

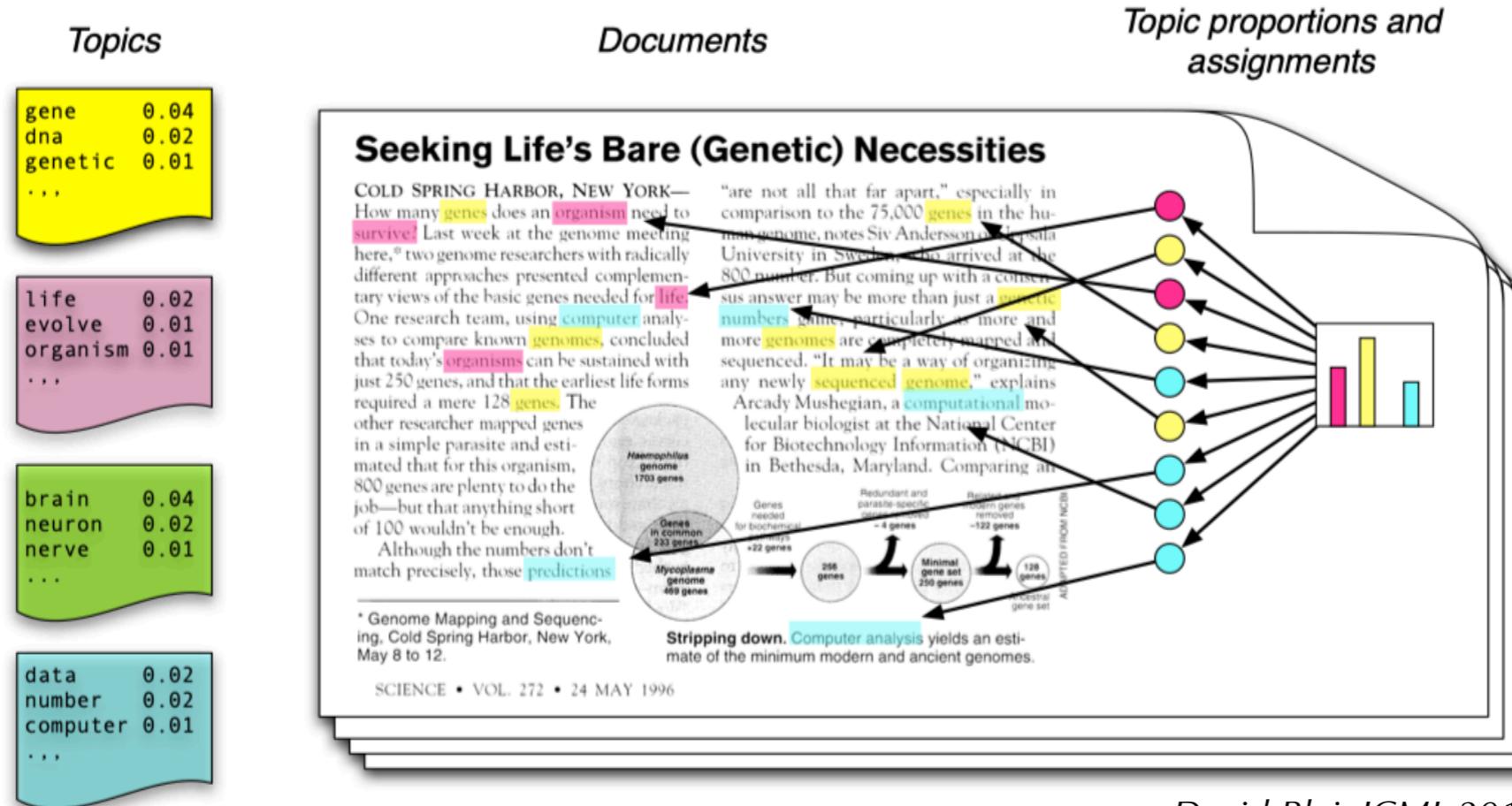
Topic Models (I)

Runzhe Yang

2019/03/20

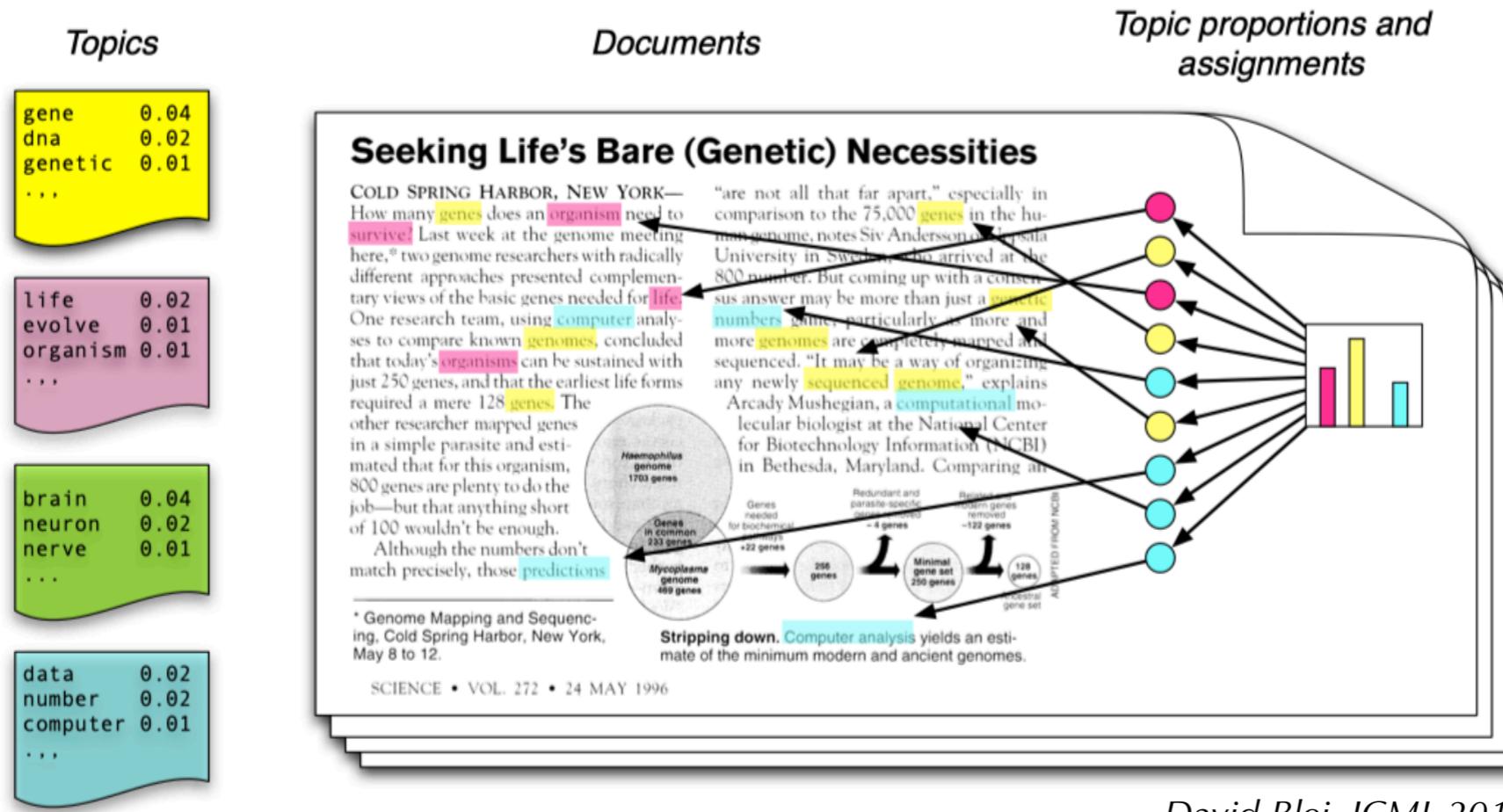
What is “topic”?

A **hidden structure** that helps determine what **words** are likely to appear in a **document**.



David Blei, ICML 2012

What is “topic”?

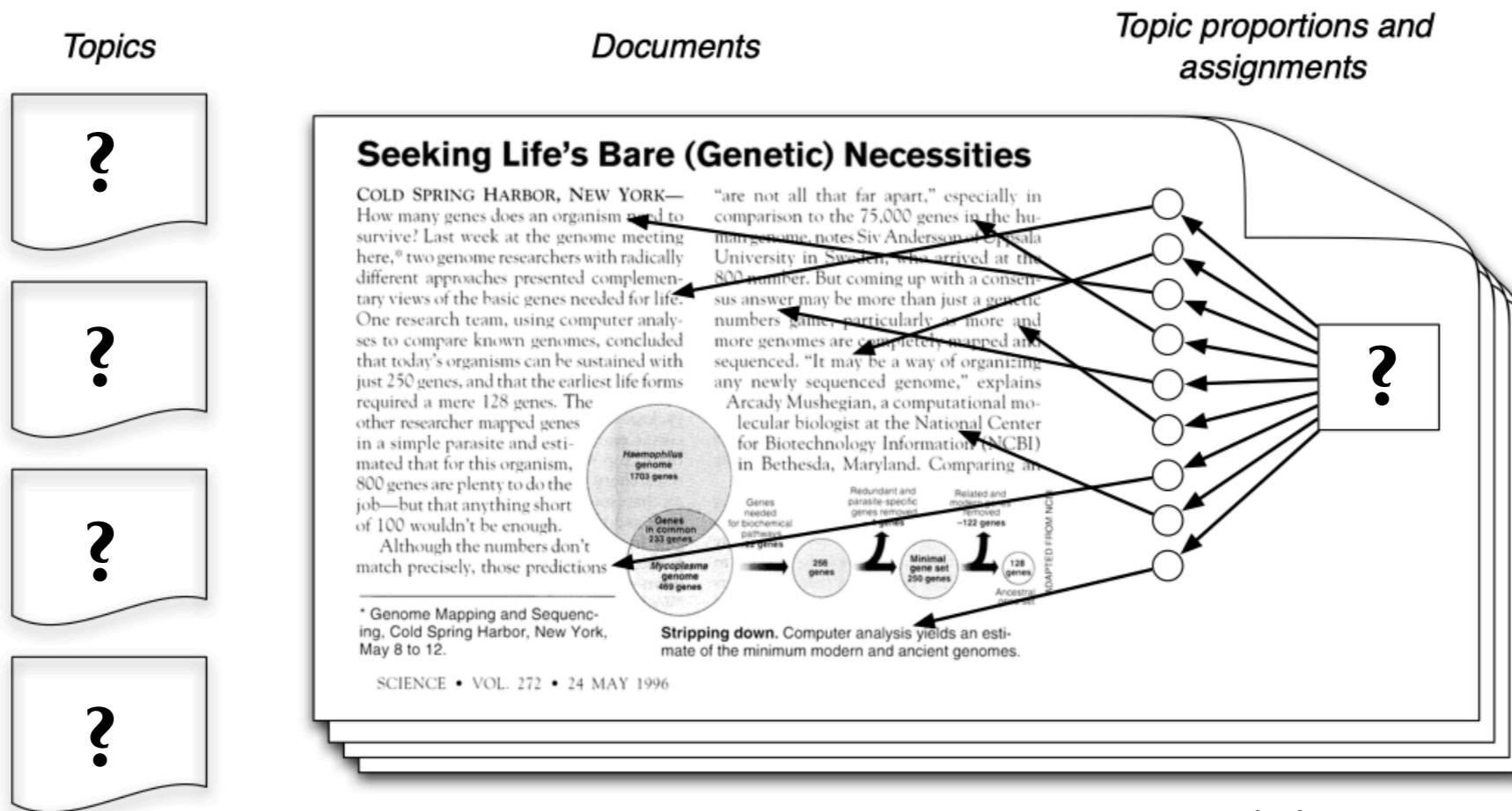


David Blei, ICML 2012

- Each **topic** is a distribution over words
- Each **document** is a mixture of corpus-wide topics
- Each **word** is drawn from one of those **topics**

What is “topic modeling”?

We only observe the **documents**.



The goal of **topic modeling** is
to infer the **hidden structures**. —

- 1) what're the topics of the corpus?
- 2) what're the topics of this document?

