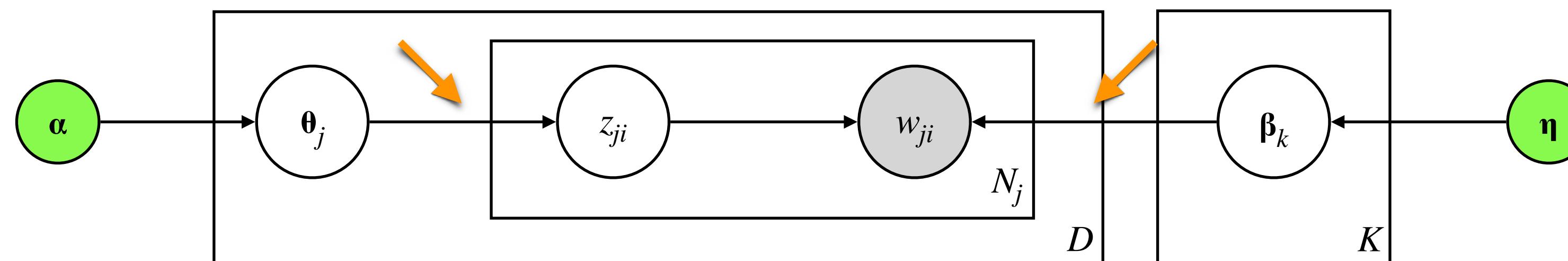
LDA – Multinomial distribution (Mult)



What is the probability of a set of outcomes for an event that has multiple outcomes?

- Roll a 6-sided dice 5 times. What is the probability of getting a “3” 1 time and a “6” 4 times?

LDA – Multinomial distribution (Mult)



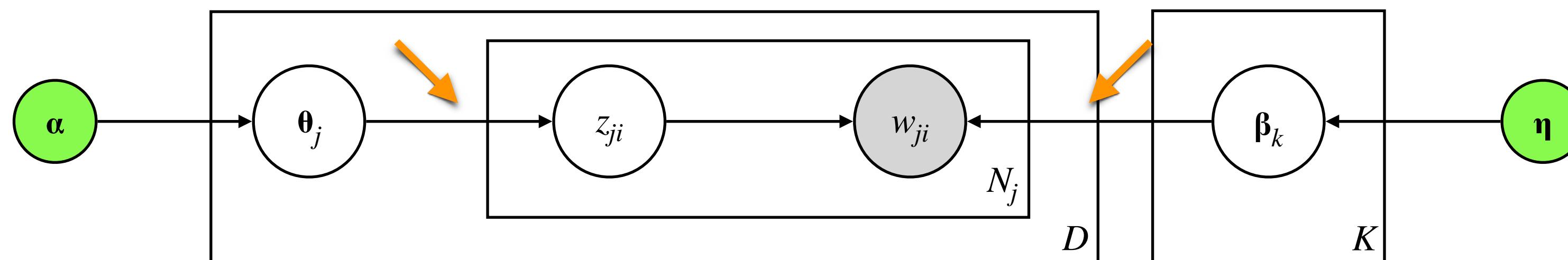
What is the probability of a set of outcomes for an event that has multiple outcomes?

- Roll a 6-sided dice 5 times. What is the probability of getting a “3” 1 time and a “6” 4 times?

$$\frac{5!}{1!4!} \cdot \left(\frac{1}{6}\right) \cdot \left(\frac{1}{6}\right)^4 \approx 0.00064$$

#ways to get 1 “3” and 4 “6”s prob. of 1 “3” prob. of 4 “6”s

LDA – Multinomial distribution (Mult)



What is the probability of a set of outcomes for an event that has multiple outcomes?

- Roll a 6-sided dice 5 times. What is the probability of getting a “3” 1 time and a “6” 4 times?

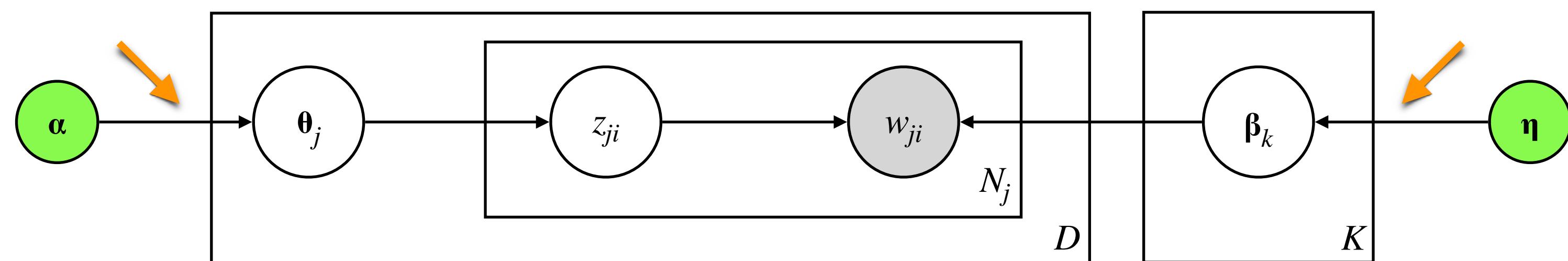
$$\frac{5!}{1!4!} \cdot \left(\frac{1}{6}\right) \cdot \left(\frac{1}{6}\right)^4 \approx 0.00064$$

#ways to get 1 “3” and 4 “6”s prob. of 1 “3” prob. of 4 “6”s

[youtube.com/watch?
v=5A_H1eHbOCY](https://youtube.com/watch?v=5A_H1eHbOCY)
10 min explanation by
Prof. John Tsitsiklis

Formally: $p(n_1, \dots, n_k) = \frac{n!}{n_1! \cdot \dots \cdot n_k!} \cdot p_1^{n_1} \cdot \dots \cdot p_k^{n_k}$ given $n, \{p_1, \dots, p_k\}$

LDA – Dirichlet distribution (Dir)



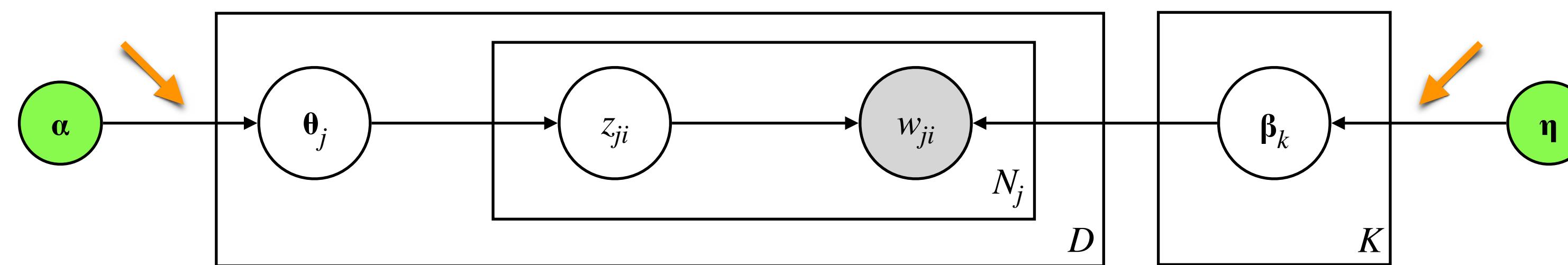
Exponential family distribution over the simplex (positive vectors with elements that sum up to 1), essentially a distribution over multinomial distributions

$$p(\boldsymbol{\theta}|\boldsymbol{\alpha}) = \frac{\Gamma\left(\sum_{k=1}^K \alpha_k\right)}{\prod_{k=1}^K \Gamma(\alpha_k)} \cdot \prod_{k=1}^K \theta_k^{\alpha_k-1} \text{ where } \Gamma(n) = (n-1)!$$

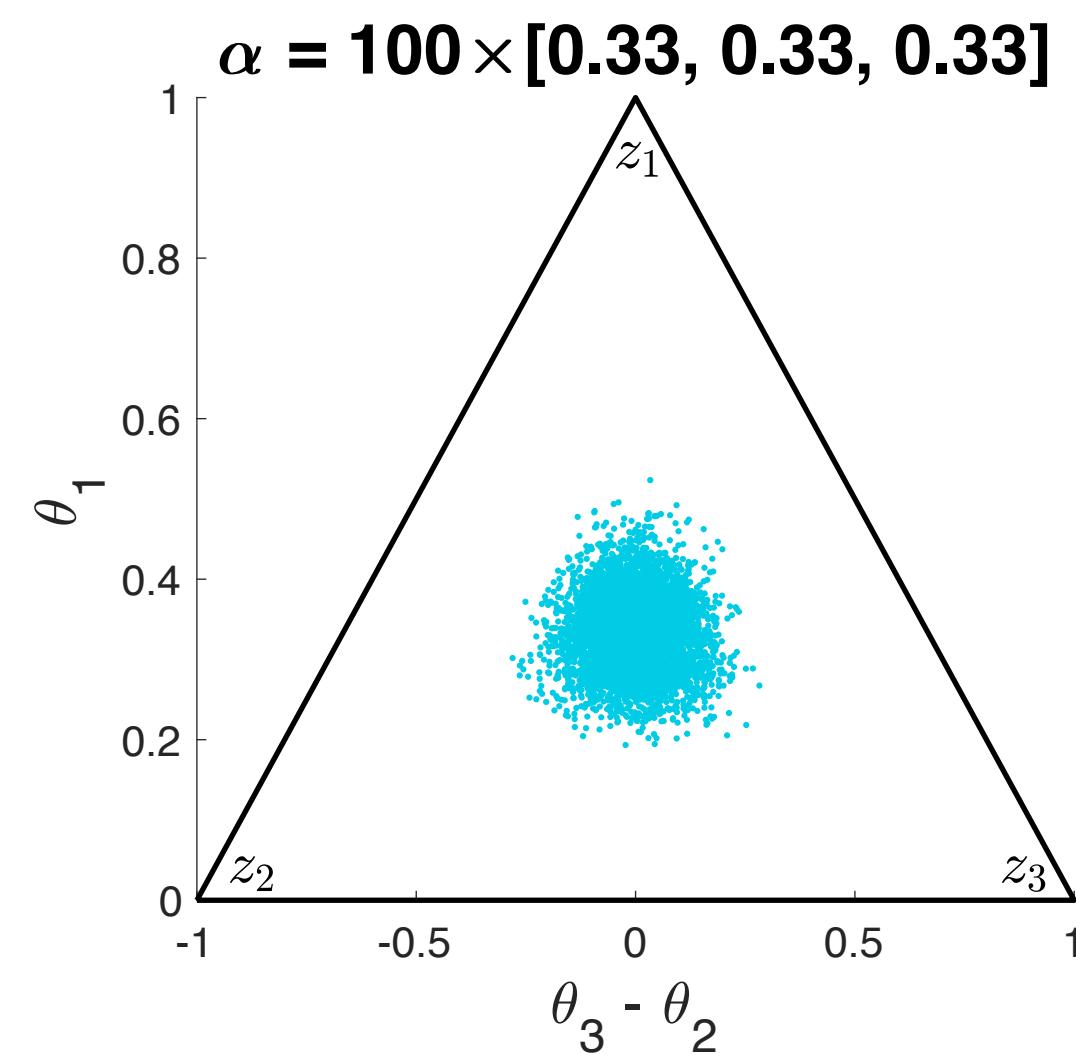
Parameter α controls the mean shape and sparsity of θ (same applies on η for β)

Note: α is a vector of K (= number of topics) parameters for θ and η has V parameters for β , where V is the size of the vocabulary (unique terms across all D documents)

LDA – Dirichlet distribution (Dir)

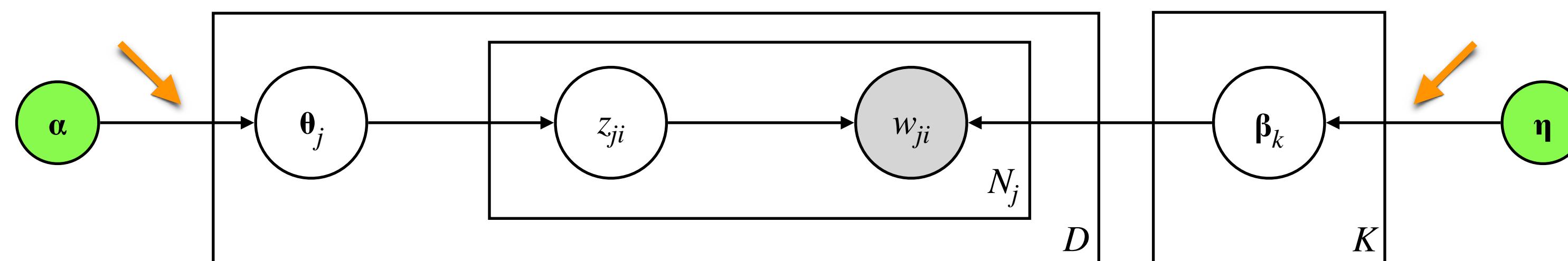


Assume a simplex $\Theta = [\theta_1, \theta_2, \theta_3]$ across $K = 3$ topics with $0 \leq \theta_i \leq 1$. How do different values for α affect the Θ produced by the Dirichlet distribution? Let's plot 5,000 samples for different α 's.

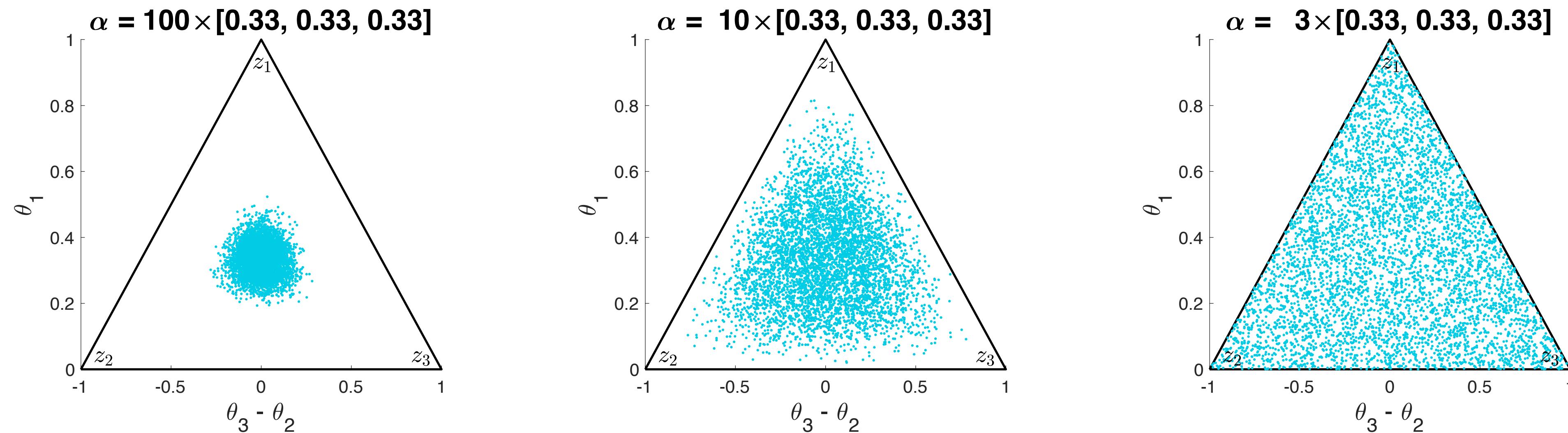


Large values of α lead to more dense Θ 's

LDA – Dirichlet distribution (Dir)

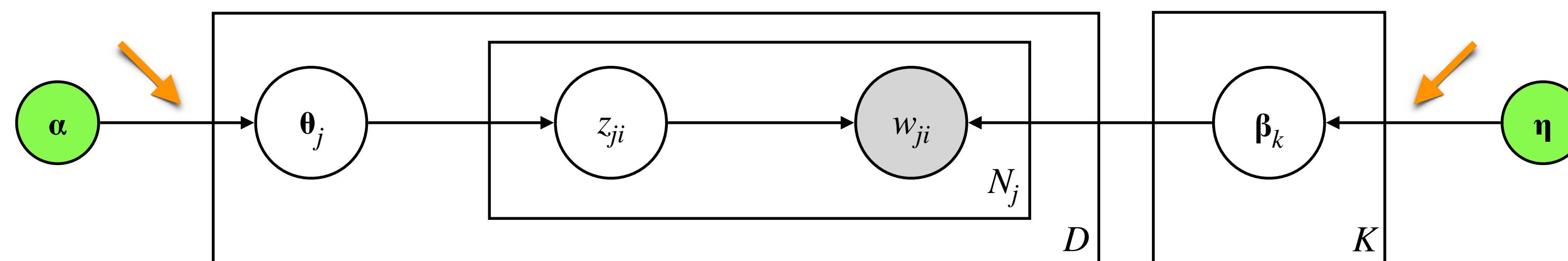


Assume a simplex $\theta = [\theta_1, \theta_2, \theta_3]$ across $K = 3$ topics with $0 \leq \theta_i \leq 1$. How do different values for α affect the θ produced by the Dirichlet distribution? Let's plot 5,000 samples for different α 's.

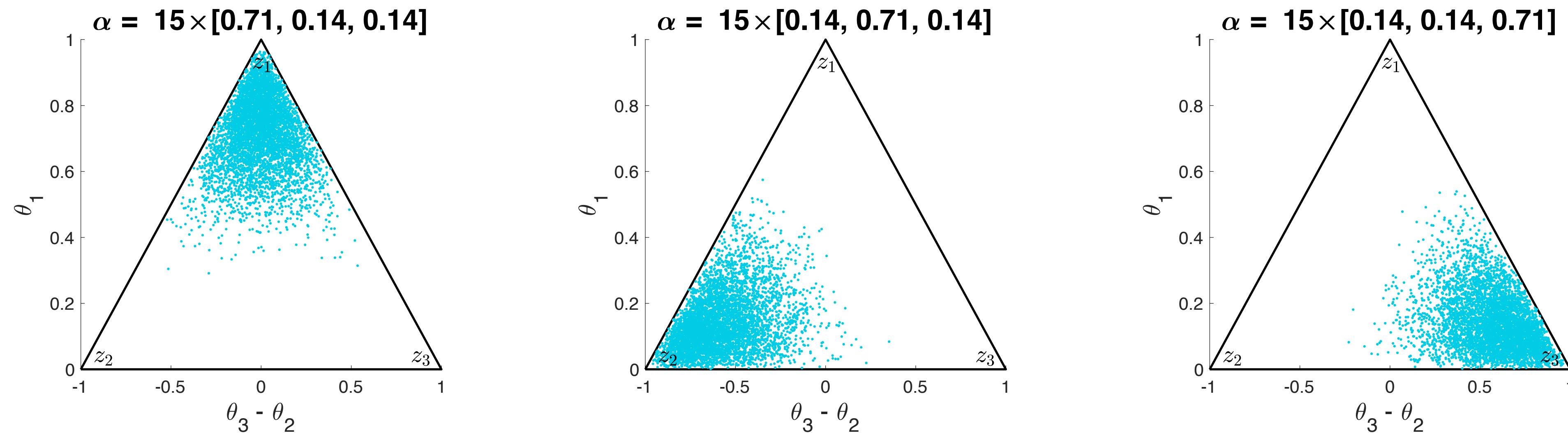


Large values of α lead to more dense θ 's

LDA – Dirichlet distribution (Dir)

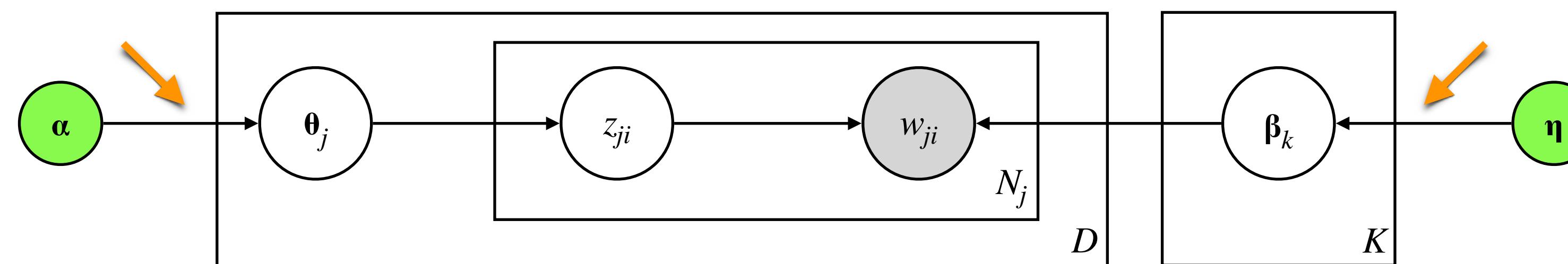


Assume a simplex $\theta = [\theta_1, \theta_2, \theta_3]$ across $K = 3$ topics with $0 \leq \theta_i \leq 1$. How do different values for α affect the θ produced by the Dirichlet distribution? Let's plot 5,000 samples for different α 's.

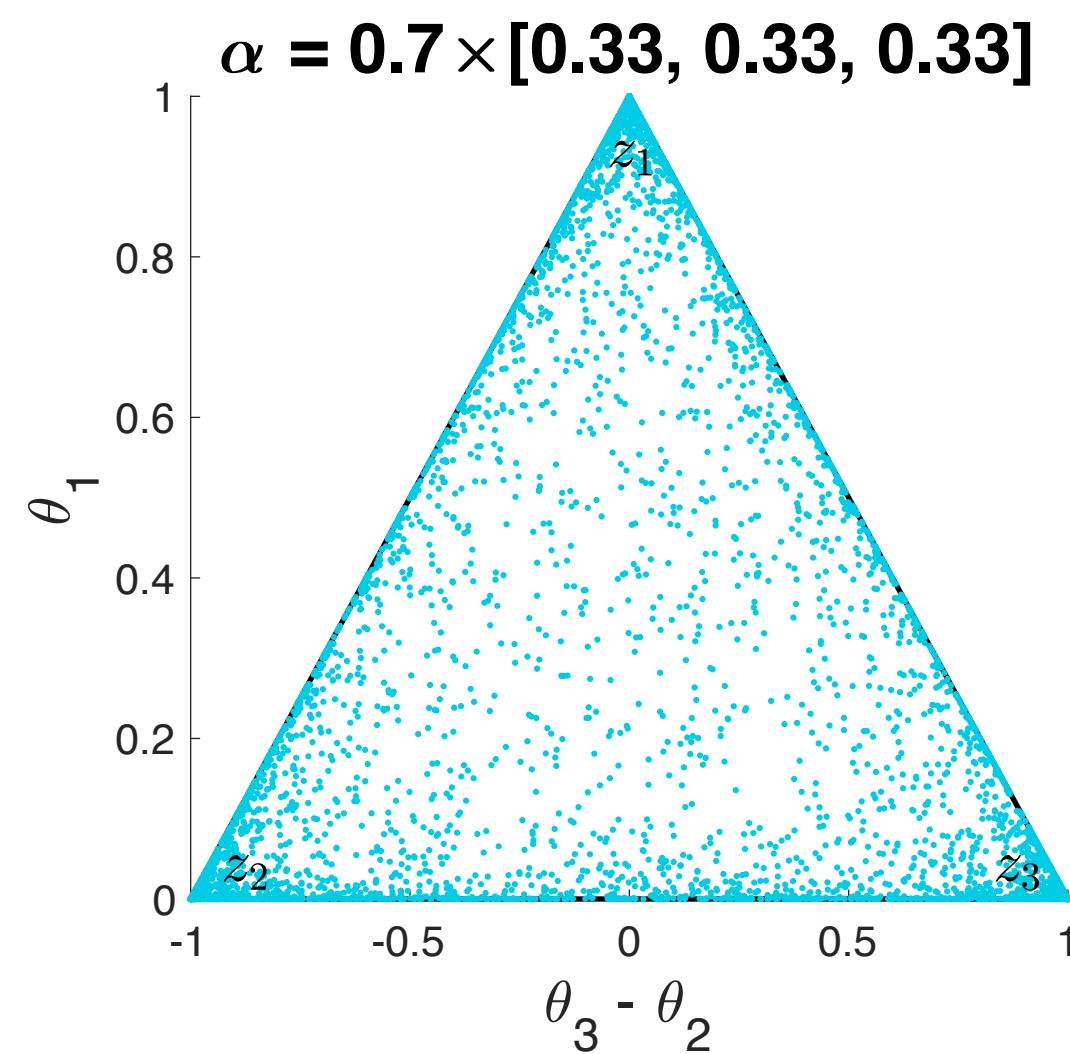


Imbalance in α shapes the focus of the distribution

LDA – Dirichlet distribution (Dir)

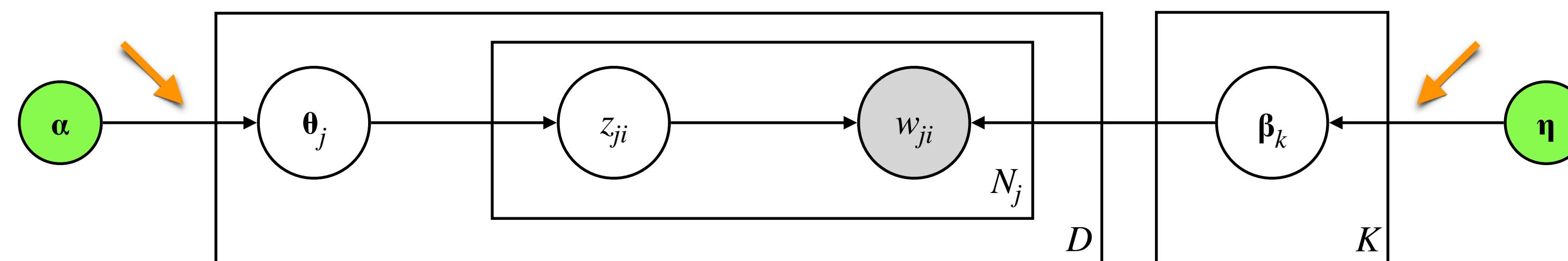


Assume a simplex $\theta = [\theta_1, \theta_2, \theta_3]$ across $K = 3$ topics with $0 \leq \theta_i \leq 1$. How do different values for α affect the θ produced by the Dirichlet distribution? Let's plot 5,000 samples for different α 's.

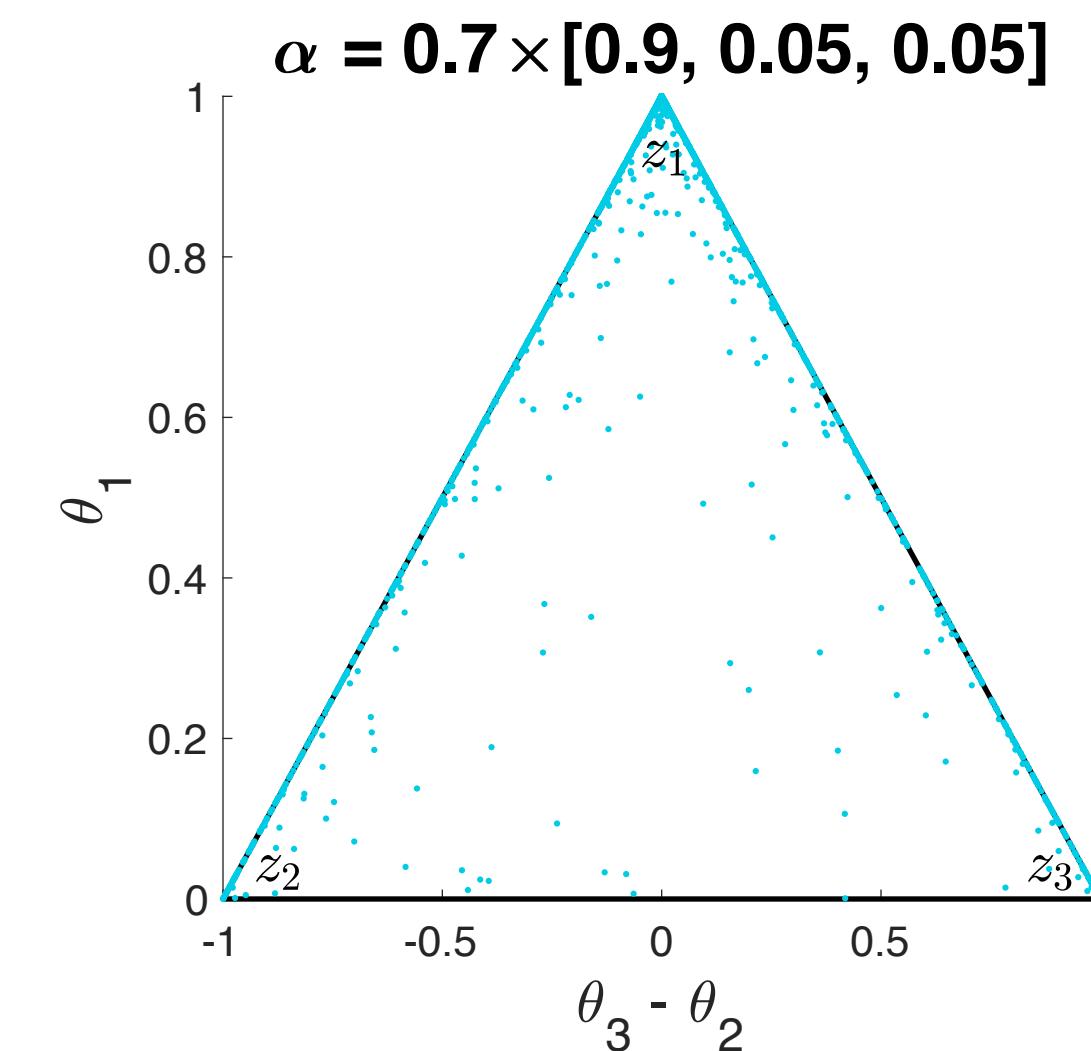
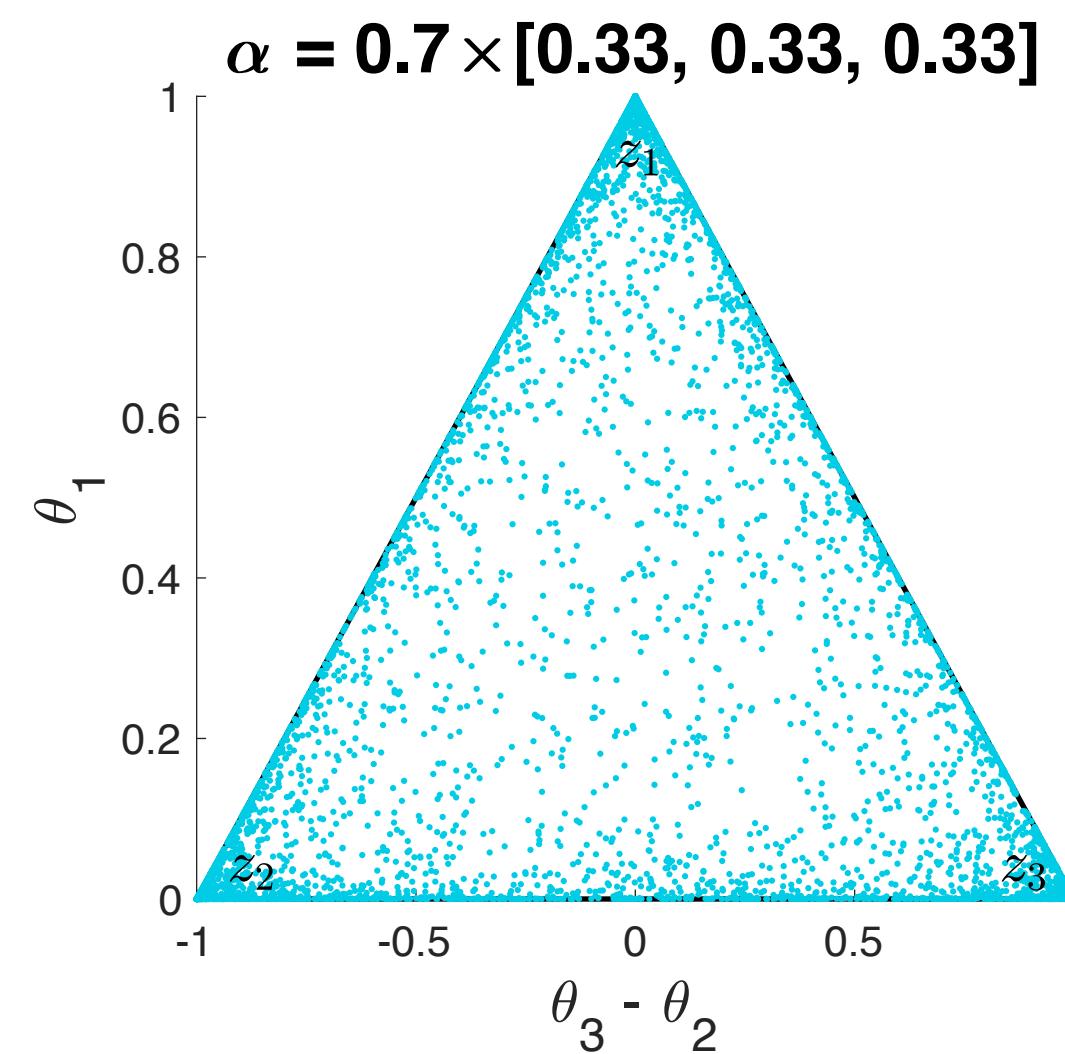


Values of $\alpha < 1$ create increasingly sparse outputs

LDA – Dirichlet distribution (Dir)

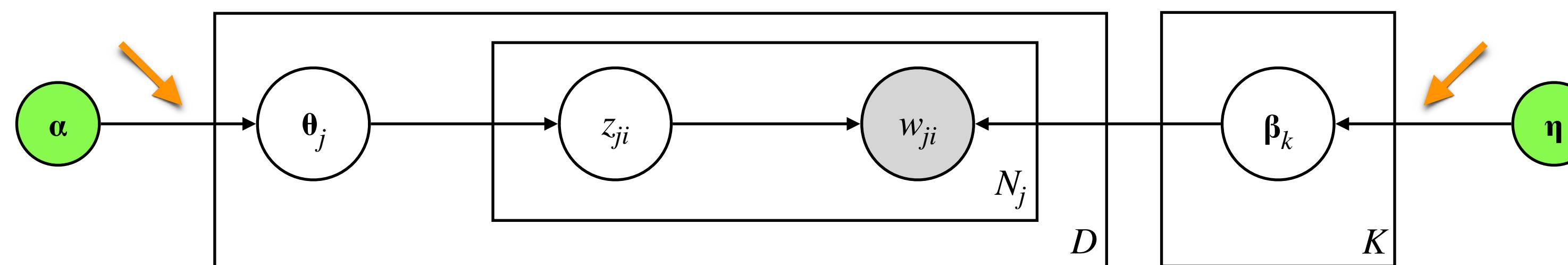


Assume a simplex $\theta = [\theta_1, \theta_2, \theta_3]$ across $K = 3$ topics with $0 \leq \theta_i \leq 1$. How do different values for α affect the θ produced by the Dirichlet distribution? Let's plot 5,000 samples for different α 's.