Outline:

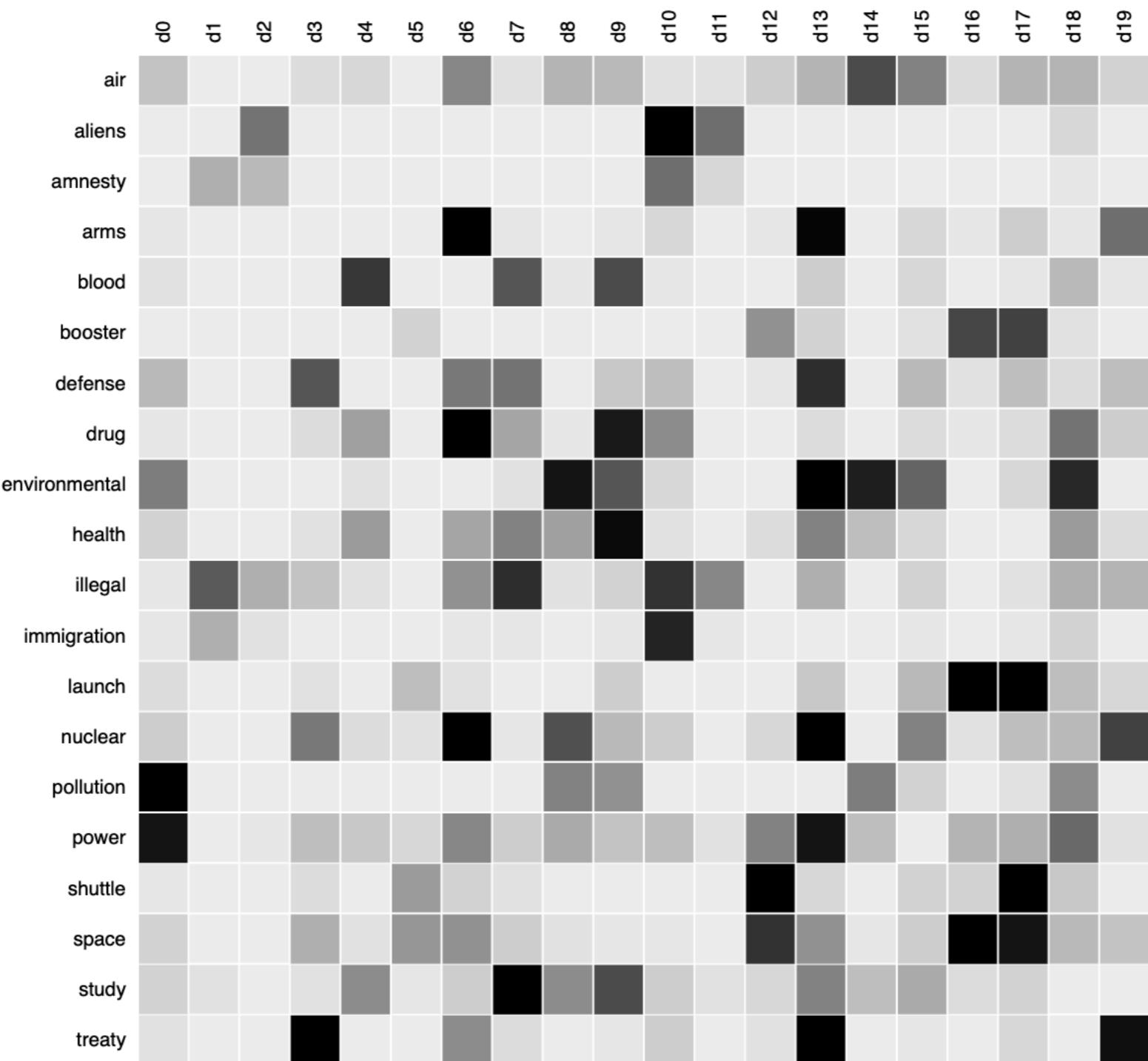
- **Classic Methods Overview**
 - Latent Semantic Indexing (LSI) - SVD
 - Probabilistic Latent Semantic Indexing (pLSI)
 - Equivalence to NMF & Autoencoder (?)
 - Latent Dirichlet Allocation (LDA)
 - Neural Topic Modeling (Gaussian prior)
 - **A comparison of SVD / NMF / LDA**
- **Community Detection & Topic Modeling**
 - ~~Stochastic Block Model (SBM)~~
 - ~~Hierarchical Stochastic Block Model (hSBM)~~
- **Beyond Probabilistic Generative Models**
 - Disynaptic Neural Network
 - Connection with Correlation Game
 - **Disynaptic Neural Topic Modelling?**

Outline:

- **Classic Methods Overview**
 - Latent Semantic Indexing (LSI) - SVD
 - Probabilistic Latent Semantic Indexing (pLSI)
 - Equivalence to NMF & Autoencoder (?)
 - Latent Dirichlet Allocation (LDA)
 - Neural Topic Modeling (Gaussian prior)
 - **A comparison of SVD / NMF / LDA**
- **Community Detection & Topic Modeling**
 - Stochastic Block Model (SBM)
 - Hierarchical Stochastic Block Model (hSBM)
- **Beyond Probabilistic Generative Models**
 - Disynaptic Neural Network
 - Correlation Game
 - **Disynaptic Neural Topic Modelling?**

Please stop and correct me
whenever you think I am wrong!

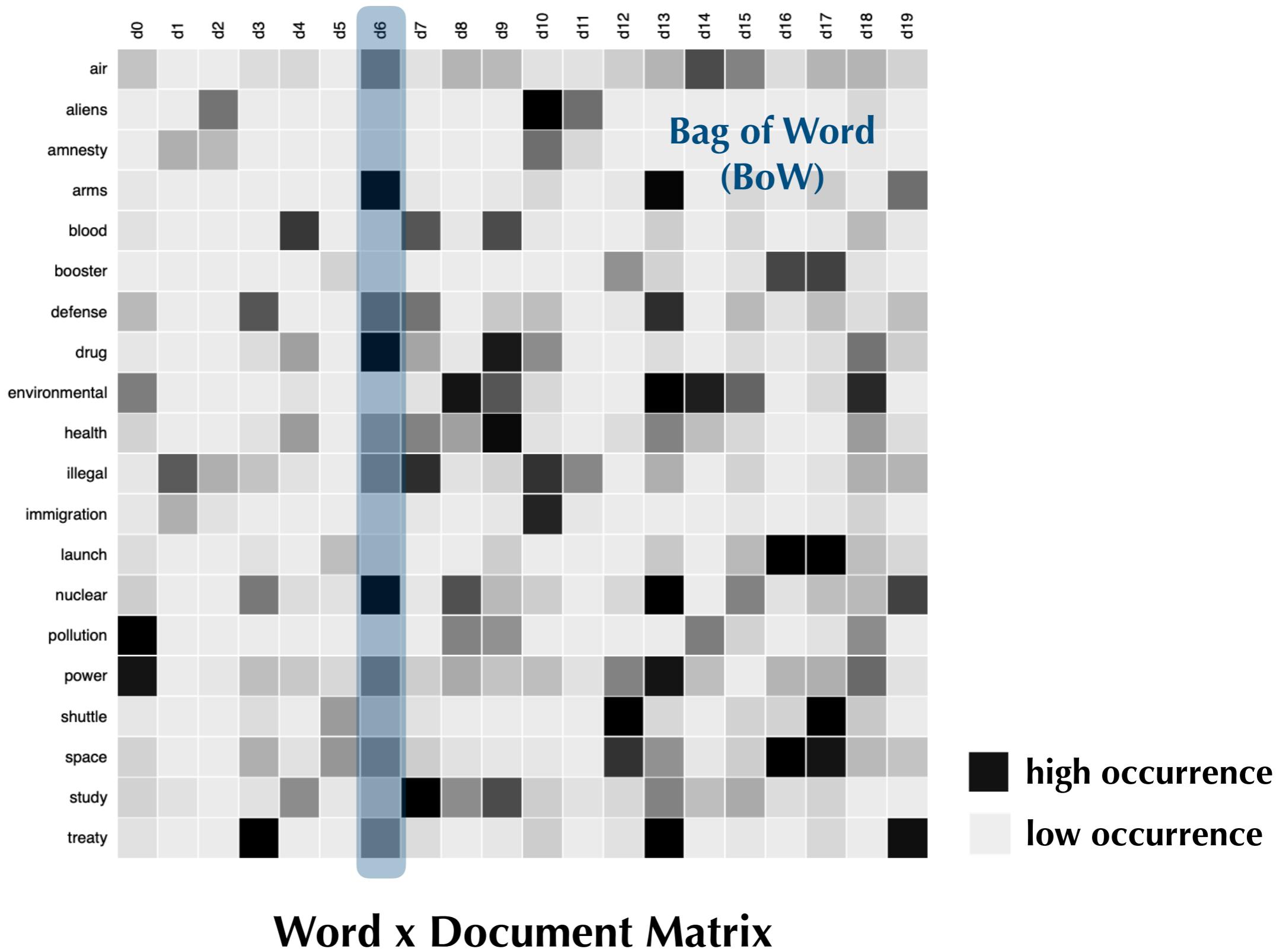
Latent Semantic Indexing (LSI)



Word x Document Matrix

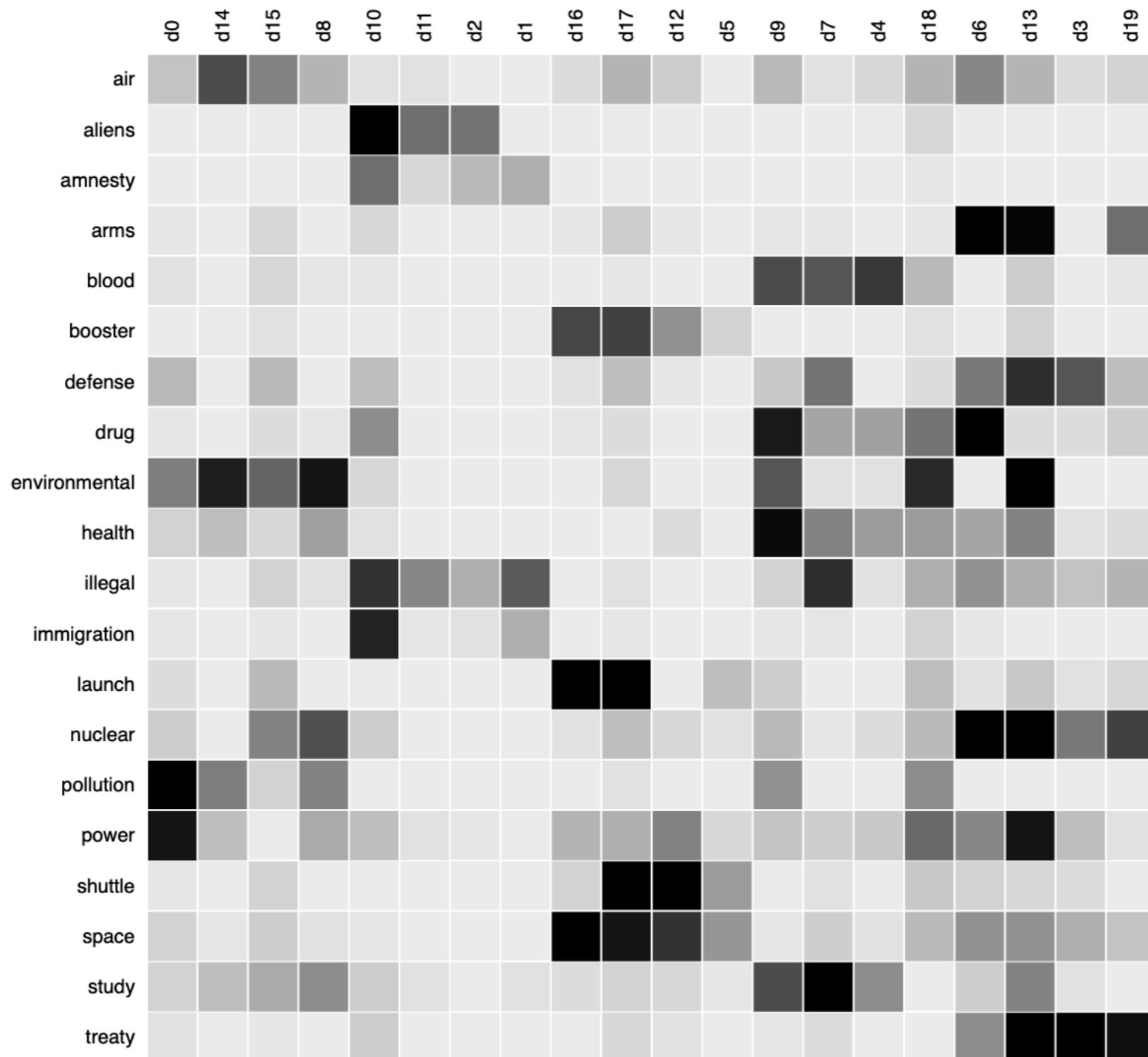
Deerwester et al., *Indexing by Latent Semantic Analysis*, 1990

Latent Semantic Indexing (LSI)