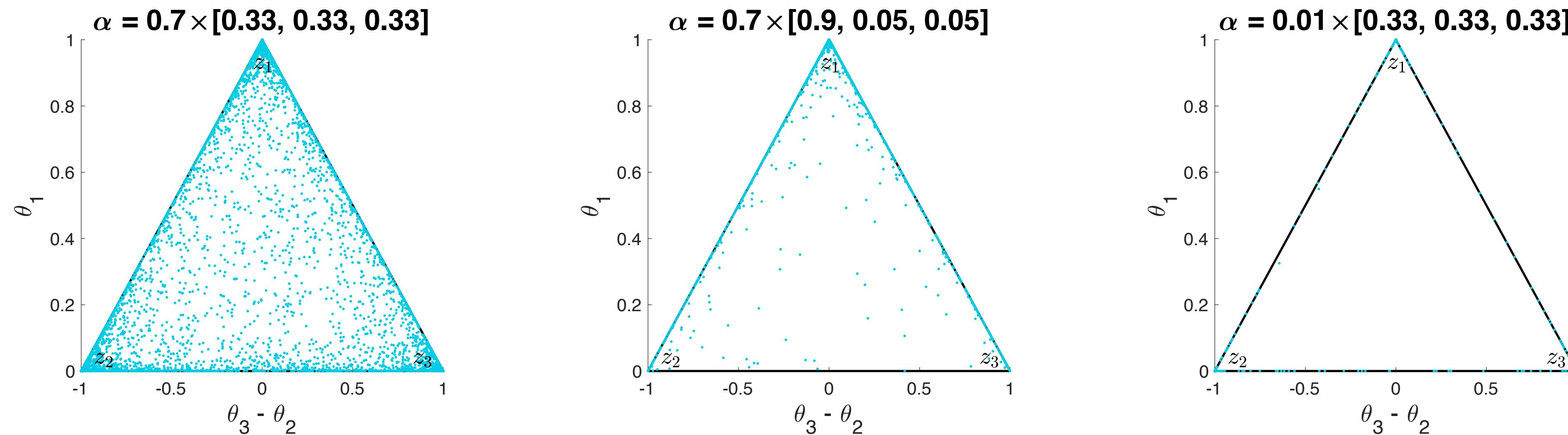


Values of $\alpha < 1$ create increasingly sparse outputs

LDA – Dirichlet distribution (Dir)

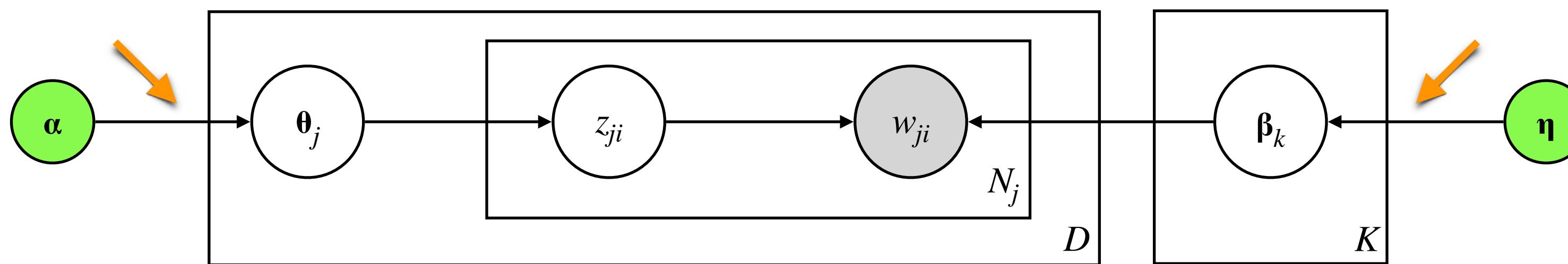


Assume a simplex $\theta = [\theta_1, \theta_2, \theta_3]$ across $K = 3$ topics with $0 \leq \theta_i \leq 1$. How do different values for α affect the θ produced by the Dirichlet distribution? Let's plot 5,000 samples for different α 's.



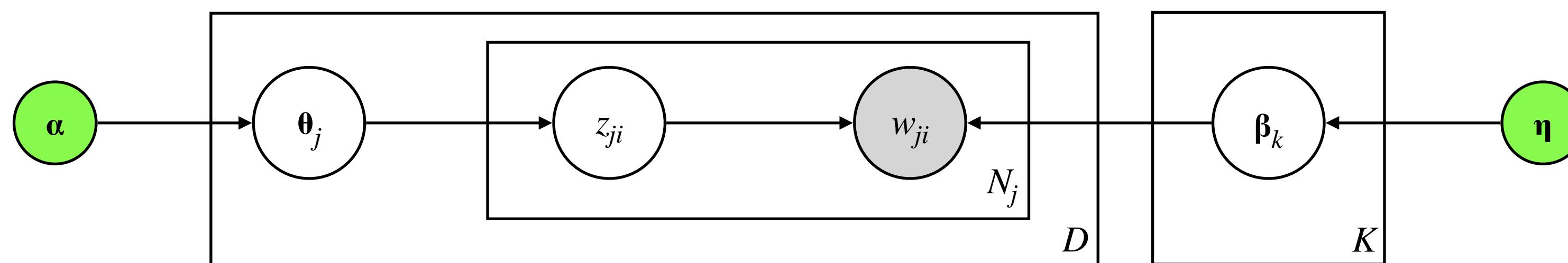
Values of $\alpha < 1$ create increasingly sparse outputs

LDA – Why do we combine Mult and Dir distributions?



- ▶ The Dirichlet distribution is **conjugate** to the Multinomial distribution
- ▶ Posterior $p(\beta|\eta, \mathbf{w})$ and prior $p(\beta|\eta)$ belong to the same distribution family as the prior (Dirichlet) given that $p(\mathbf{w}|\beta)$ is a Multinomial and $p(\beta|\eta)$ a Dirichlet
- ▶ Abstracting the math, observed data (\mathbf{w}) are adding to our prior intuition (η) about how words relate with topics

LDA – Inference



Joint probability distribution

$$p(\mathbf{W}, \Theta, \mathbf{B}, \mathbf{Z} | \alpha, \eta) = \prod_{k=1}^K p(\beta_k | \eta) \prod_{j=1}^D p(\theta_j | \alpha) \left(\prod_{i=1}^{N_j} p(z_{ji} | \theta_j) p(w_{ji} | \mathbf{B}, z_{ji}) \right)$$

We are interested in this posterior

$$p(\Theta, \mathbf{B}, \mathbf{Z} | \mathbf{W}, \alpha, \eta) = \frac{p(\Theta, \mathbf{B}, \mathbf{Z}, \mathbf{W} | \alpha, \eta)}{\int_{\mathbf{B}} \int_{\Theta} \sum_z p(\Theta, \mathbf{B}, \mathbf{Z}, \mathbf{W} | \alpha, \eta)}$$

can't compute → approximate inference

LDA – Inference, Gibbs sampling

- ▶ Initialise probabilities randomly or uniformly
- ▶ In each step, replace the value of one of the variables by a value drawn from the distribution of that variable conditioned on the values of the remaining variables
- ▶ Repeat until convergence

For $t = 1, \dots, T$:

$$\text{Sample } x_1^{(t+1)} \sim p(x_1 | x_2^{(t)}, \dots, x_N^{(t)})$$

$$\text{Sample } x_2^{(t+1)} \sim p(x_2 | x_1^{(t+1)}, x_3^{(t)}, \dots, x_N^{(t)})$$

...

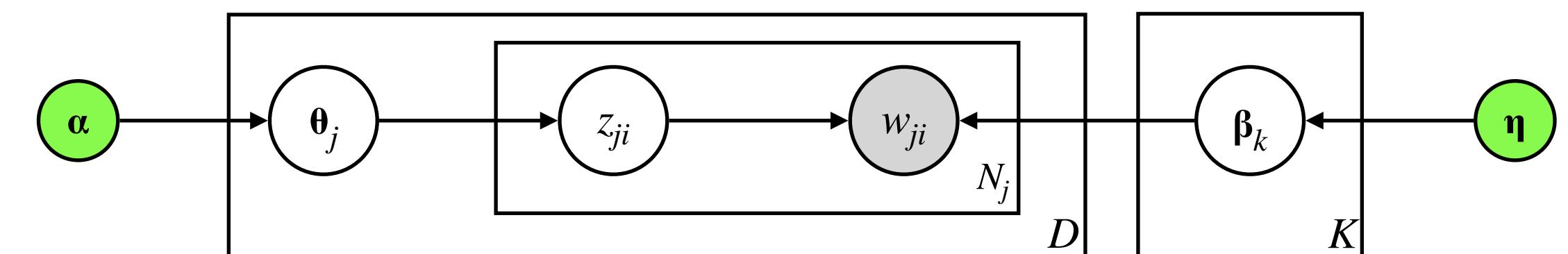
$$\text{Sample } x_j^{(t+1)} \sim p(x_j | x_1^{(t+1)}, x_2^{(t+1)}, \dots, x_{j-1}^{(t+1)}, x_{j+1}^{(t)}, \dots, x_N^{(t)})$$

...

$$\text{Sample } x_N^{(t+1)} \sim p(x_N | x_1^{(t+1)}, \dots, x_{N-1}^{(t+1)})$$

LDA – Inference, Gibbs sampling

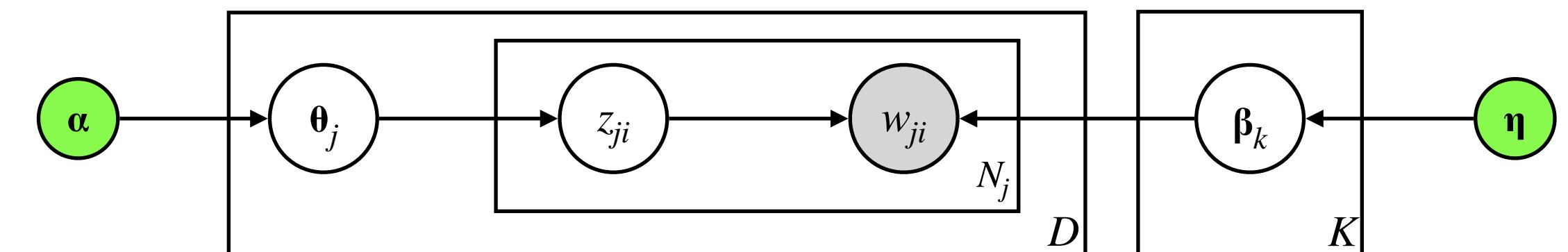
- ▶ Initialise probabilities randomly or uniformly
- ▶ Go over each word i in every document j (w_{ji})
- ▶ Estimate the probability of assigning w_{ji} to each topic, conditioned on the topic assignments ($\mathbf{z}_{j,-i}$) of all other words $\mathbf{w}_{j,-i}$ (notation indicating the exclusion of w_{ji})



$$p(z_{ji} = k | \mathbf{z}_{j,-i}, \mathbf{W}, \boldsymbol{\alpha}, \boldsymbol{\eta}) \propto \frac{n_{j,k,-i} + \alpha_k}{\sum_{k'=1}^K n_{j,k',-i} + \alpha_{k'}} \cdot \frac{m_{k,w_{ji},-i} + \eta_{w_{ji}}}{\sum_{\nu=1}^V m_{k,\nu,-i} + \eta_{\nu}}$$

LDA – Inference, Gibbs sampling

- ▶ Initialise probabilities randomly or uniformly
- ▶ Go over each word i in every document j (w_{ji})
- ▶ Estimate the probability of assigning w_{ji} to each topic, conditioned on the topic assignments ($\mathbf{z}_{j,-i}$) of all other words $\mathbf{w}_{j,-i}$ (notation indicating the exclusion of w_{ji})



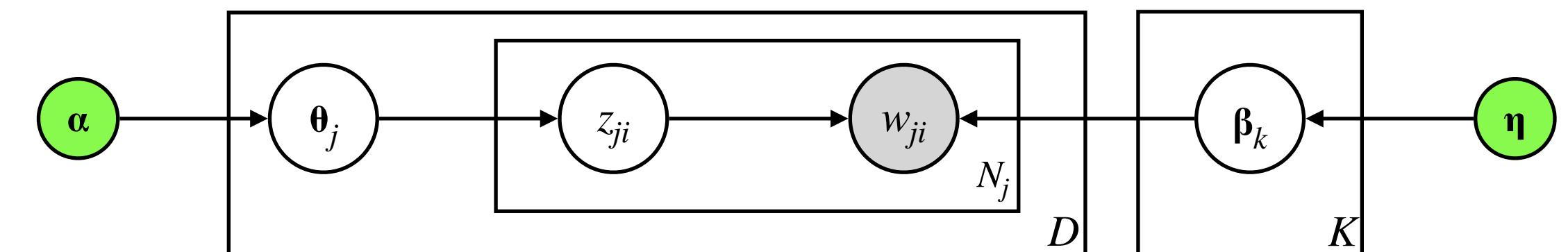
$$p(z_{ji} = k | \mathbf{z}_{j,-i}, \mathbf{W}, \alpha, \eta) \propto \frac{n_{j,k,-i} + \alpha_k}{\sum_{k'=1}^K n_{j,k',-i} + \alpha_{k'}} \cdot \frac{m_{k,w_{ji},-i} + \eta_{w_{ji}}}{\sum_{\nu=1}^V m_{k,\nu,-i} + \eta_{\nu}}$$

topic k is assigned to a word in document j without counting the current word

word w_{ji} is associated with topic k in all documents without counting the current instance of w_{ji}

LDA – Inference, Gibbs sampling

- ▶ Initialise probabilities randomly or uniformly
- ▶ Go over each word i in every document j (w_{ji})
- ▶ Estimate the probability of assigning w_{ji} to each topic, conditioned on the topic assignments ($\mathbf{z}_{j,-i}$) of all other words $\mathbf{w}_{j,-i}$ (notation indicating the exclusion of w_{ji})



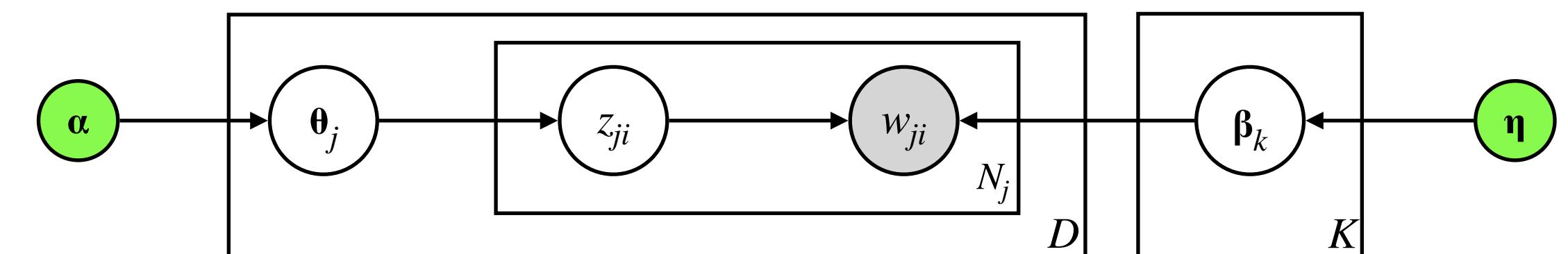
$$p(z_{ji} = k | \mathbf{z}_{j,-i}, \mathbf{W}, \alpha, \eta) \propto \frac{n_{j,k,-i} + \alpha_k}{\sum_{k'=1}^K n_{j,k',-i} + \alpha_{k'}} \cdot \frac{m_{k,w_{ji},-i} + \eta_{w_{ji}}}{\sum_{\nu=1}^V m_{k,\nu,-i} + \eta_{\nu}}$$

How much does document
 j “like” topic k ?

How much does topic
 k “like” word w_{ji} ?

LDA – Inference, Gibbs sampling

- ▶ Initialise probabilities randomly or uniformly
- ▶ Go over each word i in every document j (w_{ji})
- ▶ Estimate the probability of assigning w_{ji} to each topic, conditioned on the topic assignments ($\mathbf{z}_{j,-i}$) of all other words $\mathbf{w}_{j,-i}$ (notation indicating the exclusion of w_{ji})



$$p(z_{ji} = k | \mathbf{z}_{j,-i}, \mathbf{W}, \boldsymbol{\alpha}, \boldsymbol{\eta}) \propto \frac{n_{j,k,-i} + \alpha_k}{\sum_{k'=1}^K n_{j,k',-i} + \alpha_{k'}} \cdot \frac{m_{k,w_{ji},-i} + \eta_{w_{ji}}}{\sum_{\nu=1}^V m_{k,\nu,-i} + \eta_{\nu}}$$

- ▶ From the above conditional distribution, sample a topic and set it as the new topic assignment z_{ji} of w_{ji}

How much does topic k “like” word w_{ji} ?

How much does document j “like” topic k ?

LDA – Gibbs sampling, toy example

- Consider $K = 3$ topics
- ...

LDA – Gibbs sampling, toy example

- Consider $K = 3$ topics
- Sampling from document j (word order doesn't matter)
- ...

document j

z_{ji}	?	?	?	?	?
w_{ji}	Brexit	deficit	Europe	market	single

LDA – Gibbs sampling, toy example

- Consider $K = 3$ topics
- Sampling from document j (*word order doesn't matter*)
- **Randomly assign topics to all words in document j (and all other docs)**
- ...

document j

z_{ji}	3	?	?	?	?
w_{ji}	Brexit	deficit	Europe	market	single

LDA – Gibbs sampling, toy example

- Consider $K = 3$ topics
- Sampling from document j (word order doesn't matter)
- **Randomly assign topics to all words in document j (and all other docs)**
- ...

document j	z_{ji}	3	2	?	?	?
	w_{ji}	Brexit	deficit	Europe	market	single

LDA – Gibbs sampling, toy example

- Consider $K = 3$ topics
- Sampling from document j (*word order doesn't matter*)
- **Randomly assign topics to all words in document j (and all other docs)**
- ...

document j

z_{ji}	3	2	3	1	1
w_{ji}	Brexit	deficit	Europe	market	single

LDA – Gibbs sampling, toy example

- Consider $K = 3$ topics
- Sampling from document j (*word order doesn't matter*)
- Randomly assign topics to all words in document j (*and all other docs*)
- **Update the word-topic counts for all documents**
- ...

document j

z_{ji}	3	2	3	1	1
w_{ji}	Brexit	deficit	Europe	market	single

word-topic counts across all documents

words / topics	1	2	3
Brexit	100	30	2
deficit	10	60	0
Europe	95	5	2
market	50	70	5
single	50	15	90
...

LDA – Gibbs sampling, toy example

- Consider $K = 3$ topics
- Sampling from document j (*word order doesn't matter*)
- Randomly assign topics to all words in document j (*and all other docs*)
- Update the word-topic counts for all documents
- **Sample the first word (“Brexit”) in document j ; unassign it from topic 3 and decrement its count in the word-topic counts**
- ...

document j

z_{ji}	3	2	3	1	1
w_{ji}	Brexit	deficit	Europe	market	single

word-topic counts
across all documents

words / topics	1	2	3
Brexit	100	30	2
deficit	10	60	0
Europe	95	5	2
market	50	70	5
single	50	15	90
...

LDA – Gibbs sampling, toy example

- Consider $K = 3$ topics
- Sampling from document j (*word order doesn't matter*)
- Randomly assign topics to all words in document j (*and all other docs*)
- Update the word-topic counts for all documents
- **Sample the first word (“Brexit”) in document j ; unassign it from topic 3 and decrement its count in the word-topic counts**
- ...

document j

z_{ji}	?	2	3	1	1
w_{ji}	Brexit	deficit	Europe	market	single

word-topic counts
across all documents

words / topics	1	2	3
Brexit	100	30	2 – 1
deficit	10	60	0
Europe	95	5	2
market	50	70	5
single	50	15	90
...

LDA – Gibbs sampling, toy example

- Consider $K = 3$ topics
- Sampling from document j (*word order doesn't matter*)
- Randomly assign topics to all words in document j (*and all other docs*)
- Update the word-topic counts for all documents
- **Sample the first word (“Brexit”) in document j ; unassign it from topic 3 and decrement its count in the word-topic counts**
- ...

document j

z_{ji}	?	2	3	1	1
w_{ji}	Brexit	deficit	Europe	market	single

word-topic counts
across all documents

words / topics	1	2	3
Brexit	100	30	1
deficit	10	60	0
Europe	95	5	2
market	50	70	5
single	50	15	90
...

LDA – Gibbs sampling, toy example

- ...
- Randomly assign topics to all words in document j (*and all other docs*)
- Update the word-topic counts for all documents
- Sample the first word (“Brexit”) in document j ; unassign it from topic 3 and decrement its count in the word-topic counts
- **What are the revised topic proportions in document j ?**
- ...

document j	z_{ji}	?	2	3	1	1
	w_{ji}	Brexit	deficit	Europe	market	single

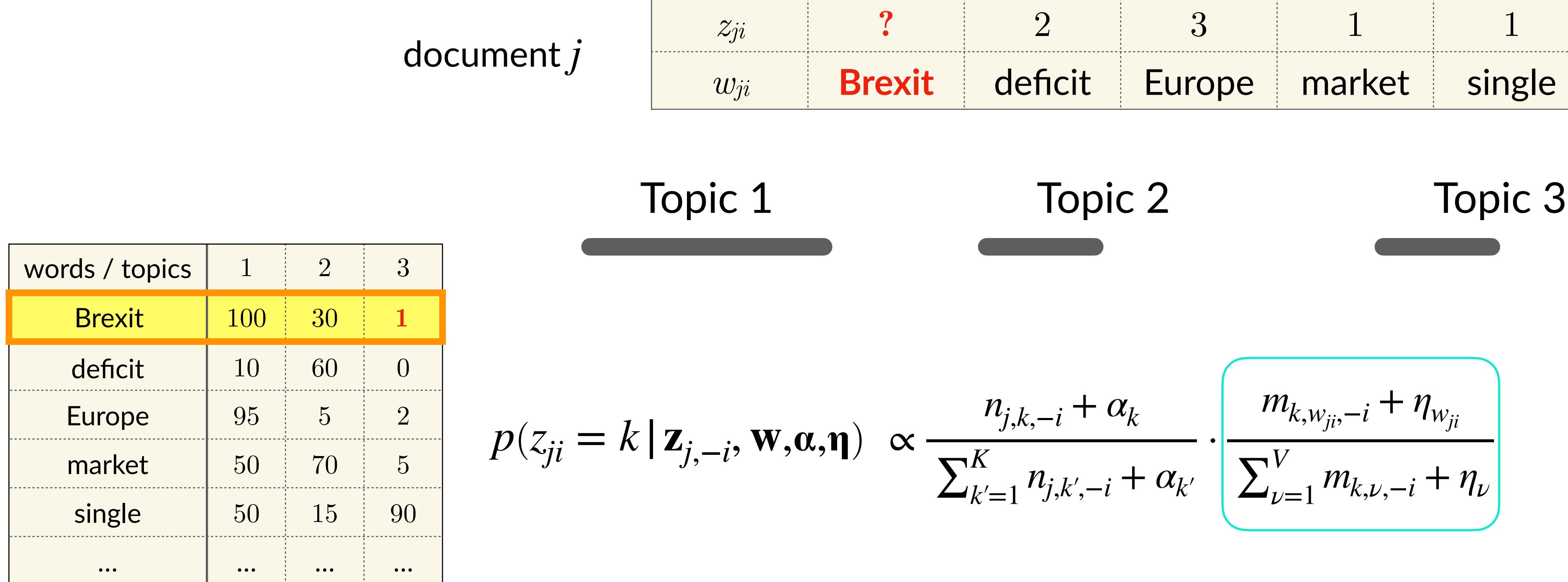
Topic 1 Topic 2 Topic 3



$$p(z_{ji} = k | \mathbf{z}_{j,-i}, \mathbf{w}, \boldsymbol{\alpha}, \boldsymbol{\eta}) \propto \frac{n_{j,k,-i} + \alpha_k}{\sum_{k'=1}^K n_{j,k',-i} + \alpha_{k'}} \cdot \frac{m_{k,w_{ji},-i} + \eta_{w_{ji}}}{\sum_{\nu=1}^V m_{k,\nu,-i} + \eta_{\nu}}$$

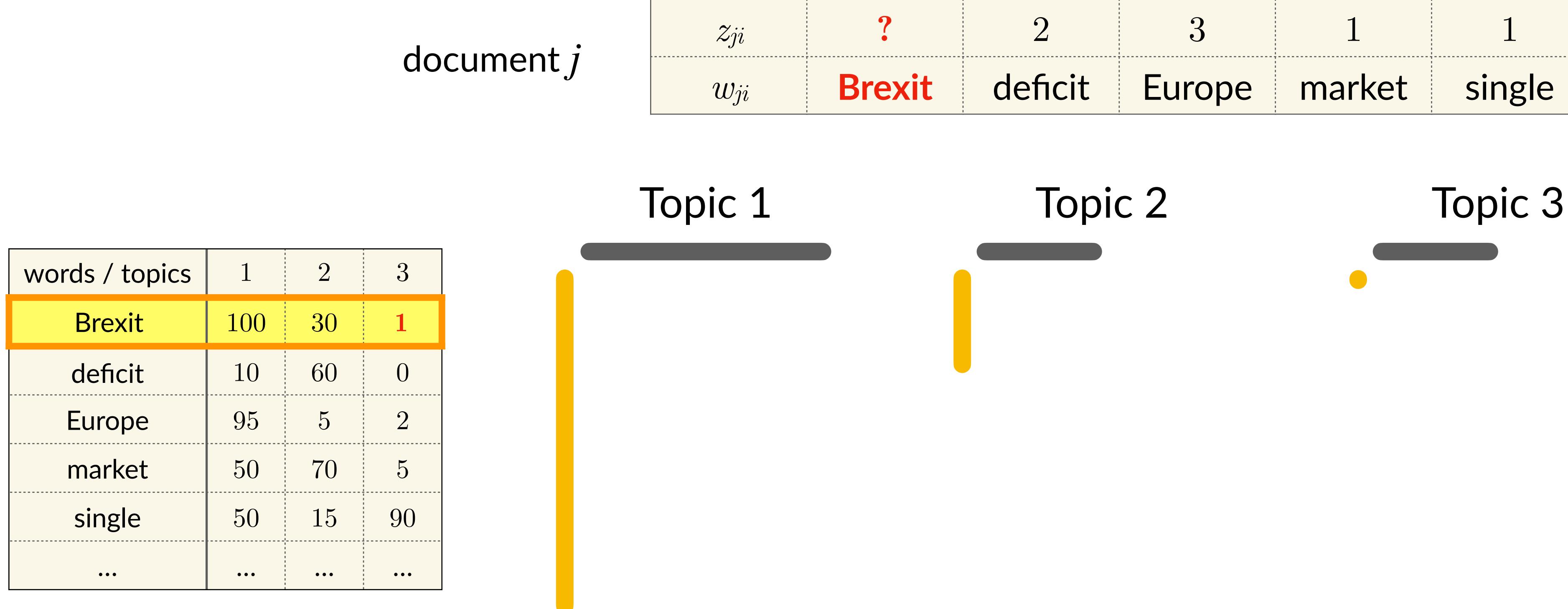
LDA – Gibbs sampling, toy example

- ...
- Update the word-topic counts for all documents
- Sample the first word (“Brexit”) in document j ; unassign it from topic 3 and decrement its count in the word-topic counts
- What are the revised topic proportions in document j ?
- **How much does each topic “like” the word Brexit?**
- ...



LDA – Gibbs sampling, toy example

- ...
- Update the word-topic counts for all documents
- Sample the first word (“Brexit”) in document j ; unassign it from topic 3 and decrement its count in the word-topic counts
- What are the revised topic proportions in document j ?
- **How much does each topic “like” the word Brexit?**
- ...



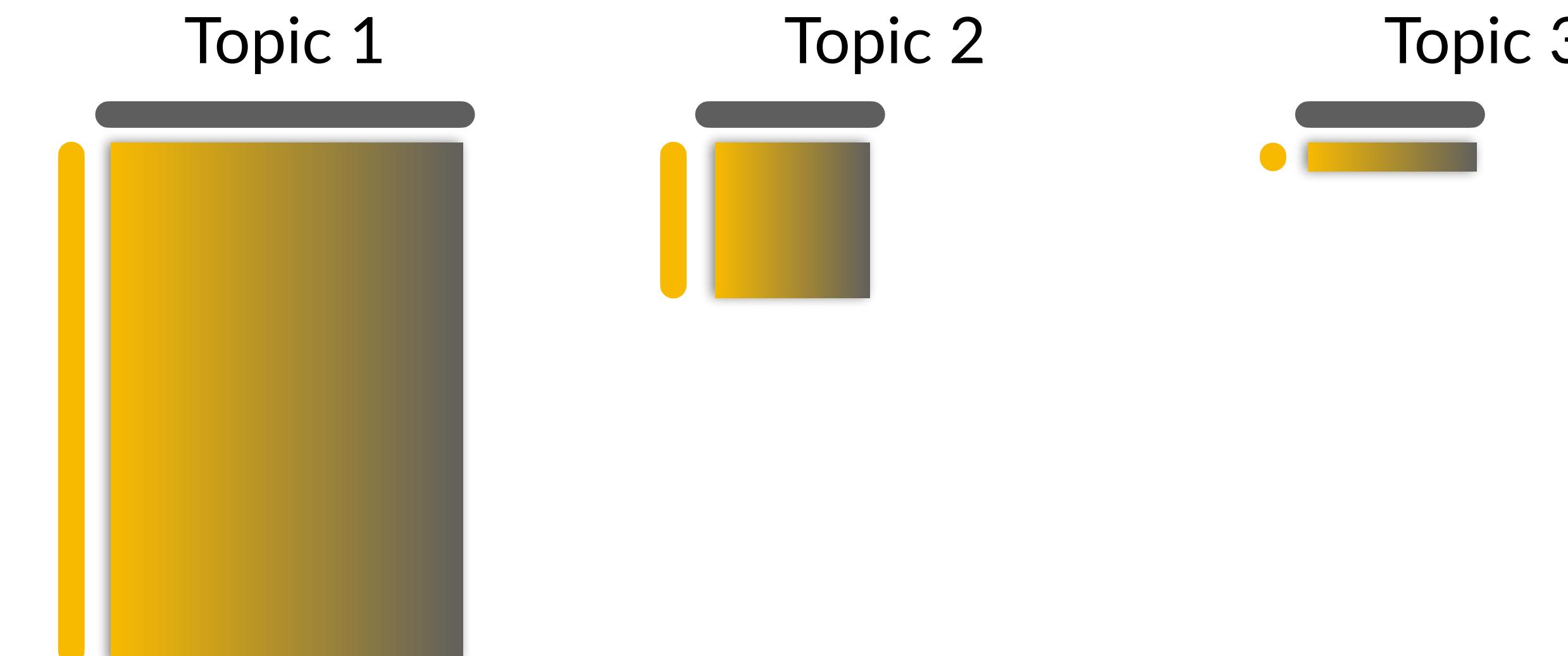
LDA – Gibbs sampling, toy example

- ...
- Update the word-topic counts for all documents
- Sample the first word (“Brexit”) in document j ; unassign it from topic 3 and decrement its count in the word-topic counts
- What are the revised topic proportions in document j ?
- **How much does each topic “like” the word Brexit?**
- ...

document j

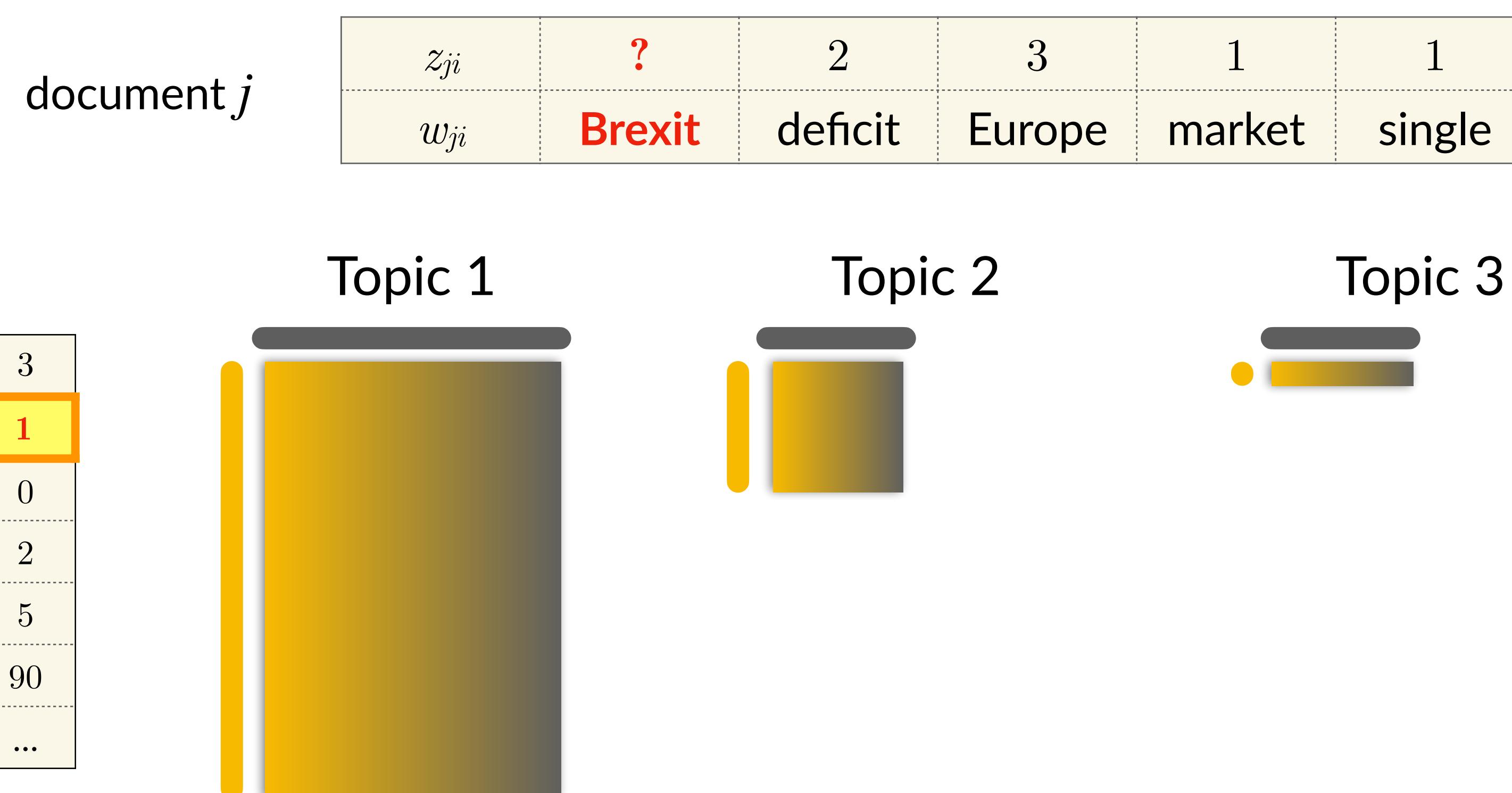
z_{ji}	?	2	3	1	1
w_{ji}	Brexit	deficit	Europe	market	single

words / topics	1	2	3
Brexit	100	30	1
deficit	10	60	0
Europe	95	5	2
market	50	70	5
single	50	15	90
...