Deerwester et al., Indexing by Latent Semantic Analysis, 1990

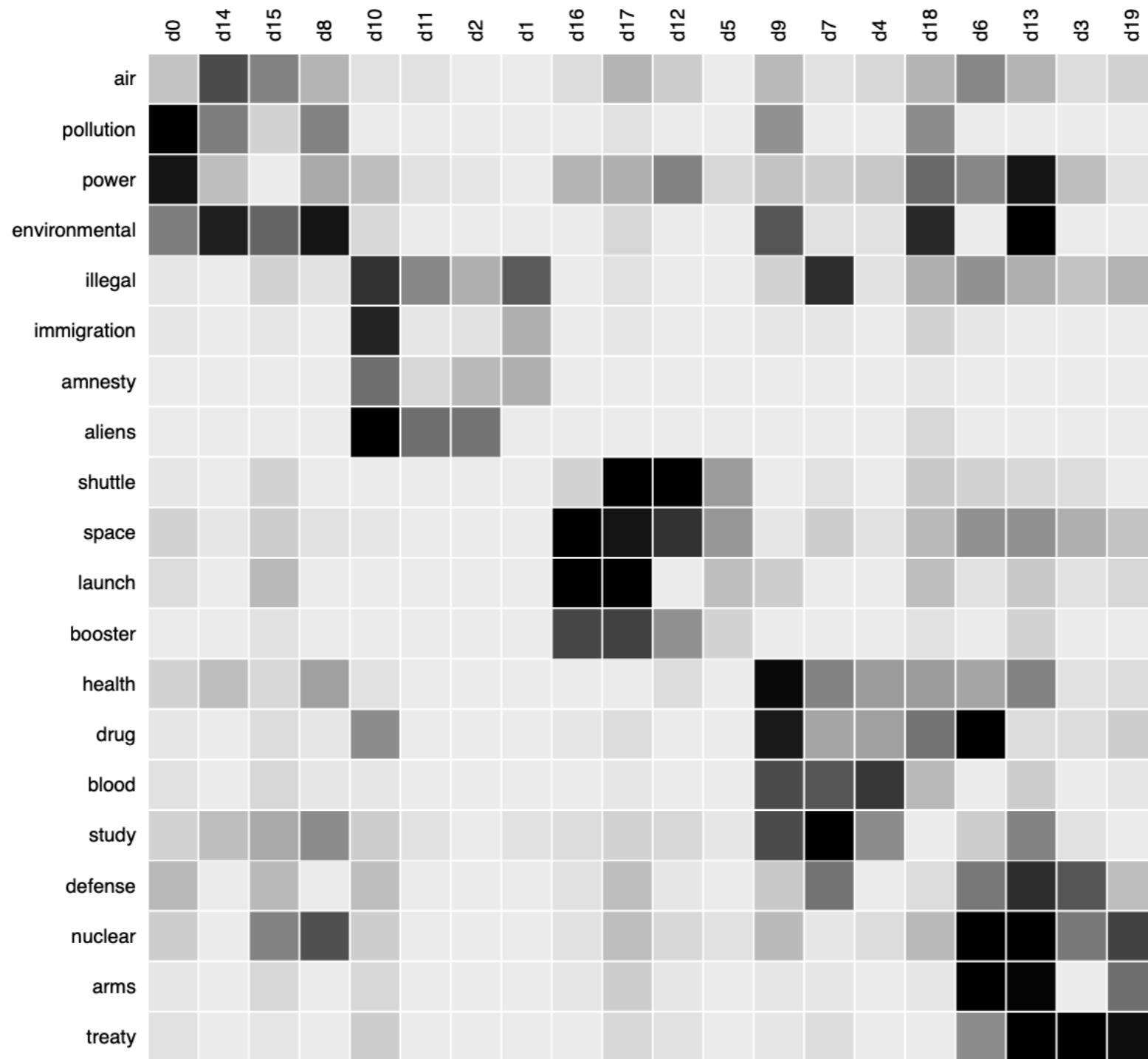
Latent Semantic Indexing (LSI)



Word x Document Matrix - Sort by columns (docs with similar words)

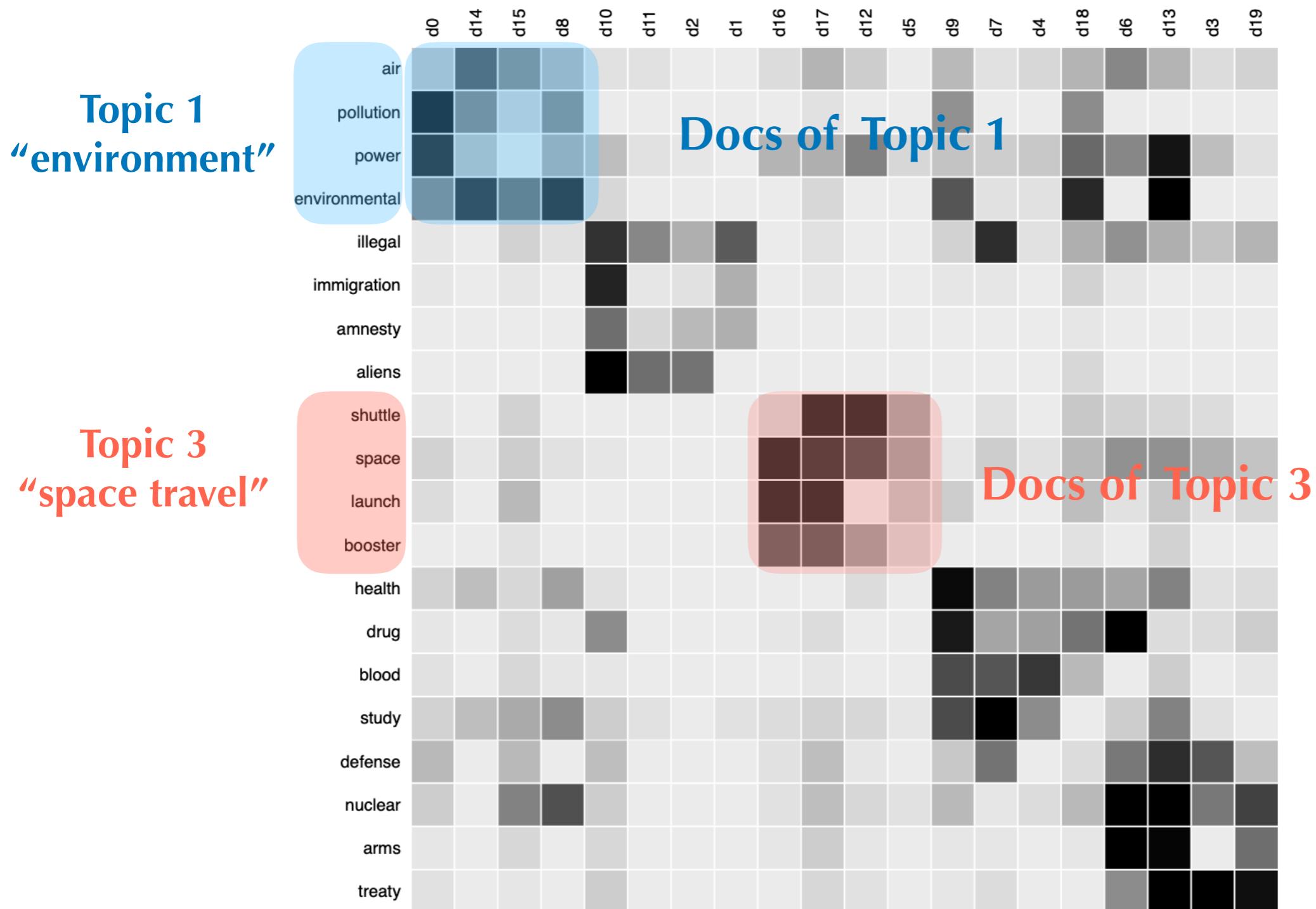
Deerwester et al., Indexing by Latent Semantic Analysis, 1990

Latent Semantic Indexing (LSI)



Word x Document Matrix - Sort by rows (words occurring in similar docs)

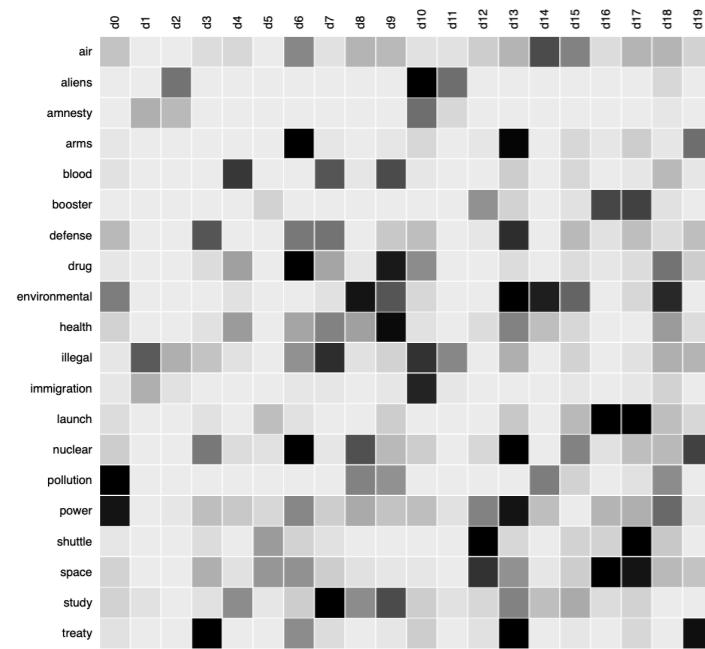
Latent Semantic Indexing (LSI)



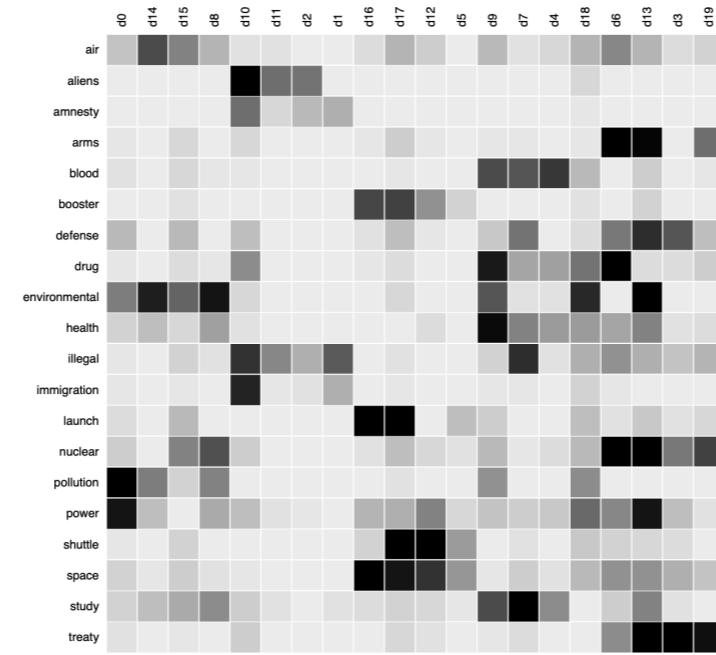
Word x Document Matrix - Sort by rows (words occurring in similar docs)

Latent Semantic Indexing (LSI)

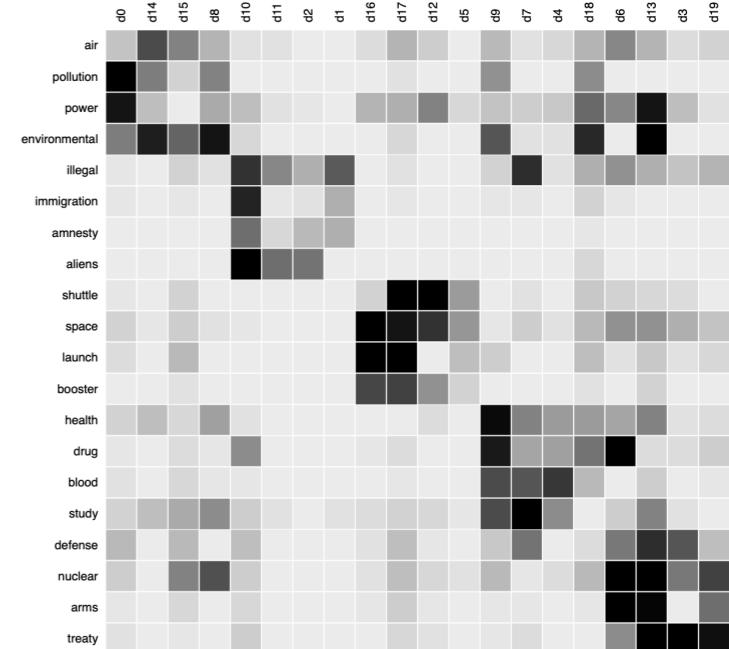
How to sort / cluster?



Word x Document
Matrix



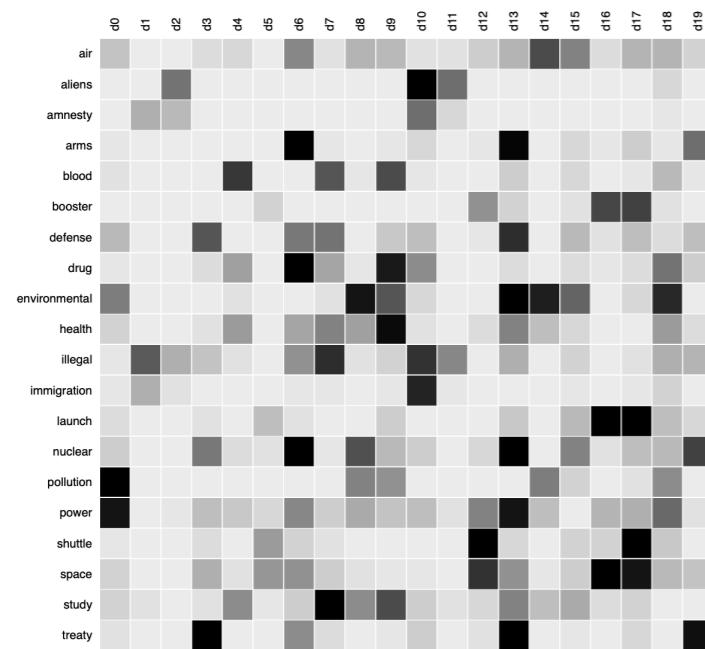
Sort by columns
(docs with similar words)



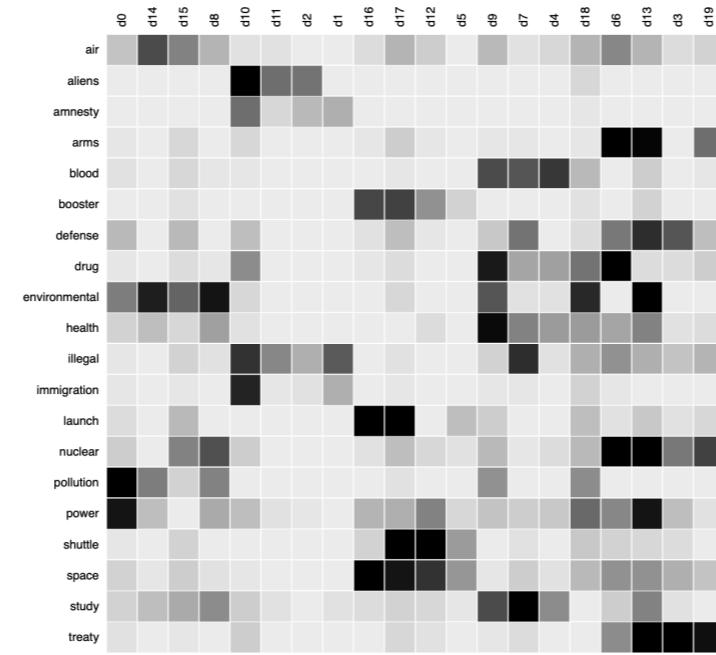
Sort by rows
(words in similar docs)

Latent Semantic Indexing (LSI)

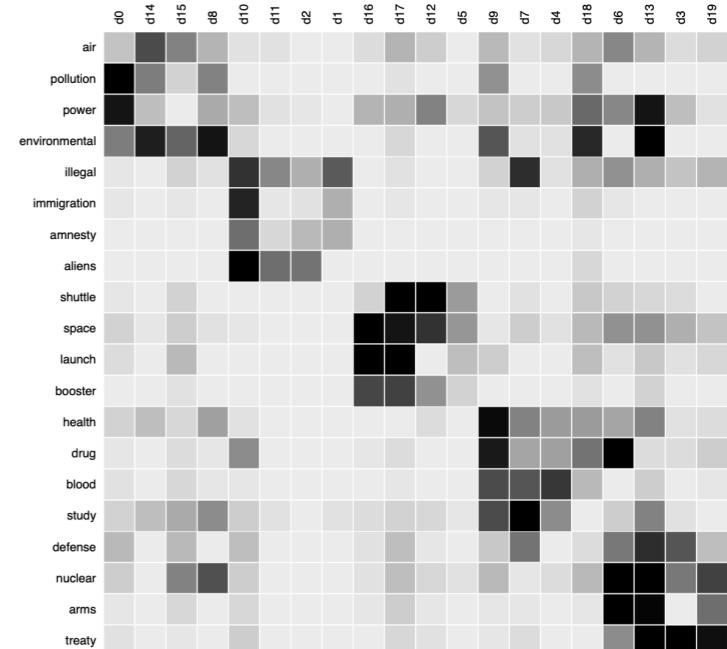
How to sort / cluster?



Word x Document
Matrix



Sort by columns
(docs with similar words)



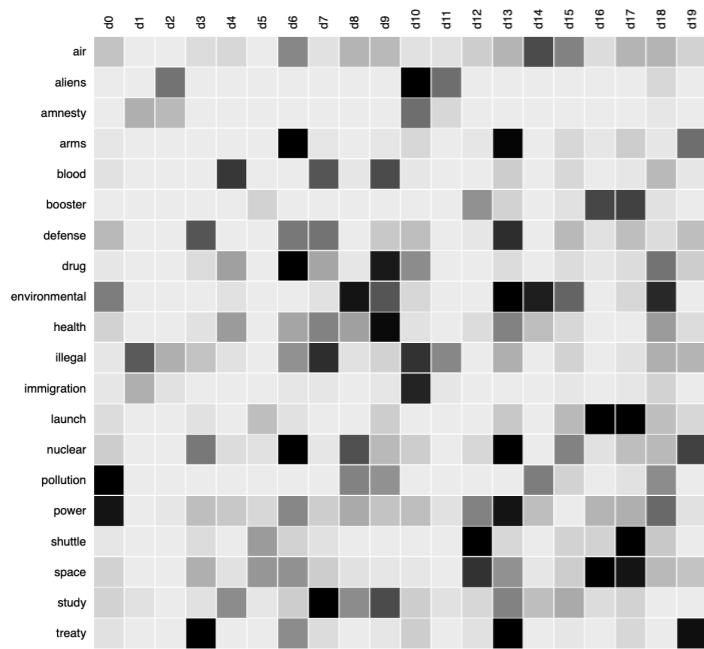
Sort by rows
(words in similar docs)

Directly use the cosine similarity of columns / rows?

Expensive! Noisy! Overly Sparse!

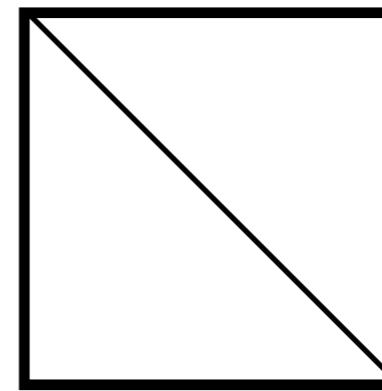
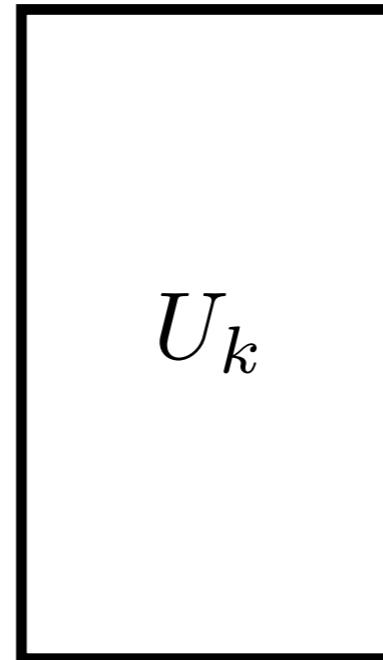
Latent Semantic Indexing (LSI)

How to sort / cluster? - Use SVD

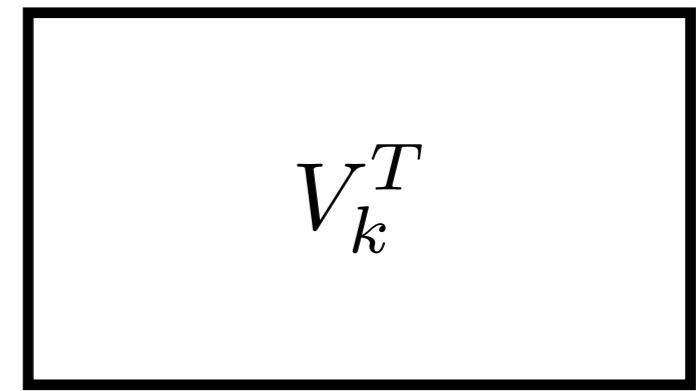


**Word x Document
Matrix**

X



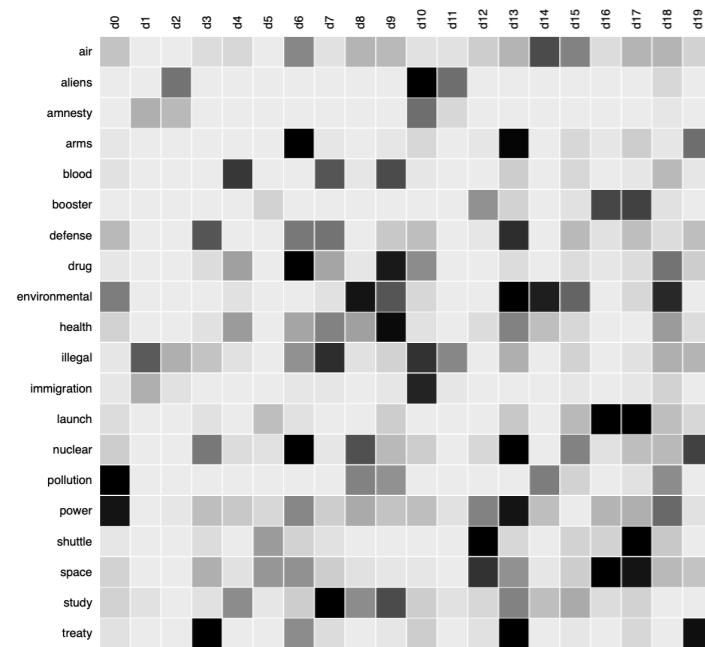
$k \times k$



$$X \approx U_k \Sigma_k V_k^T$$

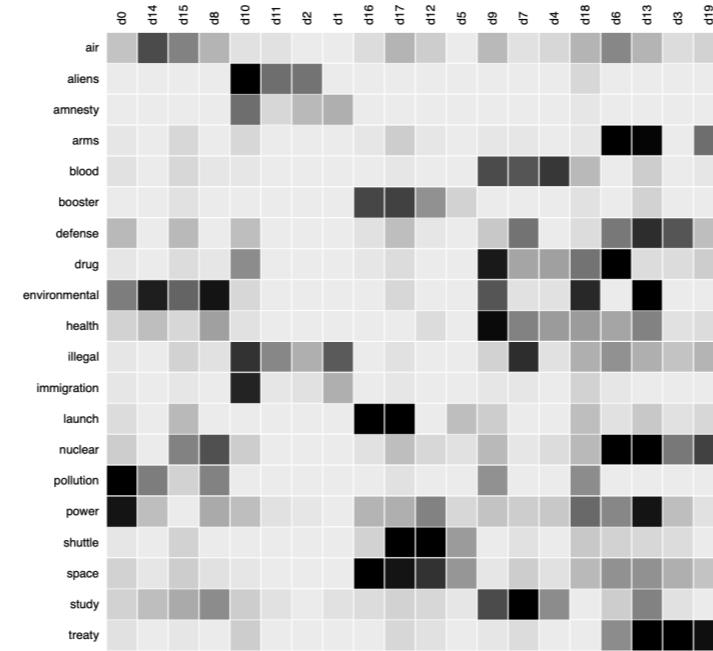
Latent Semantic Indexing (LSI)

How to sort / cluster? - Use SVD



**Word x Document
Matrix**

$$X \approx U_k \Sigma_k V_k^T$$



**Sort by columns
(docs with similar words)**

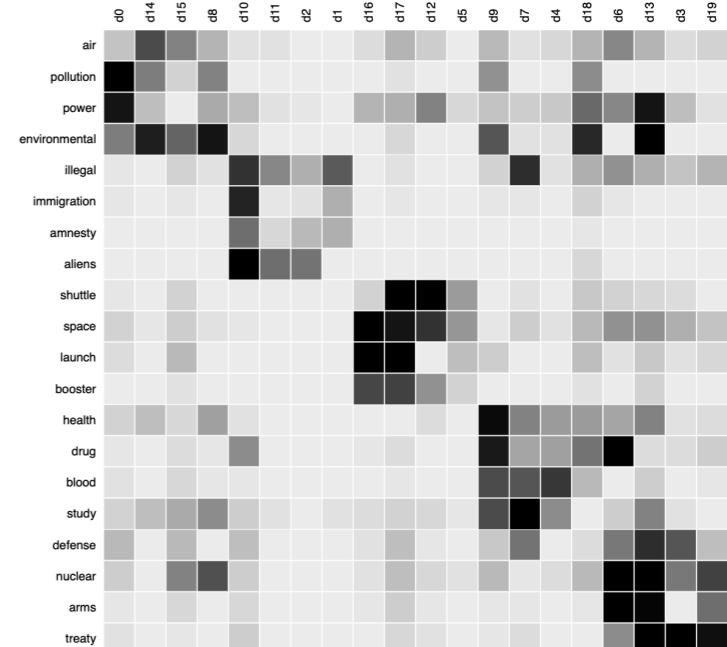
$$X^T X$$

Document
Correlation

sort

$$V_k^T$$

k x Document



**Sort by rows
(words in similar docs)**

$$X X^T$$

Word
Correlation

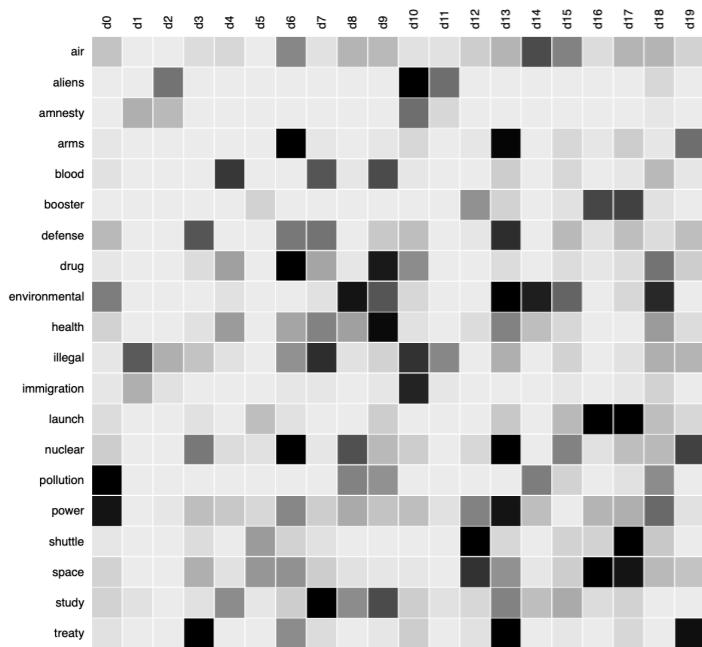
sort

$$U_k$$

Word x k

Latent Semantic Indexing (LSI)