LDA – Gibbs sampling, toy example

- ...
- Sample the first word (“Brexit”) in document j ; unassign it from topic 3 and decrement its count in the word-topic counts
- What are the revised topic proportions in document j ?
- How much does each topic “like” the word Brexit?
- **Sample from the revised conditional distribution** $p(z_{ji} = k | \mathbf{z}_{j,-i}, \mathbf{W}, \boldsymbol{\alpha}, \boldsymbol{\eta})$
- ...



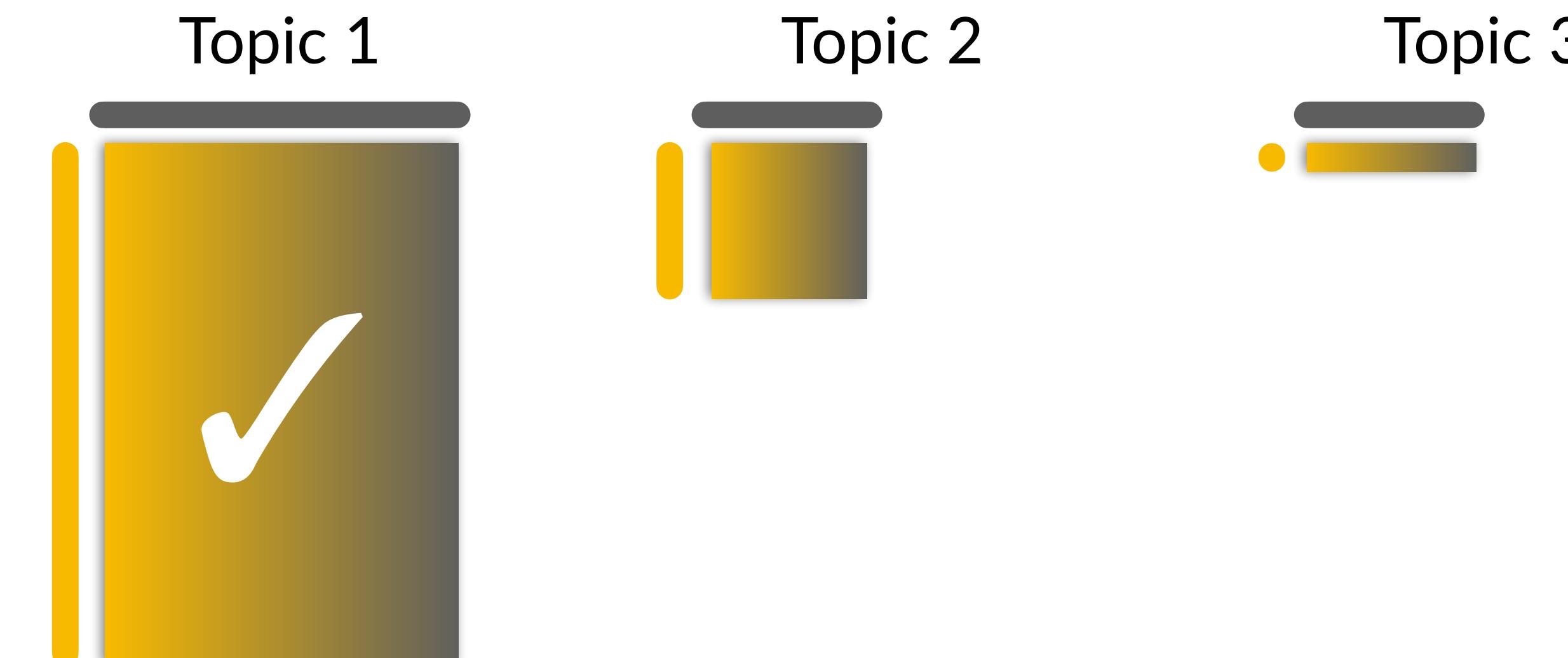
LDA – Gibbs sampling, toy example

- ...
- Sample the first word (“Brexit”) in document j ; unassign it from topic 3 and decrement its count in the word-topic counts
- What are the revised topic proportions in document j ?
- How much does each topic “like” the word Brexit?
- **Sample from the revised conditional distribution** $p(z_{ji} = k | \mathbf{z}_{j,-i}, \mathbf{W}, \boldsymbol{\alpha}, \boldsymbol{\eta})$
- ...

document j

z_{ji}	?	2	3	1	1
w_{ji}	Brexit	deficit	Europe	market	single

words / topics	1	2	3
Brexit	100	30	1
deficit	10	60	0
Europe	95	5	2
market	50	70	5
single	50	15	90
...

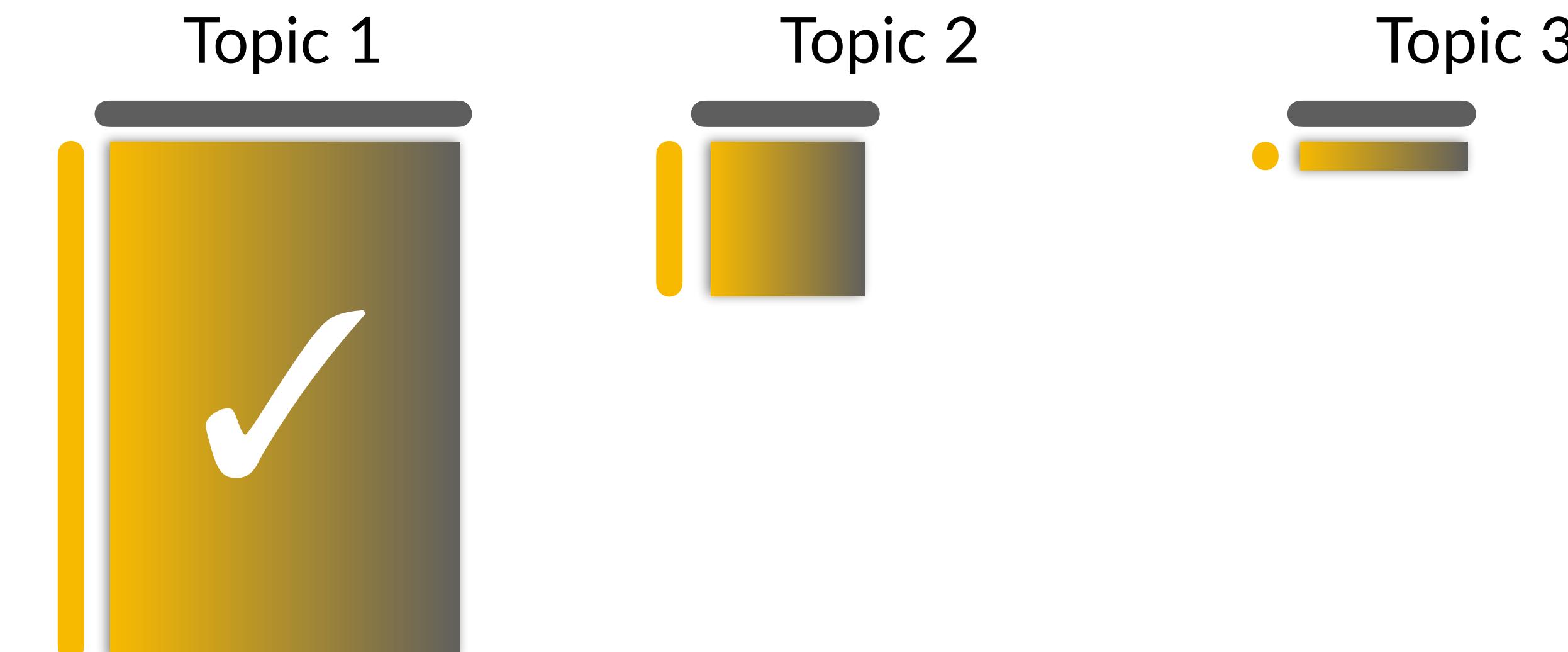


LDA – Gibbs sampling, toy example

- ...
- What are the revised topic proportions in document j ?
- How much does each topic “like” the word Brexit?
- Sample from the revised conditional distribution $p(z_{ji} = k | \mathbf{z}_{j,-i}, \mathbf{W}, \alpha, \eta)$
- **Assign the sampled topic to the word “Brexit” and update counts**
- ...

document j	z_{ji}	1	2	3	1	1
	w_{ji}	Brexit	deficit	Europe	market	single

words / topics	1	2	3
Brexit	101	30	1
deficit	10	60	0
Europe	95	5	2
market	50	70	5
single	50	15	90
...



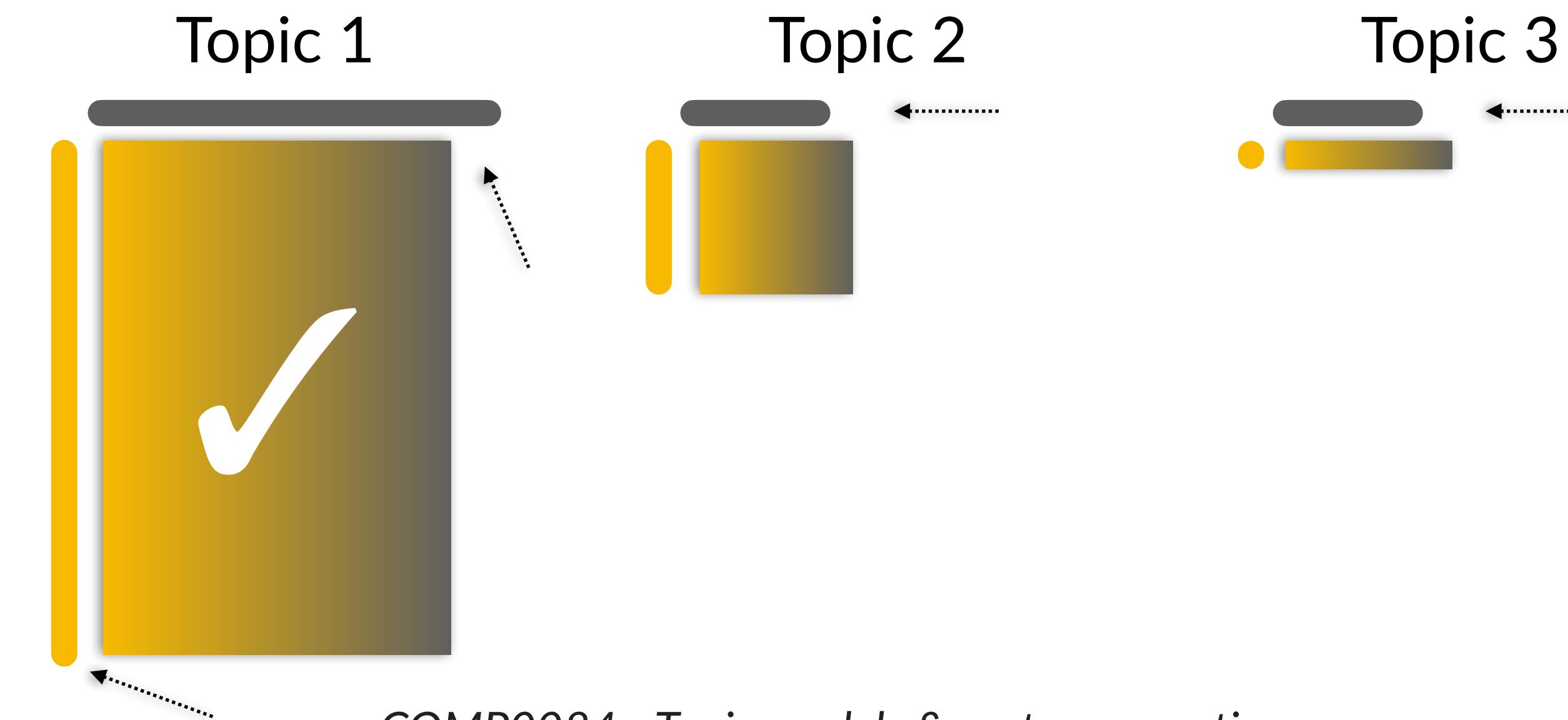
LDA – Gibbs sampling, toy example

- ...
- What are the revised topic proportions in document j ?
- How much does each topic “like” the word Brexit?
- Sample from the revised conditional distribution $p(z_{ji} = k | \mathbf{z}_{j,-i}, \mathbf{W}, \alpha, \eta)$
- **Assign the sampled topic to the word “Brexit” and update counts**
- ...

document j

z_{ji}	1	2	3	1	1
w_{ji}	Brexit	deficit	Europe	market	single

words / topics	1	2	3
Brexit	101	30	1
deficit	10	60	0
Europe	95	5	2
market	50	70	5
single	50	15	90
...



LDA online demo

mimno.infosci.cornell.edu/jsLDA/jslida.html

- ▶ **It depends on what the topics are for!**
- ▶ If they are generated for an end task with a measure-able performance, then we it makes sense to use this metric, i.e. the **performance of the end task** as a proxy for the value of the topic (Note: LDA tends to underperform in such settings)
- ▶ Compute the **probability of generating held-out documents** (*the higher the better*)
- ▶ **Word intrusion:** Show words from topics to human judges (*crowdsourcing*) with out-of-topic words inserted (intruders). How often can they identify the word that does not belong?

- ▶ We've seen that documents can be represented as vectors of word frequencies
- ▶ **Words** can also be represented as multi-dimensional **vectors**
- ▶ **Property** to exploit: words that occur in similar contexts (co-occur) tend to have similar meanings

“You shall know a word by the company it keeps”
John Rupert (J. R.) Firth (1957)

- ▶ We've seen that documents can be represented as vectors of word frequencies
- ▶ **Words** can also be represented as multi-dimensional **vectors**
- ▶ **Property** to exploit: words that occur in similar contexts (co-occur) tend to have similar meanings

“You shall know a word by the company it keeps”

John Rupert (J. R.) Firth (1957)

- My new W is much thinner than my previous one.
- I prefer to work from remote locations using a W.
- This old W has less RAM than my new smartphone.
- With a 15-inch display, it's not a W anymore!

- ▶ Property to exploit: words that occur in similar contexts (co-occur) tend to have similar meanings
 - My new W is much thinner than my previous one.
 - I prefer to work from remote locations using a W.
 - This old W has less RAM than my new smartphone.
 - With a 15-inch display, it's not a W anymore!
- ▶ Co-occurs with: “my”, “thinner”, “remote”, “smartphone”, “RAM”, “display”
- ▶ Occurs after: “my”, “a”, “new”, “old”, “display”
- ▶ Occurs before: “has”, “RAM”, “thinner”
- ▶ W = ???

- ▶ Property to exploit: words that occur in similar contexts (co-occur) tend to have similar meanings
 - My new W is much thinner than my previous one.
 - I prefer to work from remote locations using a W.
 - This old W has less RAM than my new smartphone.
 - With a 15-inch display, it's not a W anymore!
- ▶ Co-occurs with: “my”, “thinner”, “remote”, “smartphone”, “RAM”, “display”
- ▶ Occurs after: “my”, “a”, “new”, “old”, “display”
- ▶ Occurs before: “has”, “RAM”, “thinner”
- ▶ W = laptop / notebook / tablet

Words as vectors

- ▶ Generate a **word-word** matrix
 - a.k.a. **word-context** or **word co-occurrence** matrix
 - Note: words can be “terms” in practice

Words as vectors

- ▶ Generate a **word-word** matrix
 - a.k.a. **word-context** or **word co-occurrence** matrix
 - Note: words can be “terms” in practice
- ▶ If the size of our vocabulary is V , then the size of this matrix is commonly $V \times V$
- ▶ Each **cell** of the matrix reports a count of how many times two terms **co-occur** within a predefined context

Words as vectors

- ▶ Generate a **word-word** matrix
 - a.k.a. **word-context** or **word co-occurrence** matrix
 - Note: words can be “terms” in practice
- ▶ If the size of our vocabulary is V , then the size of this matrix is commonly $V \times V$
- ▶ Each **cell** of the matrix reports a count of how many times two terms **co-occur** within a predefined context
- ▶ Possible **contexts**: entire document, a paragraph in a document, a sentence, a number of terms (window, commonly ± 4 words)

... more succinct definition of **computer** science is the study...
... analysis and study of **algorithms**, discipline of computer science...
... the arrival of Japanese **mandarin** oranges signalled the real...
... of pomelo and mandarin, **orange** has genes from both...