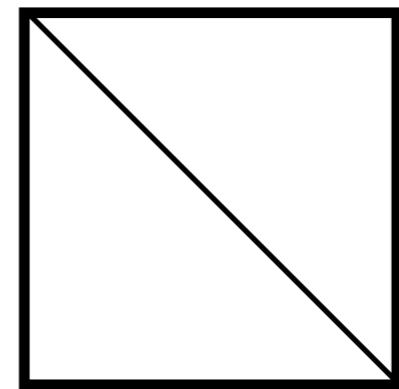
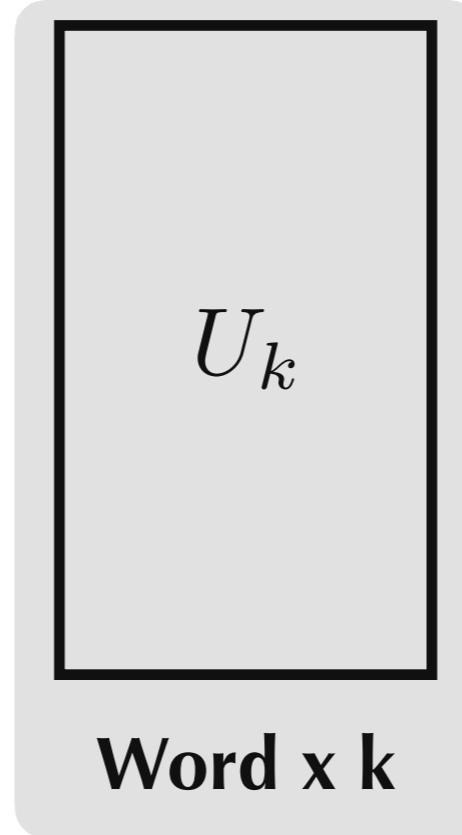
How to sort / cluster? - Use SVD



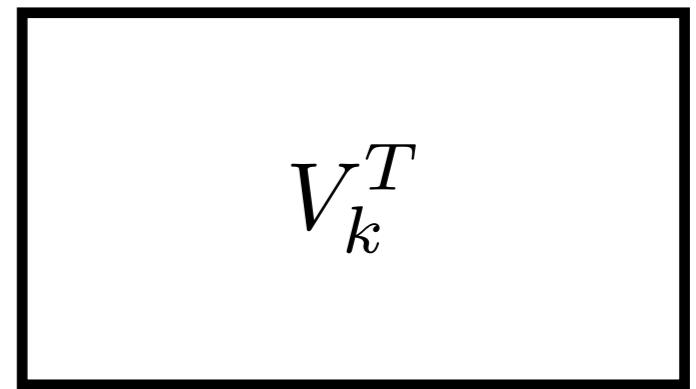
**Word x Document
Matrix**

X

Each column is a topic!



k x k

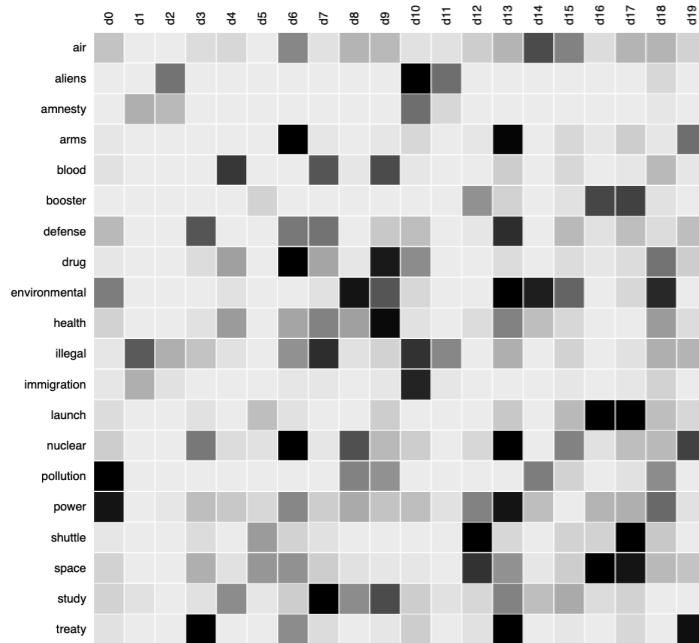


k x Document

$$X \approx U_k \Sigma_k V_k^T$$

Latent Semantic Indexing (LSI)

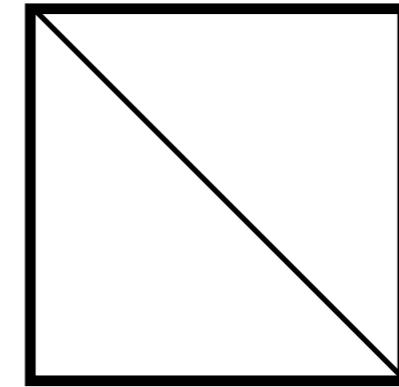
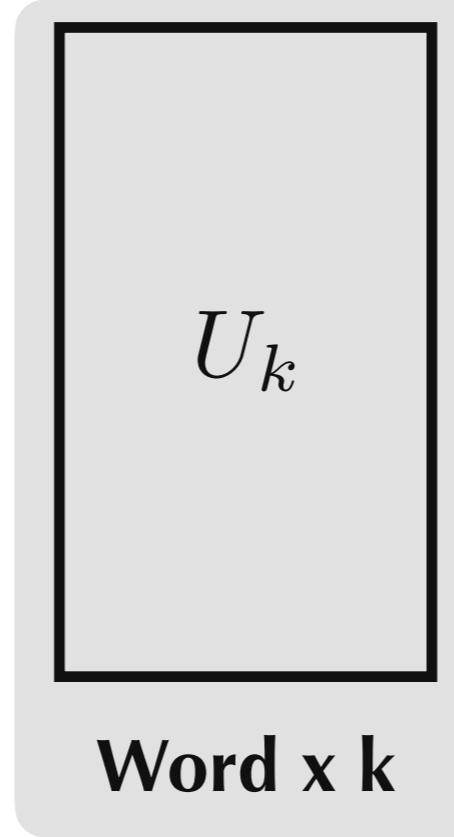
How to sort / cluster? - Use SVD



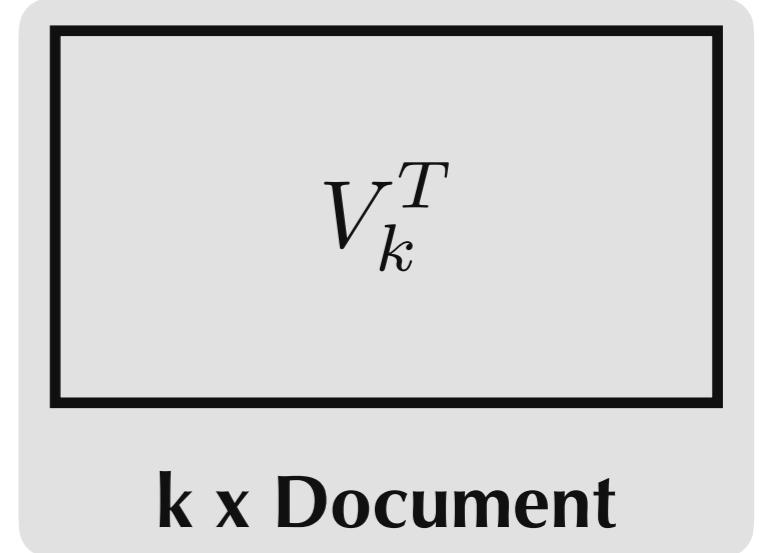
**Word x Document
Matrix**

X

Each column is a topic!



$k \times k$

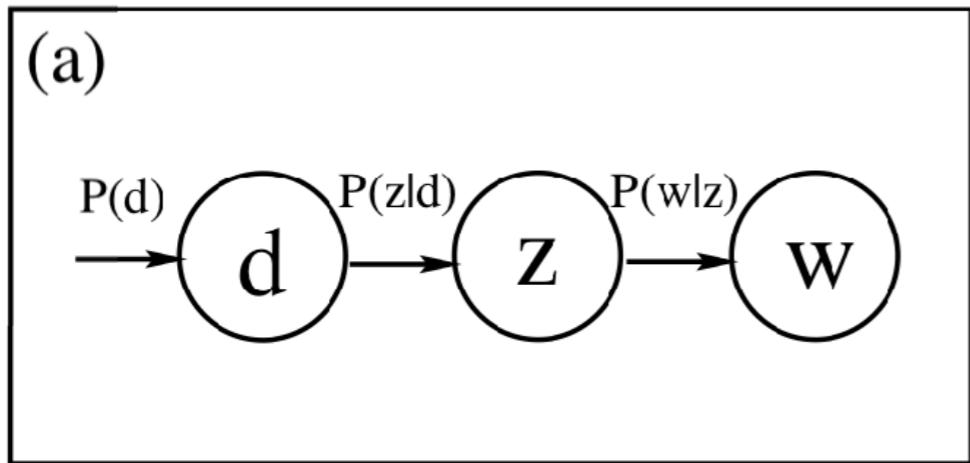


$k \times \text{Document}$

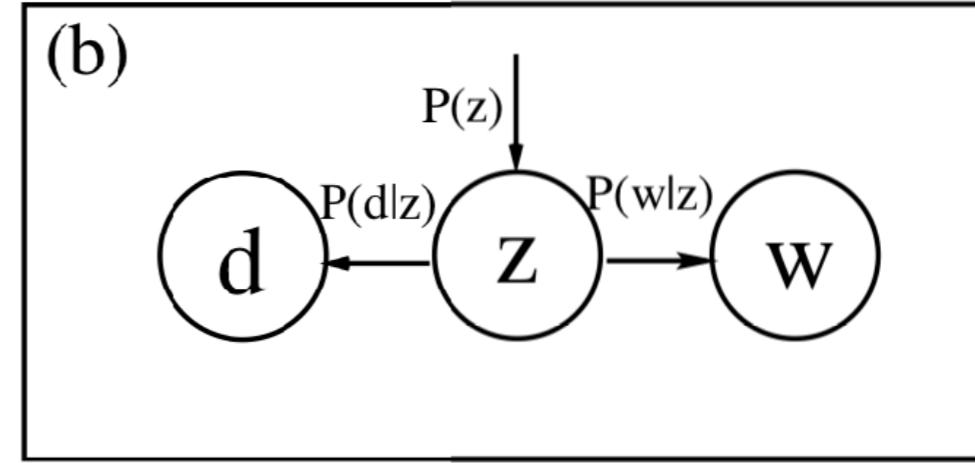
$$X \approx U_k \Sigma_k V_k^T$$

- elements of U and V can be negative, how to interpret?

Probabilistic Latent Semantic Indexing (pLSI)



Asymmetric

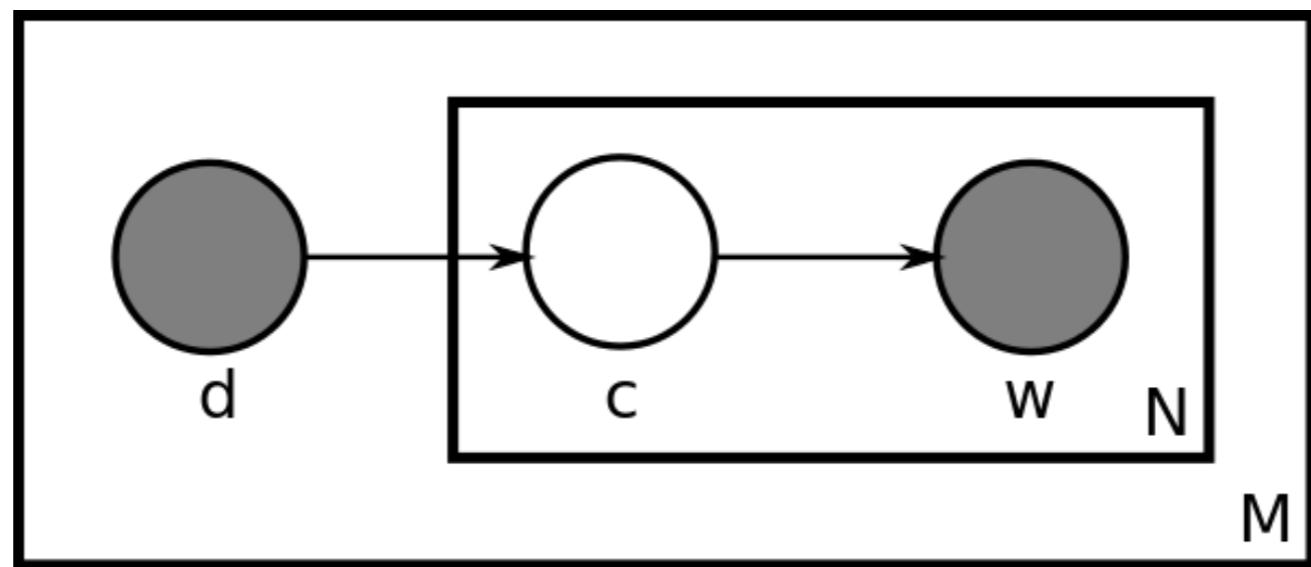


Symmetric

two equivalent probabilistic views

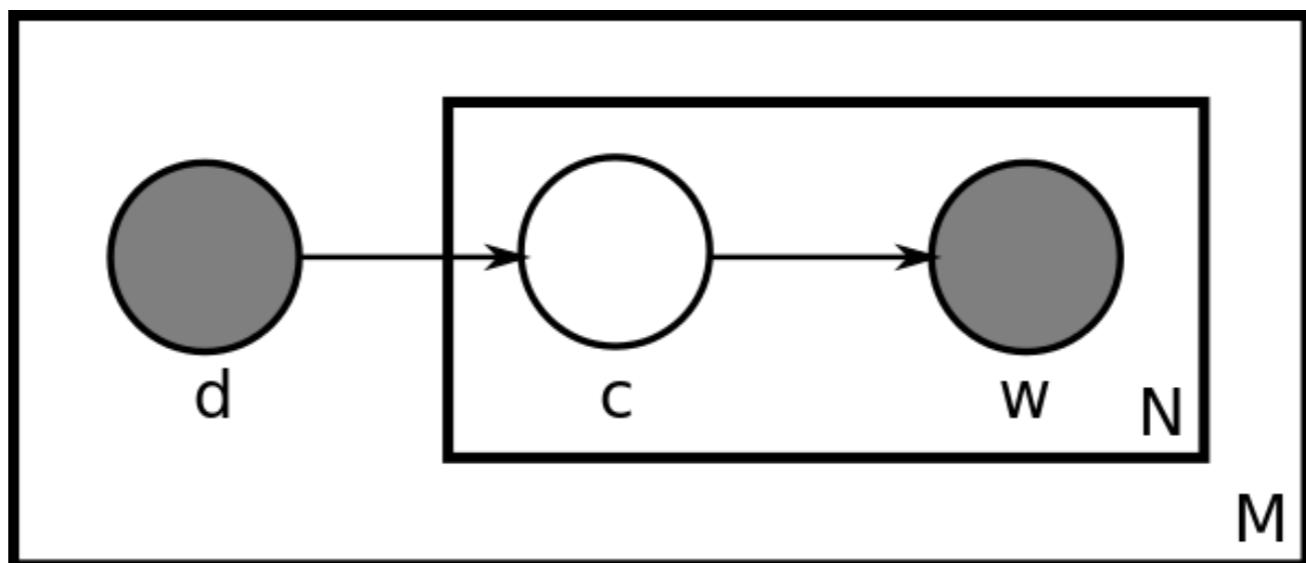
d - document
c - topic
w - word

Probabilistic Latent Semantic Indexing (pLSI)



$$P(w, d) = \sum_c P(c) P(d|c) P(w|c) = P(d) \sum_c P(c|d) P(w|c)$$

Probabilistic Latent Semantic Indexing (pLSI)

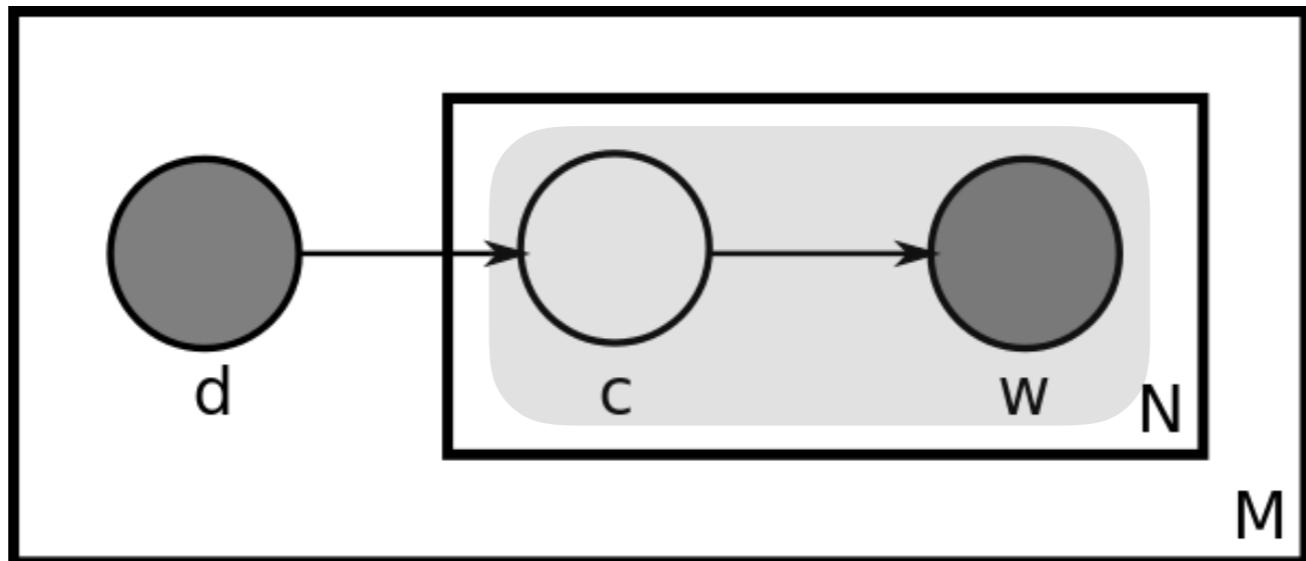


$$P(w, d) = \sum_c P(c) P(d|c) P(w|c) = P(d) \sum_c P(c|d) P(w|c)$$

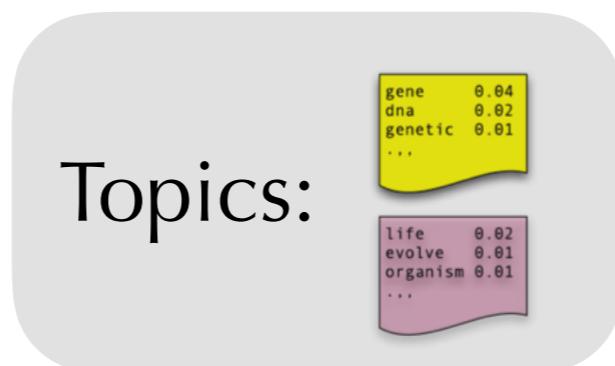
Symmetric

Asymmetric

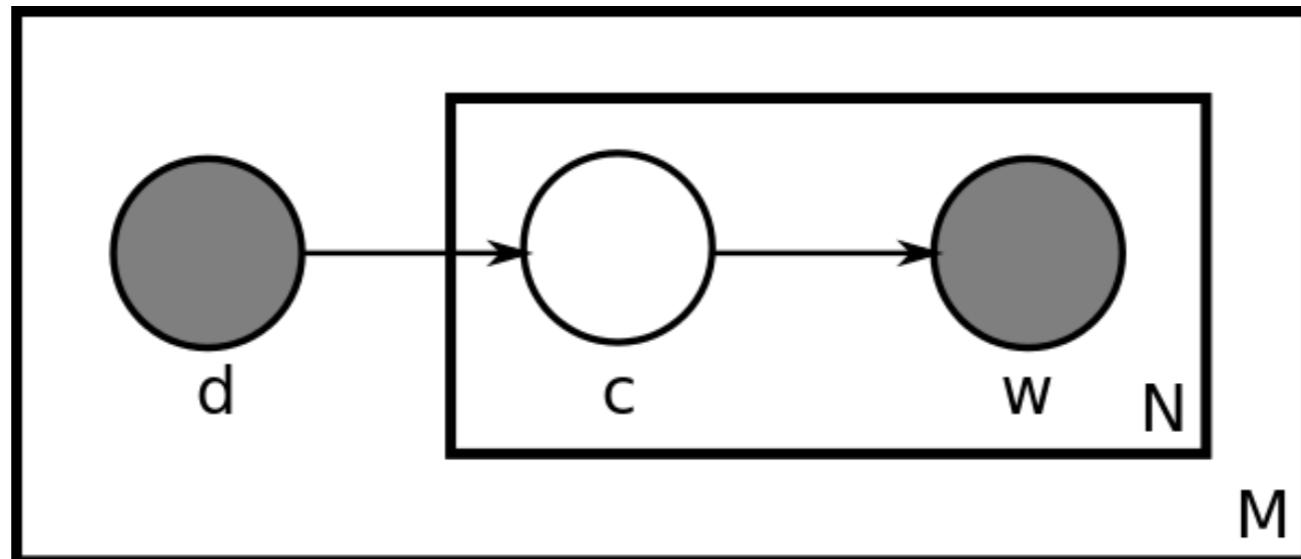
Probabilistic Latent Semantic Indexing (pLSI)



$$P(w, d) = \sum_c P(c) P(d|c) P(w|c) = P(d) \sum_c P(c|d) P(w|c)$$



Probabilistic Latent Semantic Indexing (pLSI)



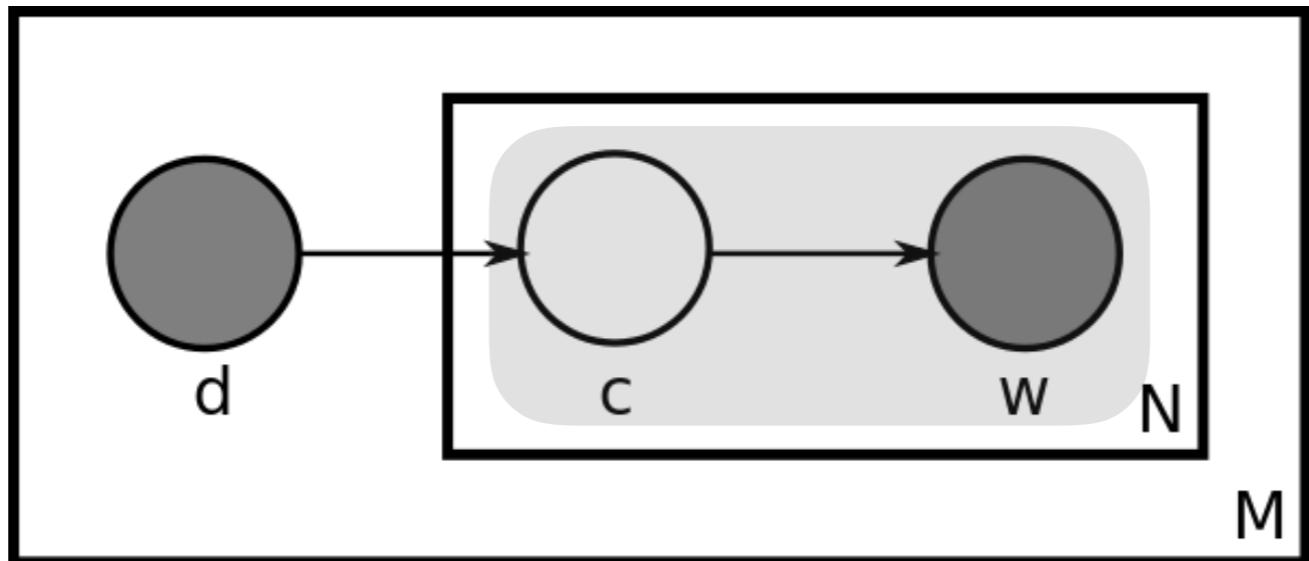
Maximum likelihood estimate of the topic c

E-Step: $P(c|d, w) = \frac{P(d)P(c|d)P(w|c)}{P(d)\sum_{c' \in C} P(w|c')P(c'|d)} = \frac{P(c|d)P(w|c)}{\sum_{c' \in C} P(w|c')P(c'|d)}$

M-Step:

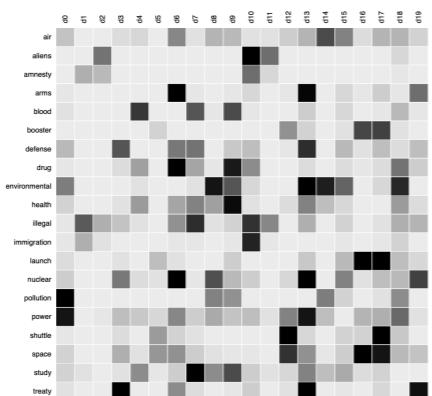
$$p(w|c) \propto \sum_{d \in \mathcal{D}} n(d, w)P(c|d, w) \quad p(c) \propto \sum_{d \in \mathcal{D}} \sum_{w \in \mathcal{W}} n(d, w)P(c|d, w)$$
$$p(d|c) \propto \sum_{w \in \mathcal{W}} n(d, w)P(c|d, w) \quad p(d|c) \propto p(d|c)p(c)$$

Probabilistic Latent Semantic Indexing (pLSI)