Words as vectors

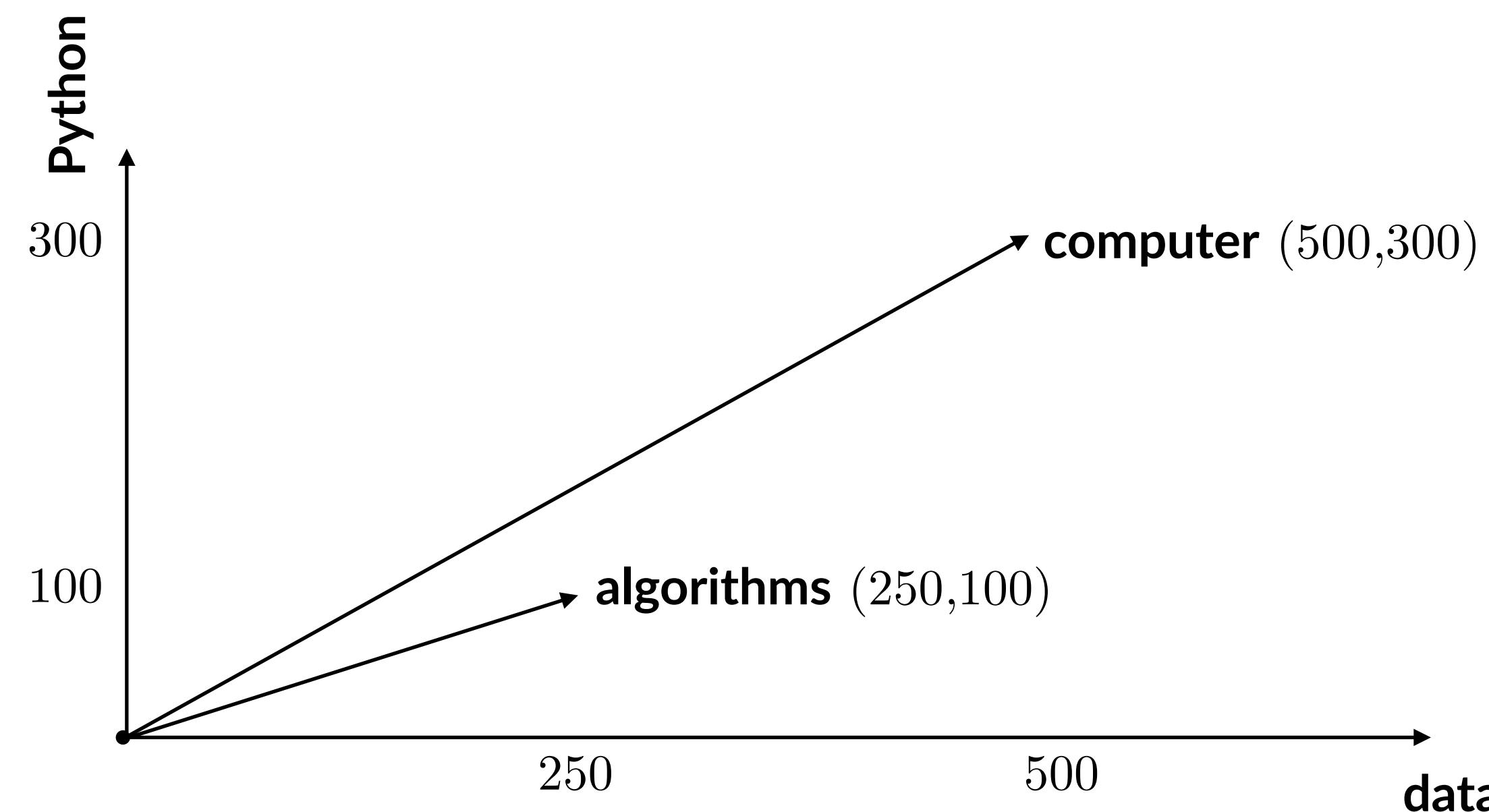
... more succinct definition of **computer** science is the study...
... analysis and study of **algorithms**, discipline of computer science...
... the arrival of Japanese **mandarin** oranges signalled the real...
... of pomelo and mandarin, **orange** has genes from both...

word-word (word co-occurrence) matrix

	...	data	...	fruit	...	Python	...
...
algorithms	...	100	...	2	...	250	...
...
computer	...	300	...	5	...	200	...
...
mandarin	...	1	...	300	...	0	...
...
orange	...	1	...	256	...	10	...
...

Words as vectors

	...	data	...	fruit	...	Python	...
...
algorithms	...	100	...	2	...	250	...
...
computer	...	300	...	5	...	200	...
...
mandarin	...	1	...	300	...	0	...
...
orange	...	1	...	256	...	10	...
...



Words as large, sparse vectors

- ▶ Recap: Word-context matrix of size $V \times V$ where V is the size of the vocabulary
- ▶ Large matrix as V is often very large ($> 100,000$ terms)
- ▶ Sparse matrix as many entries will be 0 (not all words co-occur in all contexts)

Words as large, sparse vectors

- ▶ Recap: Word-context matrix of size $V \times V$ where V is the size of the vocabulary
- ▶ Large matrix as V is often very large ($> 100,000$ terms)
- ▶ Sparse matrix as many entries will be 0 (not all words co-occur in all contexts)
- ▶ Small context window: a more **syntactic** representation
- ▶ Longer context window: a more **semantic** representation

Measuring word association – Pointwise Mutual Information (PMI)

- ▶ Raw word counts are not the best measure for word association – skewed towards frequent/infrequent words, non discriminative
- ▶ **Pointwise Mutual Information (PMI)** is a measure of how often two events co-occur, compared to what we would expect if these events were independent

Measuring word association – Pointwise Mutual Information (PMI)

- ▶ Raw word counts are not the best measure for word association – skewed towards frequent/infrequent words, non discriminative
- ▶ **Pointwise Mutual Information (PMI)** is a measure of how often two events co-occur, compared to what we would expect if these events were independent
- ▶ Centre (target) word w_i , context word c_j

$$\text{PMI}(w_i, c_j) = \log_2 \frac{p(w_i, c_j)}{p(w_i) \cdot p(c_j)}$$

- ▶ **Numerator:** How often we have seen these words together
- ▶ **Denominator:** How often we expect the words to co-occur, assuming they are independent
- ▶ **PMI:** how much more w_i, c_j co-occur than expected by chance

Positive Pointwise Mutual Information (PPMI)

- ▶ PMI ranges in $(-\infty, +\infty)$
- ▶ Negative PMI values are harder to interpret and evaluate
 - “relatedness” is easier to evaluate as opposed to “unrelatedness”
- ▶ Force positivity – Positive PMI (PPMI)

$$\text{PPMI}(w_i, c_j) = \max \left(\log_2 \frac{p(w_i, c_j)}{p(w_i) \cdot p(c_j)}, 0 \right)$$

Computing PPMI

Assume a word-context matrix \mathbf{A} of size $V \times C$; generalisation of the word-word matrix, where the C contexts may not be identical to the V target words. Let's generate a PPMI matrix from that.

$$\text{PPMI}(w_i, c_j) = \max \left(\log_2 \frac{p(w_i, c_j)}{p(w_i) \cdot p(c_j)}, 0 \right)$$

Computing PPMI

Assume a word-context matrix \mathbf{A} of size $V \times C$; generalisation of the word-word matrix, where the C contexts may not be identical to the V target words. Let's generate a PPMI matrix from that.

$$\text{PPMI}(w_i, c_j) = \max \left(\log_2 \frac{p(w_i, c_j)}{p(w_i) \cdot p(c_j)}, 0 \right)$$

$$p(w_i, c_j) = \frac{n_{ij}}{\sum_{i=1}^V \left(\sum_{j=1}^C n_{ij} \right)}$$

target word w_i co-occurs with context word c_j divided by the total count of word occurrences in the corpus

$$p(w_i) = \frac{\sum_{j=1}^C n_{ij}}{\sum_{i=1}^V \left(\sum_{j=1}^C n_{ij} \right)}$$

target word w_i appears in the corpus (sum of row i of \mathbf{A}) divided by...

$$p(c_j) = \frac{\sum_{i=1}^V n_{ij}}{\sum_{i=1}^V \left(\sum_{j=1}^C n_{ij} \right)}$$

context word c_j appears in the corpus (sum of col. j of \mathbf{A}) divided by...

Computing PPMI

Assume a word-context matrix \mathbf{A} of size $V \times C$; generalisation of the word-word matrix, where the C contexts may not be identical to the V target words. Let's generate a PPMI matrix from that.

$$\text{PPMI}(w_i, c_j) = \max \left(\log_2 \frac{p(w_i, c_j)}{p(w_i) \cdot p(c_j)}, 0 \right)$$

$$p(w_i, c_j) = \frac{n_{ij}}{\sum_{i=1}^V \left(\sum_{j=1}^C n_{ij} \right)}$$

target words
the total count

$$p(w_i) = \frac{\sum_{j=1}^C n_{ij}}{\sum_{i=1}^V \left(\sum_{j=1}^C n_{ij} \right)}$$

target word w_i appears in the corpus (sum of row i of \mathbf{A}) divided by...

$$p(c_j) = \frac{\sum_{i=1}^V n_{ij}}{\sum_{i=1}^V \left(\sum_{j=1}^C n_{ij} \right)}$$

context word c_j appears in the corpus (sum of col. j of \mathbf{A}) divided by...

We can use the PPMI matrix to measure how (semantically) similar different words are. We need a **similarity metric** for that.

Measuring word similarity – Cosine

- Dot product between word vectors \mathbf{w}, \mathbf{v} : $\mathbf{w}^\top \mathbf{v} = \sum_{i=1}^N w_i \cdot v_i$

Measuring word similarity – Cosine

- ▶ Dot product between word vectors \mathbf{w} , \mathbf{v} : $\mathbf{w}^\top \mathbf{v} = \sum_{i=1}^N w_i \cdot v_i$
- ▶ Greater values for longer vectors and for frequent words

Measuring word similarity – Cosine

- ▶ Dot product between word vectors \mathbf{w}, \mathbf{v} : $\mathbf{w}^\top \mathbf{v} = \sum_{i=1}^N w_i \cdot v_i$
- ▶ Greater values for longer vectors and for frequent words
- ▶ Normalise it by dividing with the length of the vectors! Leads to cosine similarity, i.e. the cosine of the angle (ϕ) between the two vectors

$$\text{cosine-sim}(\mathbf{w}, \mathbf{v}) = \frac{\sum_{i=1}^N w_i \cdot v_i}{\sqrt{\sum_{i=1}^N w_i^2} \cdot \sqrt{\sum_{i=1}^N v_i^2}} = \frac{\mathbf{w}^\top \mathbf{v}}{\|\mathbf{w}\|_2 \|\mathbf{v}\|_2} = \cos \phi$$

Measuring word similarity – Cosine

- ▶ Dot product between word vectors \mathbf{w}, \mathbf{v} : $\mathbf{w}^\top \mathbf{v} = \sum_{i=1}^N w_i \cdot v_i$
- ▶ Greater values for longer vectors and for frequent words
- ▶ Normalise it by dividing with the length of the vectors! Leads to cosine similarity, i.e. the cosine of the angle (ϕ) between the two vectors

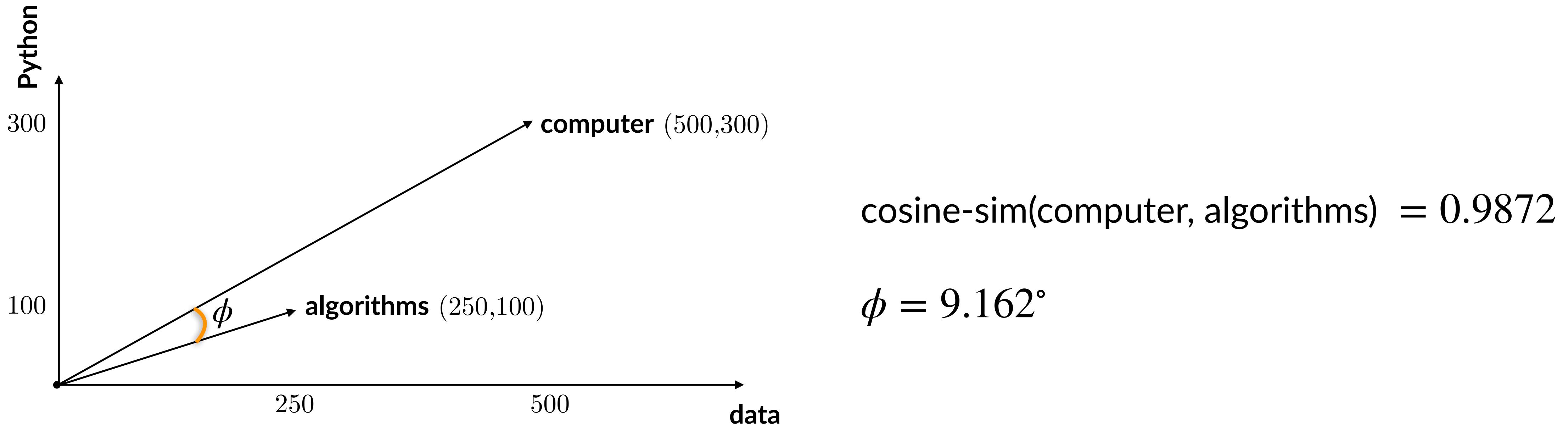
$$\text{cosine-sim}(\mathbf{w}, \mathbf{v}) = \frac{\sum_{i=1}^N w_i \cdot v_i}{\sqrt{\sum_{i=1}^N w_i^2} \cdot \sqrt{\sum_{i=1}^N v_i^2}} = \frac{\mathbf{w}^\top \mathbf{v}}{\|\mathbf{w}\|_2 \|\mathbf{v}\|_2} = \cos \phi$$

- ▶ Since \mathbf{w} and $\mathbf{v} > 0$, $\text{cosine-sim}(\mathbf{w}, \mathbf{v})$ ranges from $[0, 1]$
 - $\text{cosine-sim}(\mathbf{w}, \mathbf{v}) = 0$ means that $\phi = 90^\circ$
 - $\text{cosine-sim}(\mathbf{w}, \mathbf{v}) = 1$ means that $\phi = 0^\circ$

Measuring word similarity – Cosine

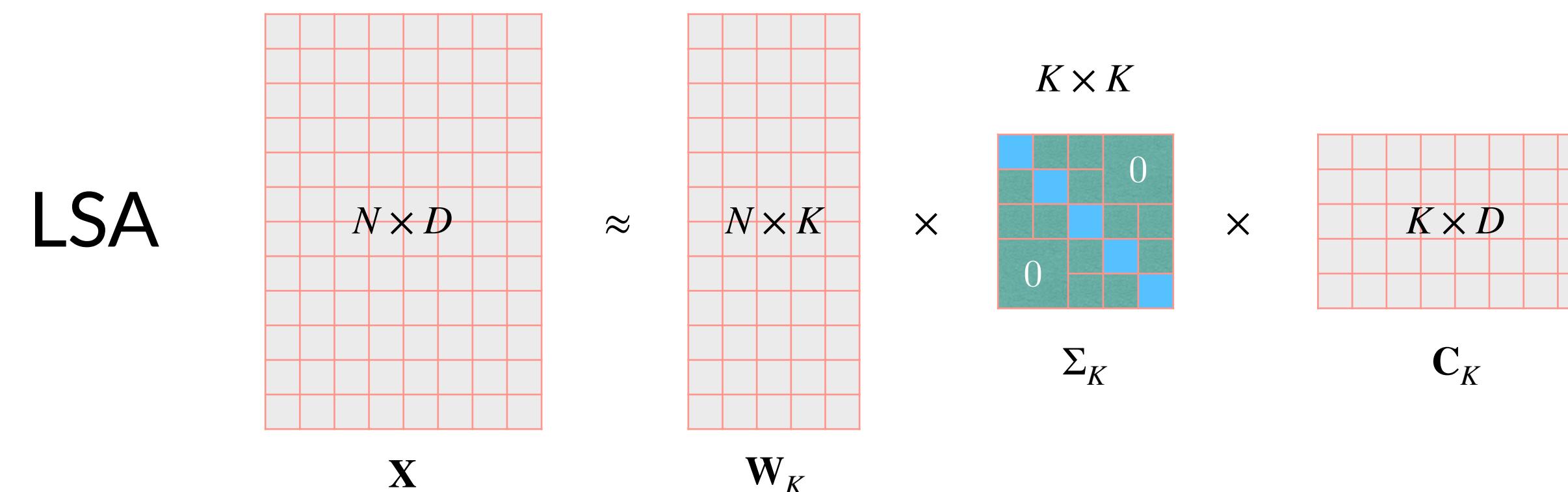
$$\text{cosine-sim}(\mathbf{w}, \mathbf{v}) = \frac{\sum_{i=1}^N w_i \cdot v_i}{\sqrt{\sum_{i=1}^N w_i^2} \cdot \sqrt{\sum_{i=1}^N v_i^2}} = \frac{\mathbf{w}^\top \mathbf{v}}{\|\mathbf{w}\|_2 \|\mathbf{v}\|_2} = \cos \phi$$

- ▶ Since \mathbf{w} and $\mathbf{v} > 0$, cosine-sim (\mathbf{w}, \mathbf{v}) ranges from $[0,1]$
 - cosine-sim (\mathbf{w}, \mathbf{v}) = 0 means that $\phi = 90^\circ$
 - cosine-sim (\mathbf{w}, \mathbf{v}) = 1 means that $\phi = 0^\circ$



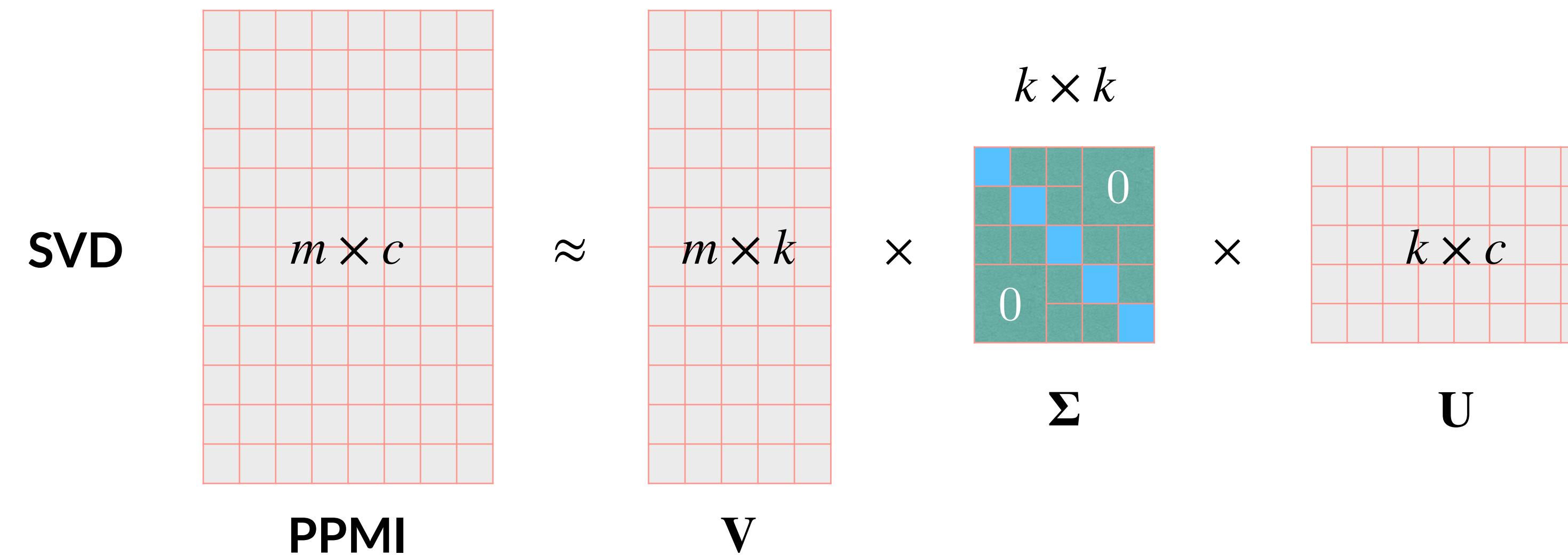
From sparse to dense word vectors

- ▶ Previously shown word representations: **long** (equal to size of the vocabulary V) and **sparse** (many 0's)
- ▶ **Short and dense** representations have advantages
 - easier to use as features in statistical learning methods
 - capture synonymy better
 - generalise better