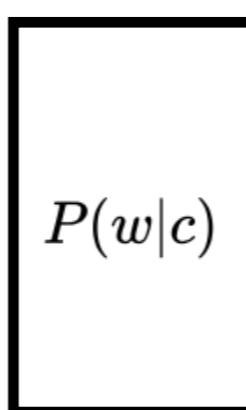
$$P(w|d) = \sum_c p(w|c)p(c|d)$$

$$P(w|d)$$

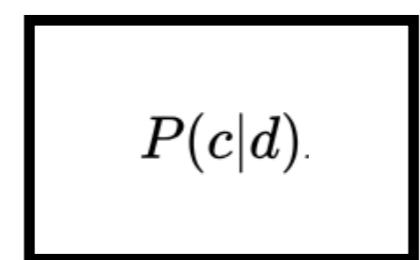


**Word x Document
Frequency Matrix**

\approx

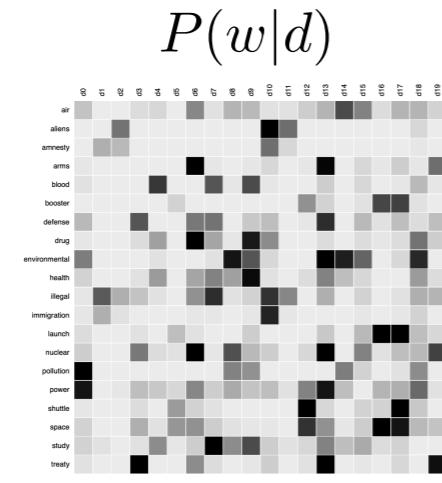


Word x k



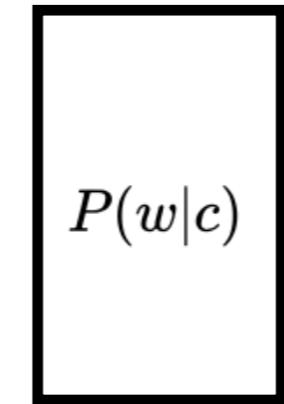
k x Document

Nonnegative Matrix Factorization (NMF)

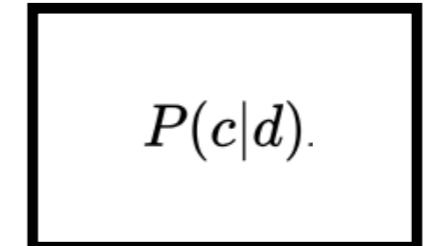


Word x Document
Frequency Matrix

\approx



Word x k



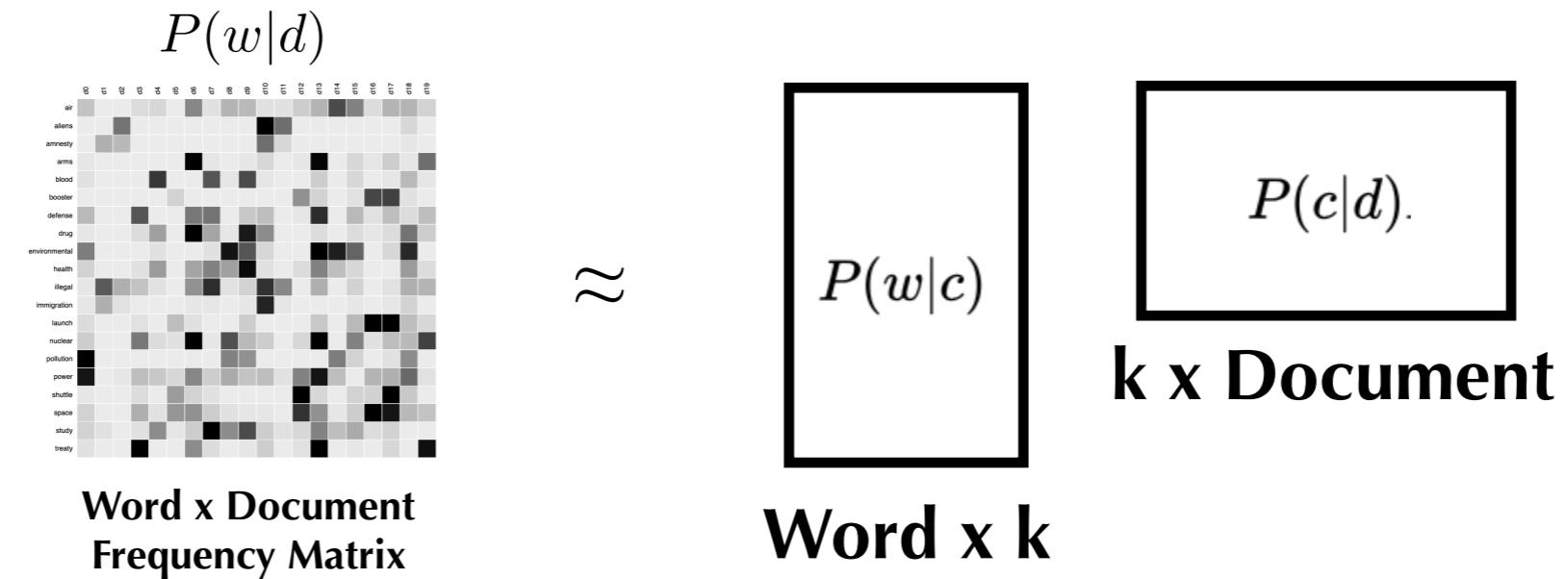
k x Document

V

$W \geq 0$

$H \geq 0$

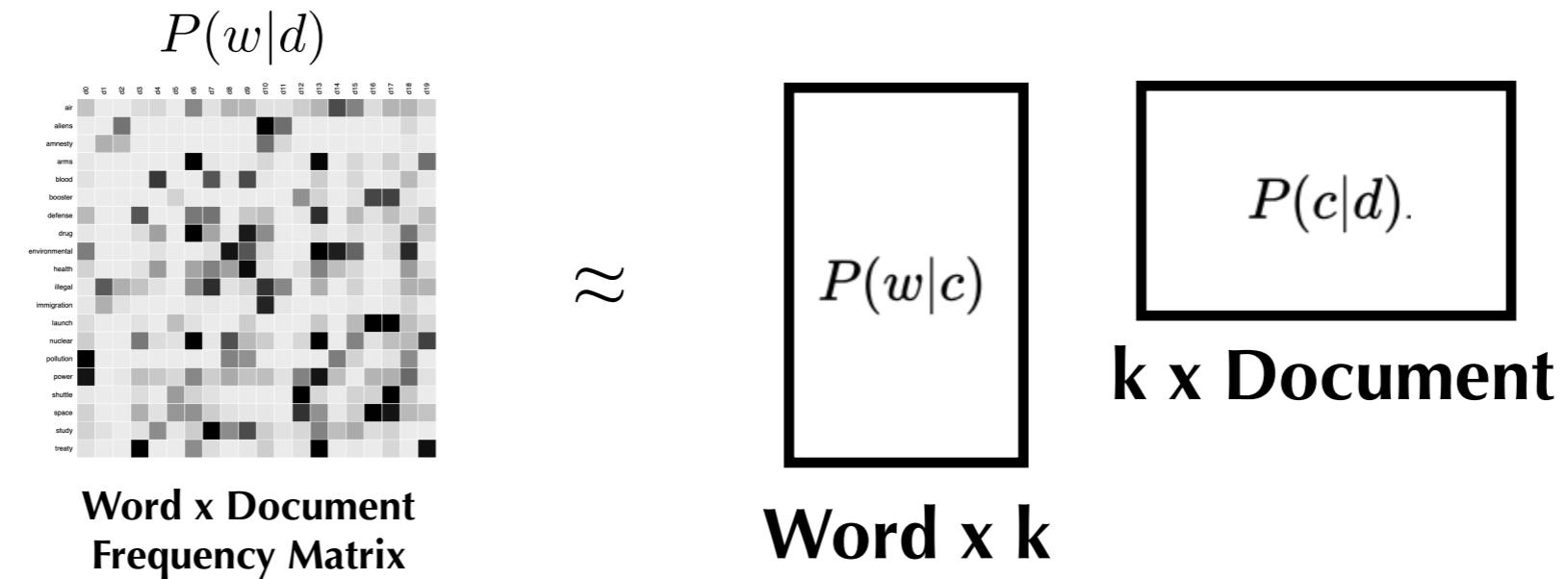
Nonnegative Matrix Factorization (NMF)



Given V , find $W \geq 0$ $H \geq 0$

to minimize $\|H - WH\|_F$ or $D(H||WH)$

Nonnegative Matrix Factorization (NMF)



Given V , find $W \geq 0$ $H \geq 0$

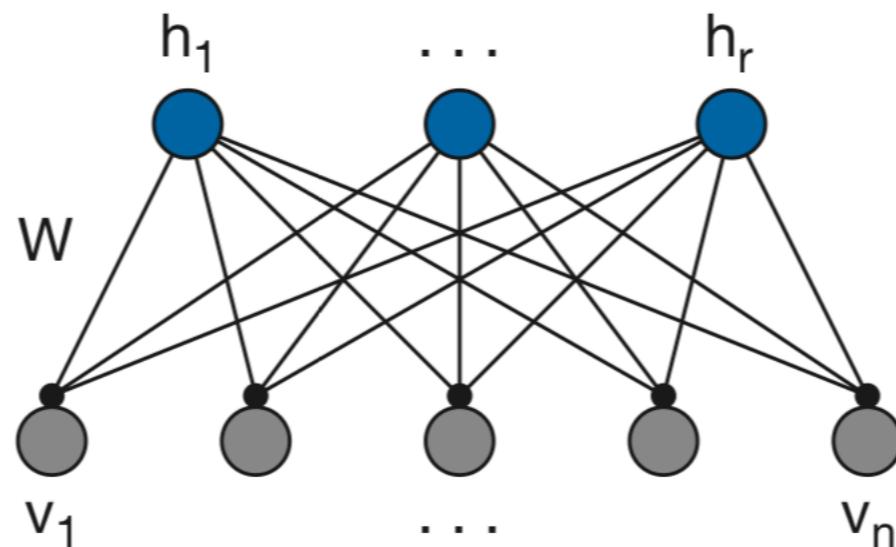
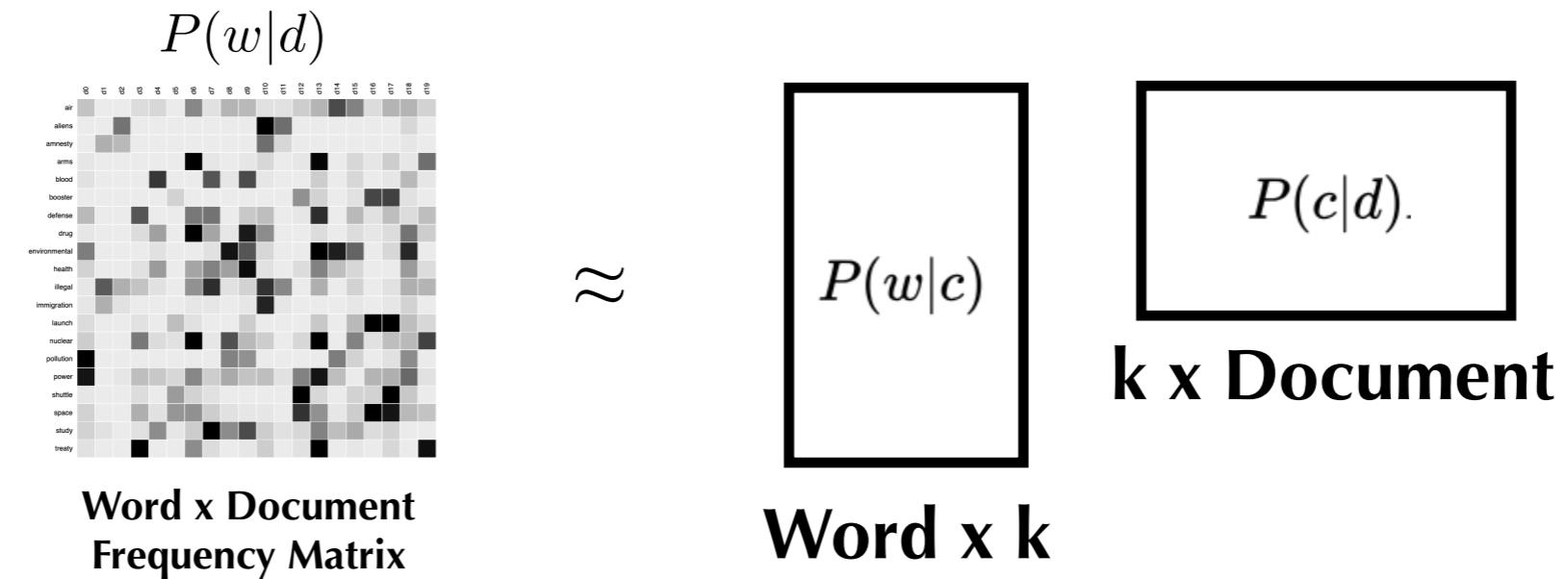
to minimize $\|H - WH\|_F$ or $D(H||WH)$

Similar to EM:

$$W_{ia} \leftarrow W_{ia} \sum_{\mu} \frac{V_{i\mu}}{(WH)_{i\mu}} H_{a\mu}$$
$$W_{ia} \leftarrow \frac{W_{ia}}{\sum_j W_{ja}}$$

$$H_{a\mu} \leftarrow H_{a\mu} \sum_i W_{ia} \frac{V_{i\mu}}{(WH)_{i\mu}}$$

Nonnegative Matrix Factorization (NMF)