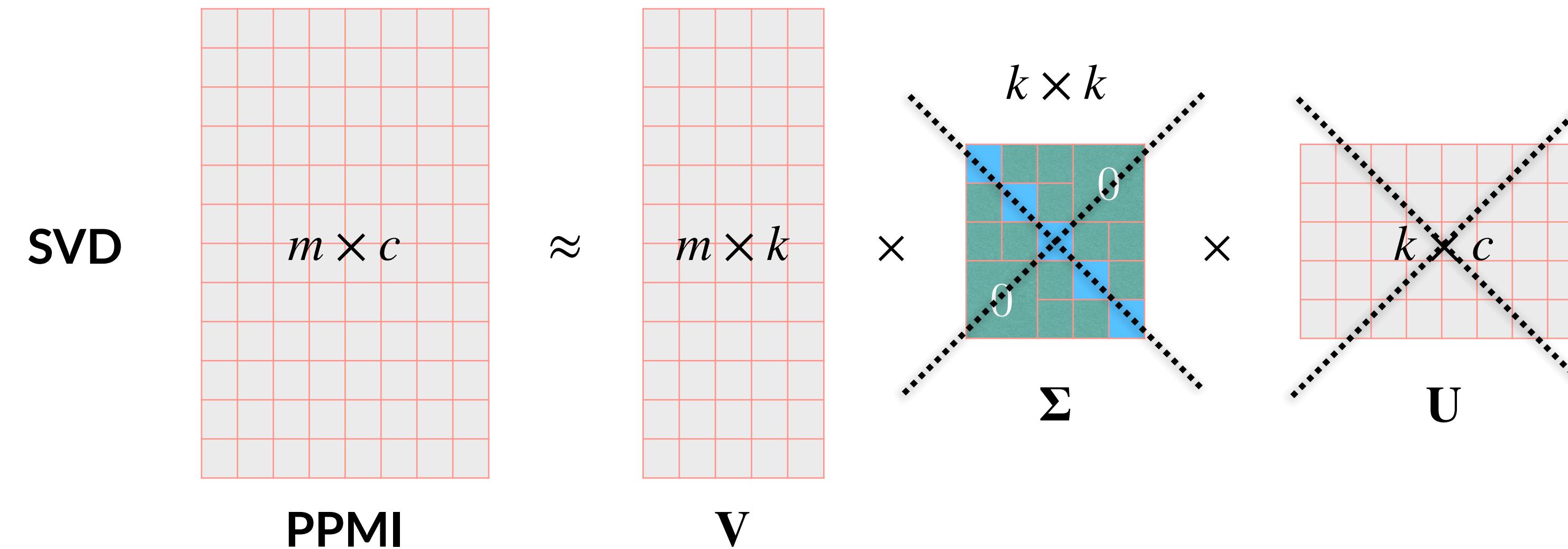


- ▶ Recall Latent Semantic Analysis (LSA), i.e. SVD on the word-document matrix, \mathbf{X} . What if we perform SVD on a word co-occurrence or a PPMI matrix?

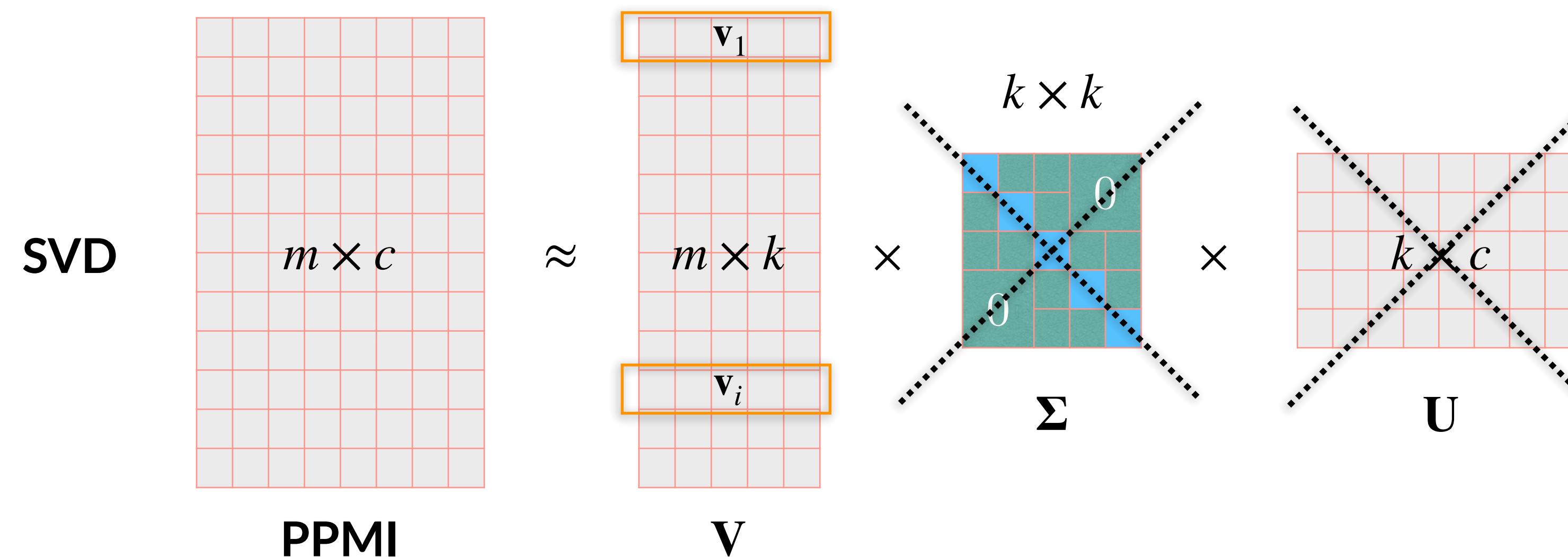
SVD on the PPMI word-context matrix



SVD on the PPMI word-context matrix



SVD on the PPMI word-context matrix



- ▶ v_i : k -dimensional vector that represents word i in our vocabulary
 - also known as a **word embedding**
 - commonly, $k = 128$ to 1024 , i.e. v_i is short and dense
- ▶ **Downside:** SVD has a significant computational cost

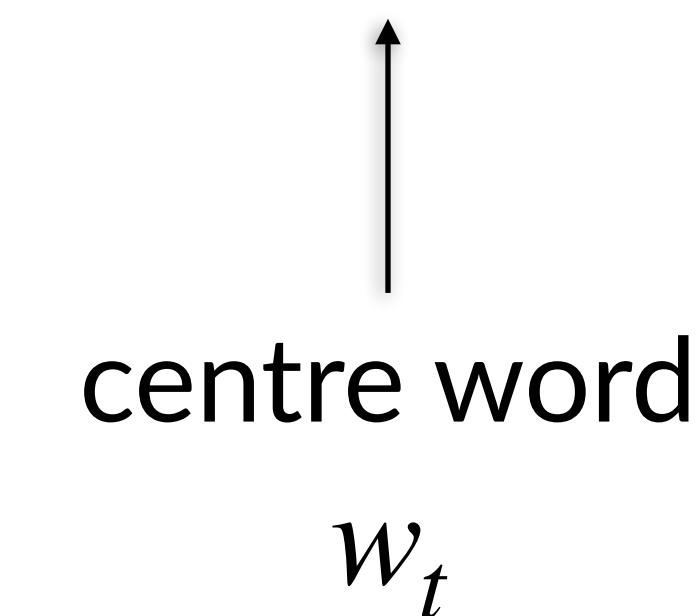
Word embeddings from prediction

- ▶ Same intuition, *different* approach
 - words with similar meanings will co-occur
 - instead of counting co-occurrences, **predict** them
- ▶ Popular example: **word2vec** — title of the software library, but there is a small family of methods behind it
- ▶ Algorithms
 - skip-gram: Predict the context (surrounding) words based on a centre word
 - CBOW (continuous bag-of-words): Predict a centre word based on the context words
- ▶ Training methods
 - Hierarchical softmax
 - Negative sampling
 - Naïve softmax

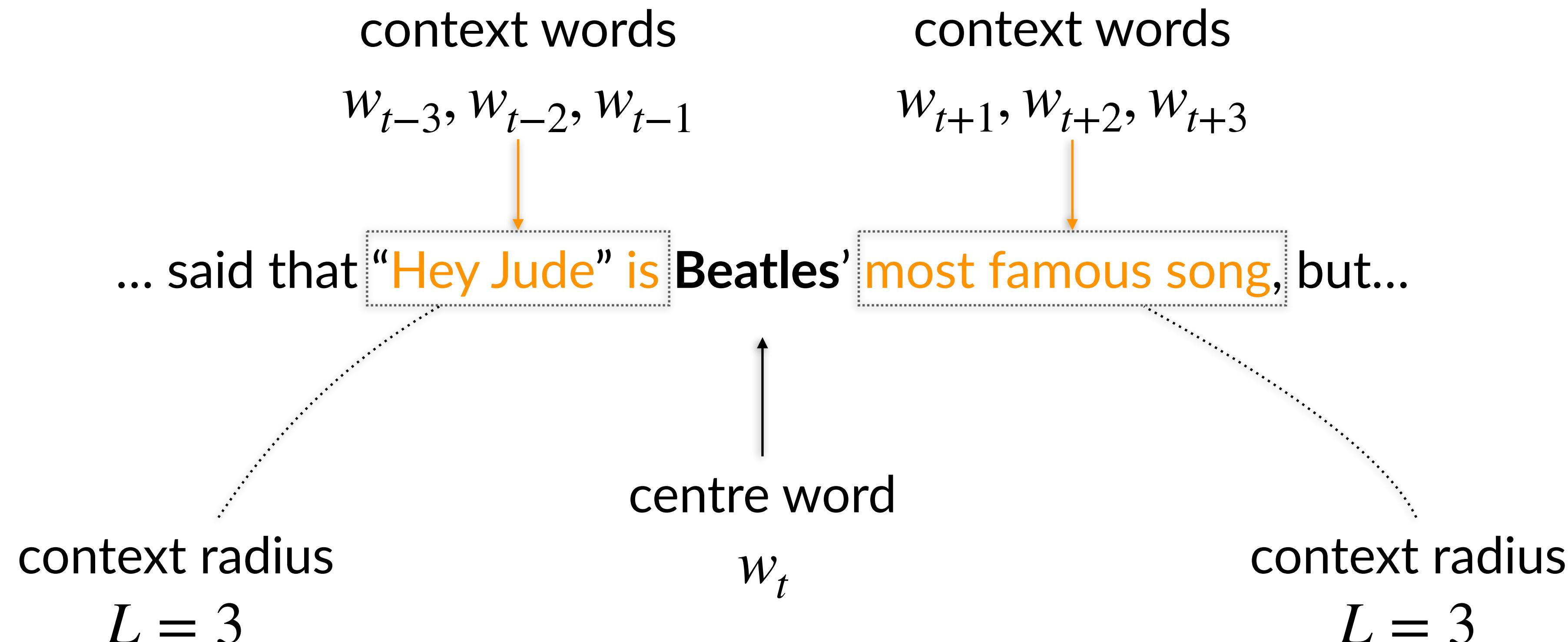
... said that “Hey Jude” is Beatles’ most famous song, but...

word2vec – skip-gram

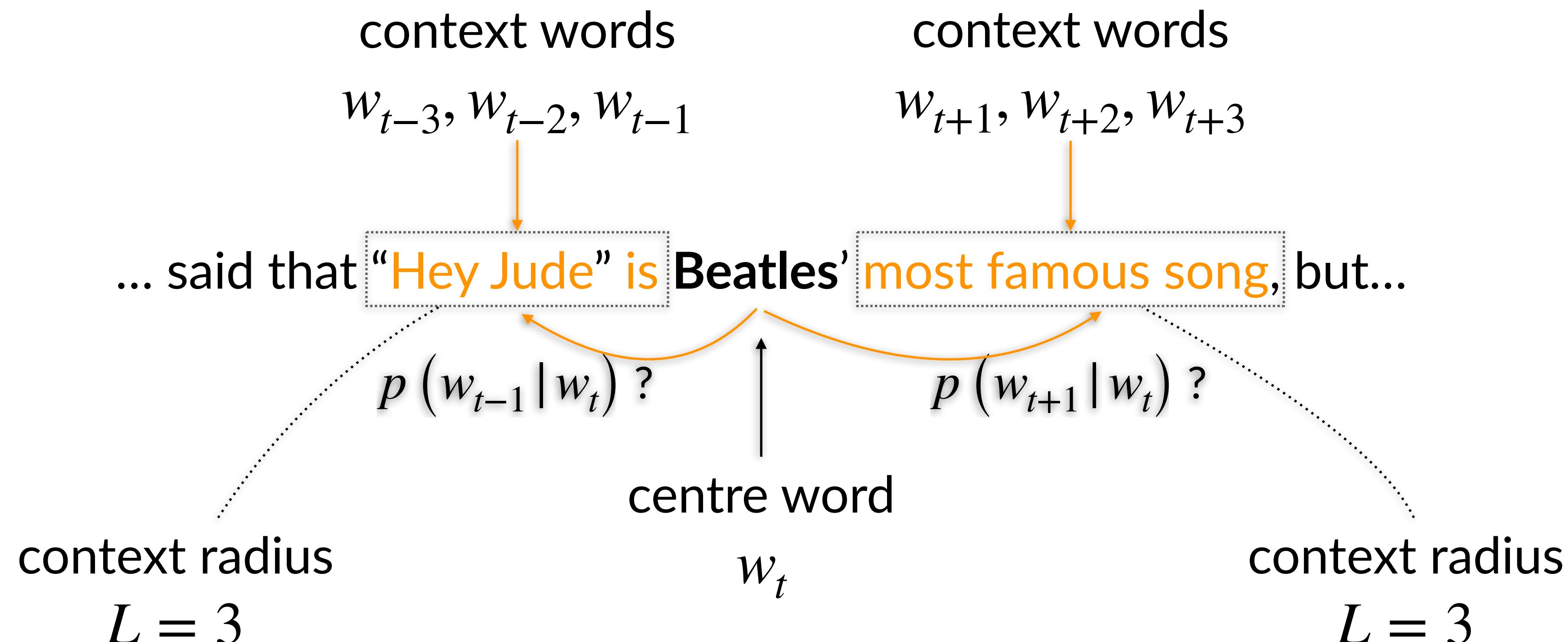
... said that “Hey Jude” is **Beatles’** most famous song, but...



word2vec – skip-gram



word2vec – skip-gram



skip-gram – Simplified objective function

For each word position t out of T , predict the context words using a fixed radius L (or a symmetric window $2L$)

Objective: Maximise the probability of any context word given the current centre word (the position of surrounding words does not matter)

$$\max \prod_{t=1}^T \prod_{i=-L, i \neq 0}^L p(w_{t+i} | w_t)$$

skip-gram – Simplified objective function

For each word position t out of T , predict the context words using a fixed radius L (or a symmetric window $2L$)

Objective: Maximise the probability of any context word given the current centre word (the position of surrounding words does not matter)

$$\max \prod_{t=1}^T \prod_{i=-L, i \neq 0}^L p(w_{t+i} | w_t)$$

Prefer to minimise things, and sums over products

Minimise the mean (across all T samples) negative log likelihood:

$$\min \frac{1}{T} \left(- \sum_{t=1}^T \sum_{i=-L, i \neq 0}^L \log(p(w_{t+i} | w_t)) \right)$$

skip-gram – Simplified objective function

For each word position t out of T , predict the context words using a fixed radius L (or a symmetric window $2L$)

Objective: Maximise the probability of any context word given the current centre word

$$\min \frac{1}{T} \left(- \sum_{t=1}^T \sum_{i=-L, i \neq 0}^L \log \left(p(w_{t+i} | w_t) \right) \right)$$

- Assume that each centre word (t) has a k -dimensional vector representation \mathbf{v}_c ; all the vectors of the m centre words are held in an $k \times m$ matrix \mathbf{V}
- Assume that each context word has a k -dimensional vector representation \mathbf{u}_x ; all the vectors of the m context words are held in an $k \times m$ matrix \mathbf{U}
- Thus, the model parameters ($2mk$) are now $\mathbf{Q} = [\mathbf{V} \ \mathbf{U}]$

$$\min_{\mathbf{Q}} \frac{1}{T} \left(- \sum_{t=1}^T \sum_{i=-L, i \neq 0}^L \log \left(p(w_{t+i} | w_t; \mathbf{Q}) \right) \right)$$

skip-gram – Simplified objective function

$$\min_{\mathbf{Q}} \frac{1}{T} \left(- \sum_{t=1}^T \sum_{i=-L, i \neq 0}^L \log \left(p(w_{t+i} | w_t; \mathbf{Q}) \right) \right)$$

We need an estimate of the probability $p(w_{t+1} | w_t)$

Use a (*bad*) measure of similarity (dot product) and normalise it using a common approach in neural networks, the **softmax** function that squashes vector elements to a $[0,1]$ range

Assuming a vocabulary of m words, for a centre word c (\mathbf{v}_c) and a context word x (\mathbf{u}_x)

$$p(x | c) = \frac{\exp(\mathbf{u}_x^\top \mathbf{v}_c)}{\sum_{w=1}^m \exp(\mathbf{u}_w^\top \mathbf{v}_c)}$$

skip-gram – In practice

$$w_t = [0 \ 0 \ \dots \ 1 \ \dots \ 0]^T$$

centre word as an one-hot vector

$$v_c = \mathbf{V} \cdot w_t$$

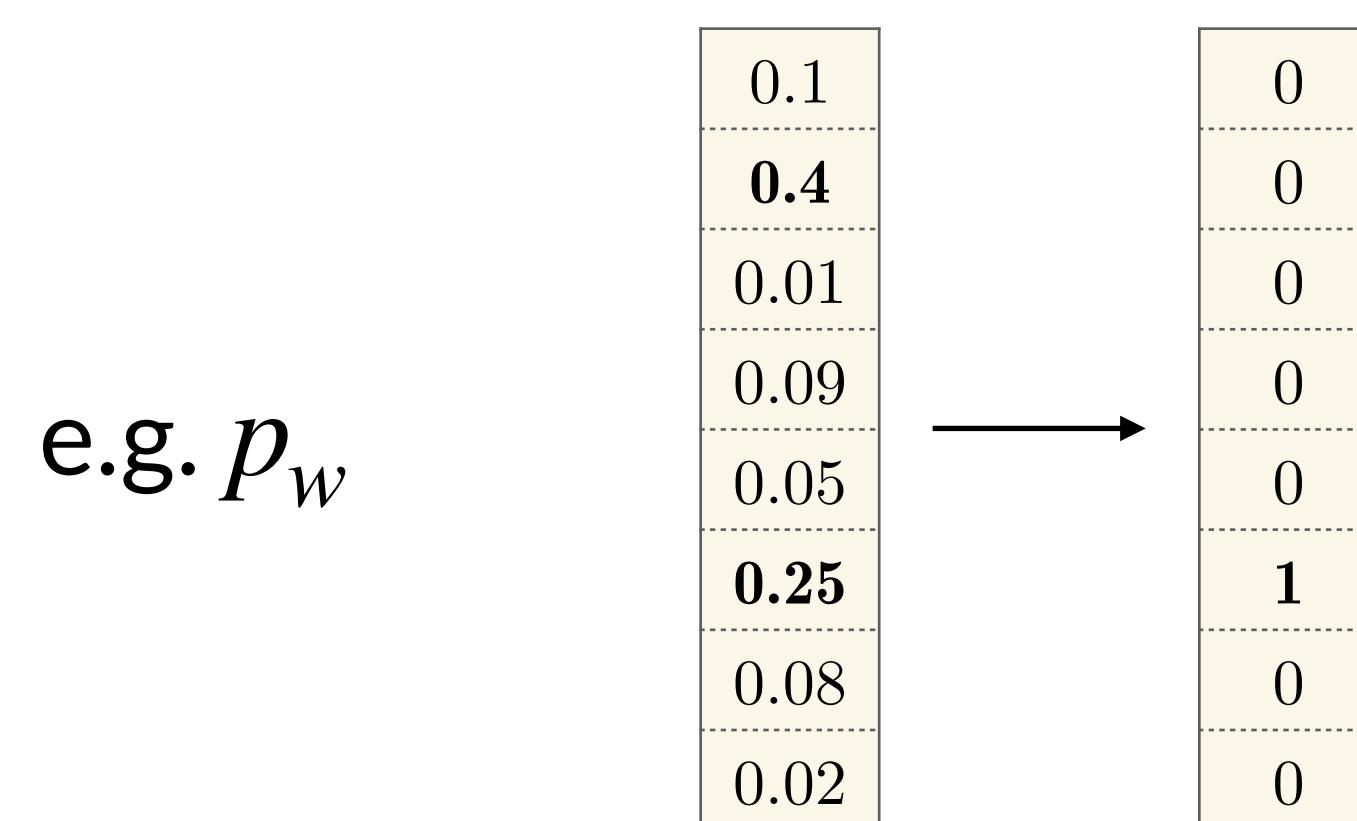
get its vector representation (embedding) from the matrix of centre word embeddings

$$o = \mathbf{U}^T \cdot v_c$$

dot product with all context word vectors m (voc. size) $\times 1$

$$p_{w_i} = \text{softmax}(o)_i$$

compute the softmax of this vector – this is the probability of word i , but we shall focus on the $2L$ context words



But we also know the correct answer!
In this case, we need to improve our
embeddings (\mathbf{V} and \mathbf{U}).

In neural nets: do error back-propagation.

skip-gram – Negative sampling

Naïve / inefficient way for parameter inference

$$J(\mathbf{Q}) = -\frac{1}{T} \sum_{t=1}^T \sum_{i=-L, i \neq 0}^L \log \left(p(w_{t+i} | w_t; \mathbf{Q}) \right)$$

Gradient descent: $\mathbf{Q}_{p+1} = \mathbf{Q}_p - \gamma \nabla_{\mathbf{Q}} J(\mathbf{Q}_p)$

Too slow and computationally expensive. Recall, the denominator is too expensive to compute (for large vocabularies; m)

$$p(x | c) = \frac{\exp (\mathbf{u}_x^\top \mathbf{v}_c)}{\sum_{w=1}^m \exp (\mathbf{u}_w^\top \mathbf{v}_c)}$$

skip-gram – Negative sampling

Naïve / inefficient way for parameter inference

$$J(\mathbf{Q}) = -\frac{1}{T} \sum_{t=1}^T \sum_{i=-L, i \neq 0}^L \log \left(p(w_{t+i} | w_t; \mathbf{Q}) \right)$$

Gradient descent: $\mathbf{Q}_{p+1} = \mathbf{Q}_p - \gamma \nabla_{\mathbf{Q}} J(\mathbf{Q}_p)$

Too slow and computationally expensive. Recall, the denominator is too expensive to compute (for large vocabularies; m)

$$p(x | c) = \frac{\exp (\mathbf{u}_x^\top \mathbf{v}_c)}{\sum_{w=1}^m \exp (\mathbf{u}_w^\top \mathbf{v}_c)}$$

Negative sampling: For each context word, sample non-neighbouring words as “negative” samples

New objective: High dot product with context words and low dot product with “negative” samples

skip-gram – Negative sampling

Naïve / inefficient way for parameter inference

$$J(\mathbf{Q}) = -\frac{1}{T} \sum_{t=1}^T \sum_{i=-L, i \neq 0}^L \log \left(p(w_{t+i} | w_t; \mathbf{Q}) \right)$$

Gradient descent: $\mathbf{Q}_{p+1} = \mathbf{Q}_p - \gamma \nabla_{\mathbf{Q}} J(\mathbf{Q}_p)$

Too slow and computationally expensive.

skip-gram – Negative sampling

Naïve / inefficient way for parameter inference

$$J(\mathbf{Q}) = -\frac{1}{T} \sum_{t=1}^T \sum_{i=-L, i \neq 0}^L \log \left(p(w_{t+i} | w_t; \mathbf{Q}) \right)$$

Gradient descent: $\mathbf{Q}_{p+1} = \mathbf{Q}_p - \gamma \nabla_{\mathbf{Q}} J(\mathbf{Q}_p)$

Too slow and computationally expensive.

Apply mini-batch gradient descent:

i.e. instead of going through all the data for computing $\nabla_{\mathbf{Q}} J(\mathbf{Q}_p)$

we use one or small subsets of the data (mini batches) to update the gradient

Word analogies with word embeddings

Note

Word embeddings
tend to carry the
biases or
stereotypes of the
corpora used to
train them!

$$\text{vector}(\text{'queen'}) \approx \text{vector}(\text{'king'}) - \text{vector}(\text{'man'}) + \text{vector}(\text{'woman'})$$

Word analogies with word embeddings

Note

Word embeddings
tend to carry the
biases or
stereotypes of the
corpora used to
train them!

b

a

a_p

b_p

$$\text{vector}(\text{'queen'}) \approx \text{vector}(\text{'king'}) - \text{vector}(\text{'man'}) + \text{vector}(\text{'woman'})$$

Word analogies with word embeddings

Note

Word embeddings tend to carry the biases or stereotypes of the corpora used to train them!

b

a

a_p

b_p

$$\text{vector}(\text{'queen'}) \approx \text{vector}(\text{'king'}) - \text{vector}(\text{'man'}) + \text{vector}(\text{'woman'})$$

Do this using the cosine similarity between word embeddings

$$b = \arg \max_{b \in V} \left(\cos \left(\mathbf{v}_b, \mathbf{v}_a - \mathbf{v}_{a_p} + \mathbf{v}_{b_p} \right) \right)$$

Compute the cosine similarity between the composite embedding $(\mathbf{v}_a - \mathbf{v}_{a_p} + \mathbf{v}_{b_p})$ and each other embedding in our vocabulary, and expect that $\mathbf{v}_b = \text{vector}(\text{'queen'})$ will have the greatest one.

Word analogies with word embeddings

Note

Word embeddings tend to carry the biases or stereotypes of the corpora used to train them!

b

a

a_p

b_p

$$\text{vector}(\text{'queen'}) \approx \text{vector}(\text{'king'}) - \text{vector}(\text{'man'}) + \text{vector}(\text{'woman'})$$

Do this using the cosine similarity between word embeddings

$$b = \arg \max_{b \in V} \left(\cos \left(\mathbf{v}_b, \mathbf{v}_a - \mathbf{v}_{a_p} + \mathbf{v}_{b_p} \right) \right)$$

Compute the cosine similarity between the composite embedding $(\mathbf{v}_a - \mathbf{v}_{a_p} + \mathbf{v}_{b_p})$ and each other embedding in our vocabulary, and expect that $\mathbf{v}_b = \text{vector}(\text{'queen'})$ will have the greatest one.

This gives rise to the word analogy

a_p is for a , what b_p is for b
or man is for king, what woman is for queen

Word analogies with word embeddings

Note

Word embeddings tend to carry the biases or stereotypes of the corpora used to train them!

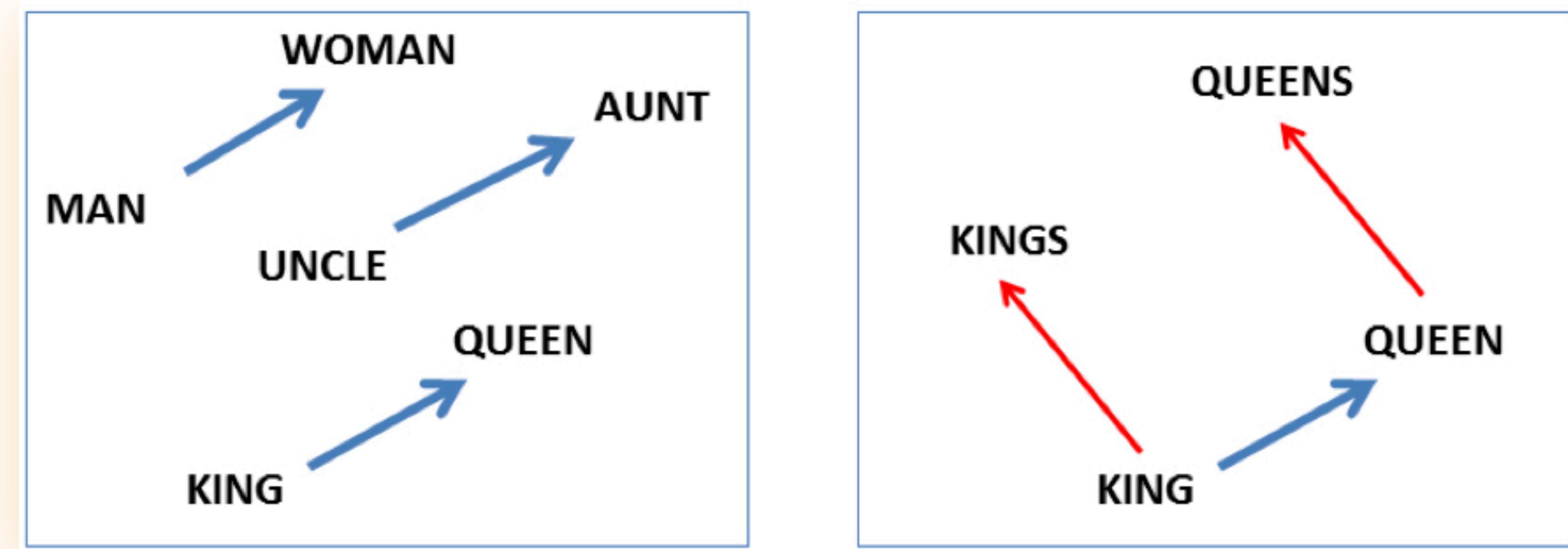
b

a

a_p

b_p

$$\text{vector}(\text{'queen'}) \approx \text{vector}(\text{'king'}) - \text{vector}(\text{'man'}) + \text{vector}(\text{'woman'})$$



This gives rise to the **word analogy**

a_p is for a , what b_p is for b
or man is for king, what woman is for queen

word2vec embeddings

- ▶ trained (*a few years back*) on 1.1 billion tweets post during 2012 to 2016, approximately geolocated in the UK
- ▶ tweets represent current trends, include informal forms of language, and are often topic-consistent
- ▶ 470,194 terms covered (size of the vocabulary)
- ▶ the dimensionality of the embedding is equal to 512
- ▶ available online at figshare.com/articles/UK_Twitter_word_embeddings_II_/5791650

Twitter word embeddings – Similarities

Top-5 most similar words using cosine similarity on word embeddings

- ▶ **Monday:** Tuesday, Thursday, Wednesday, Friday, Sunday
- ▶ **January:** February, August, October, March, June
- ▶ **red:** yellow, blue, purple, pink, green
- ▶ **we:** they, you, we've, our, us
- ▶ **espresso:** espresso, cappuccino, macchiato, latte, coffee
- ▶ **linux:** Unix, Centos, Debian, Ubuntu, Redhat
- ▶ **retweet:** rt, tweet, retweets, retweeting, rewteet
- ▶ **democracy:** democratic, dictatorship, democracies, socialism, undemocratic
- ▶ **loool:** looool, lool, looooool, loooooool, loooooooool
- ▶ **xxxx:** xxxx, xxx, xxxxxxxxx, xxxxxx, xxxxxxxx
- ▶ **enviroment:** environment, environments, env, enviro, habitats

Twitter word embeddings – Analogies

- ▶ **she** is to **her** what **he** is to ...

Twitter word embeddings – Analogies

- ▶ **she** is to **her** what **he** is to ... [**his, him, himself**]

Twitter word embeddings – Analogies

- ▶ **she** is to **her** what **he** is to ... [**his, him, himself**]
- ▶ **Rome** is to **Italy** what **London** is to ...

Twitter word embeddings – Analogies

- ▶ **she** is to **her** what **he** is to ... [**his, him, himself**]
- ▶ **Rome** is to **Italy** what **London** is to ... [**UK, Denmark, Sweden**]

Twitter word embeddings – Analogies

- ▶ **she** is to **her** what **he** is to ... [**his, him, himself**]
- ▶ **Rome** is to **Italy** what **London** is to ... [**UK, Denmark, Sweden**]
- ▶ **go** is for **went** what **do** is to... [**did, doing, happened**]

Twitter word embeddings – Analogies

- ▶ **she** is to **her** what **he** is to ... [**his, him, himself**]
- ▶ **Rome** is to **Italy** what **London** is to ... [**UK, Denmark, Sweden**]
- ▶ **go** is for **went** what **do** is to... [**did, doing, happened**]
- ▶ **big** is to **bigger** what **small** is to... [**smaller, larger, tiny**]

Twitter word embeddings – Analogies

- ▶ **she** is to **her** what **he** is to ... [**his, him, himself**]
- ▶ **Rome** is to **Italy** what **London** is to ... [**UK, Denmark, Sweden**]
- ▶ **go** is for **went** what **do** is to... [**did, doing, happened**]
- ▶ **big** is to **bigger** what **small** is to... [**smaller, larger, tiny**]
- ▶ **poet** is to **poem** what **author** is to... [**novel, excerpt, memoir**]
- ▶ **Messi** is to **football** what **Lebron** is to... [**basketball, bball, NBA**]
- ▶ **Elvis** is to **Presley** what **Aretha** is to... [**Franklin, Ruffin, Vandross**]

Twitter word embeddings – Analogies

- ▶ **she** is to **her** what **he** is to ... [**his, him, himself**]
- ▶ **Rome** is to **Italy** what **London** is to ... [**UK, Denmark, Sweden**]
- ▶ **go** is for **went** what **do** is to... [**did, doing, happened**]
- ▶ **big** is to **bigger** what **small** is to... [**smaller, larger, tiny**]
- ▶ **poet** is to **poem** what **author** is to... [**novel, excerpt, memoir**]
- ▶ **Messi** is to **football** what **Lebron** is to... [**basketball, bball, NBA**]
- ▶ **Elvis** is to **Presley** what **Aretha** is to... [**Franklin, Ruffin, Vandross**]
- ▶ **UK** is for **Brexit** what **Greece** is to... [**Grexit, Syriza, Tsipras**]
- ▶ **UK** is for **Farage** what **USA** is to... [**Trump, Farrage, Putin**]

Next lectures with me

- ▶ March 2, 5pm to 6pm, guest lecture by Dr Adam Tsakalidis
- ▶ March 23, 5pm to 6pm, guest lecture by me about some of my research