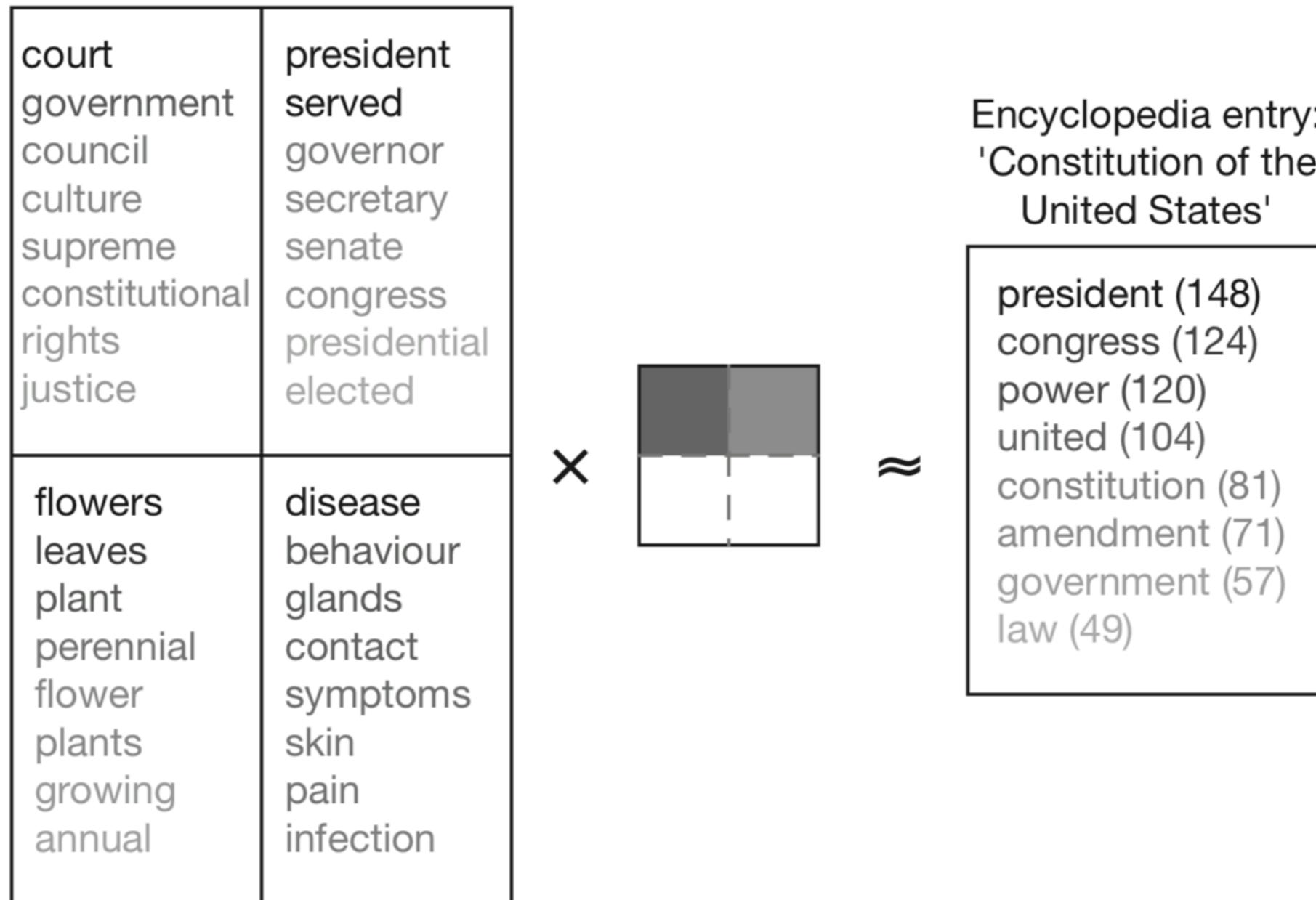


Probabilistic hidden variables model underlying non-negative matrix factorization.

The bottom layer of nodes are generated from the hidden variables.

The influence of h_a on v_i is represented by a connection with strength W_{ia} .

Nonnegative Matrix Factorization (NMF)



Non-negative matrix factorization (NMF) discovers semantic features of 30,991 articles from the Grolier encyclopedia.

Nonnegative Matrix Factorization (NMF)

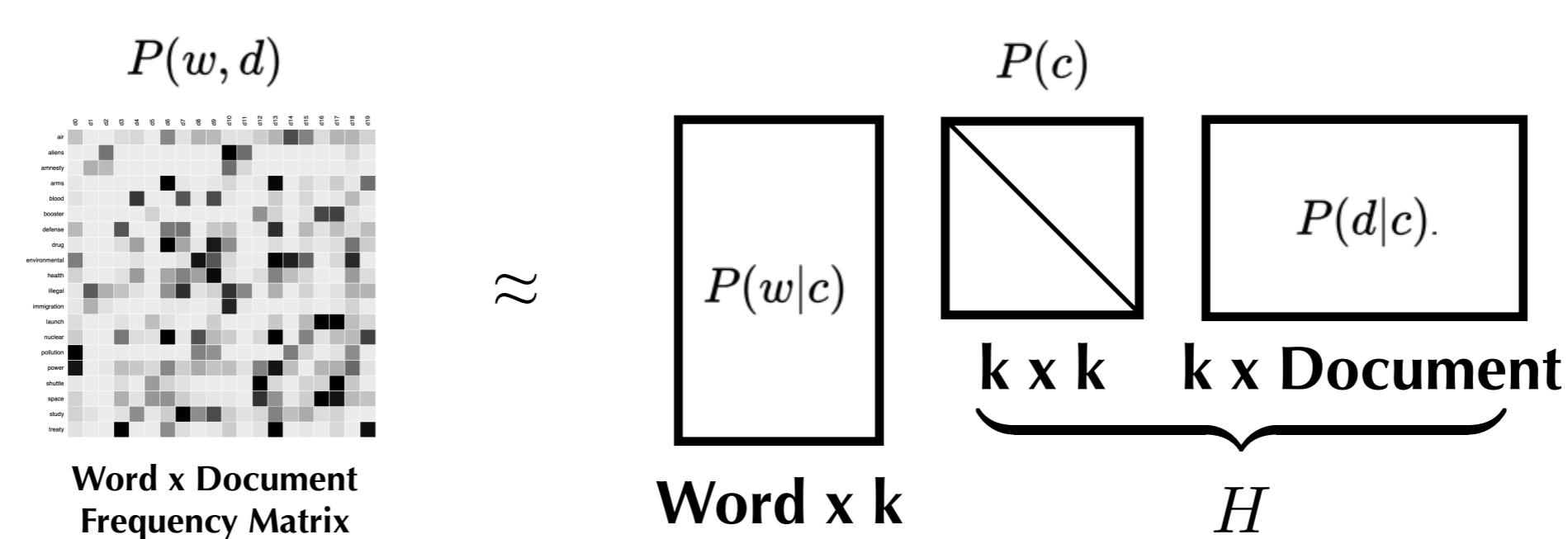
Equivalence to pLSI — H is non-normalized?

Proposition 1. The objective function of PLSI is identical to the objective function of NMF, i.e.,

$$\max J_{\text{PLSI}} \iff \min J_{\text{NMF}} \quad (6)$$

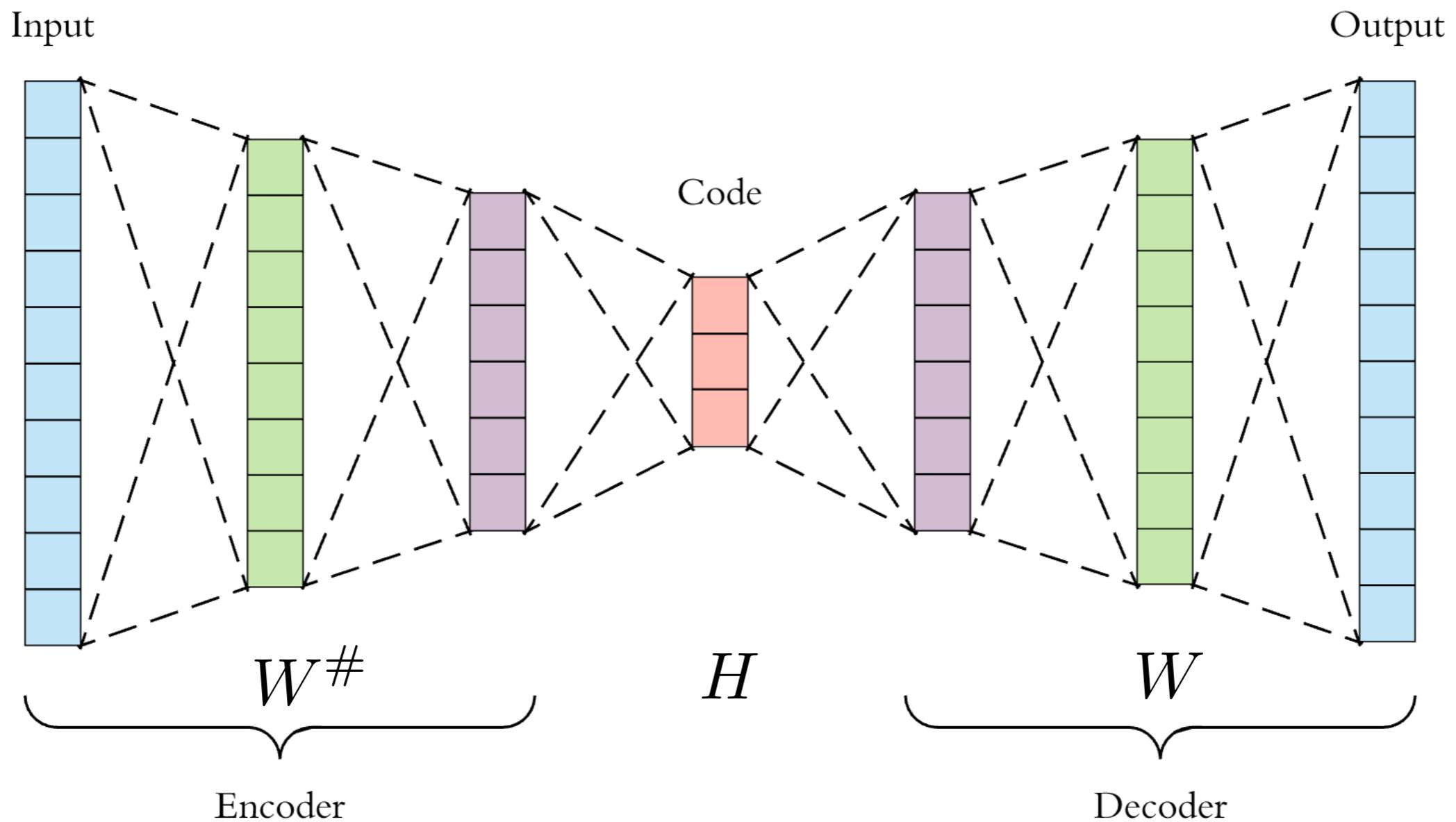
Proposition 2. Column normalized NMF of Eq.(1) is equivalent to the probability factorization of Eq.(4), i.e., $(CH^T)_{ij} = P(w_i, d_j)$.

$$P(w, d) = \sum_c P(c)P(d|c)P(w|c)$$



Autoencoder (AE)

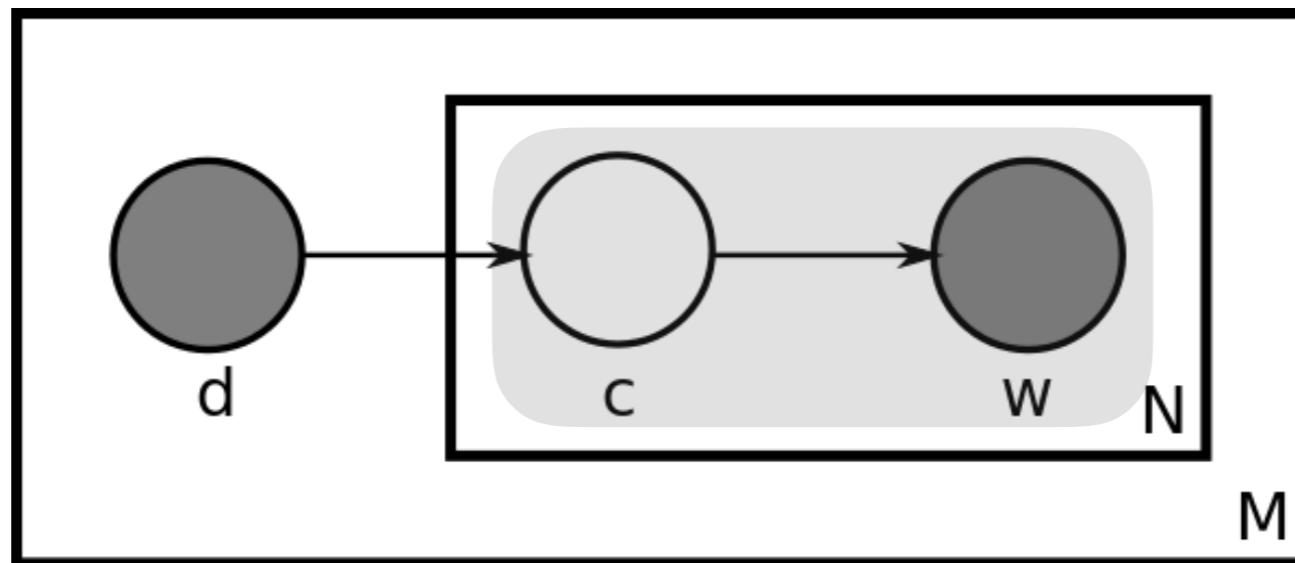
$$V \approx WH$$



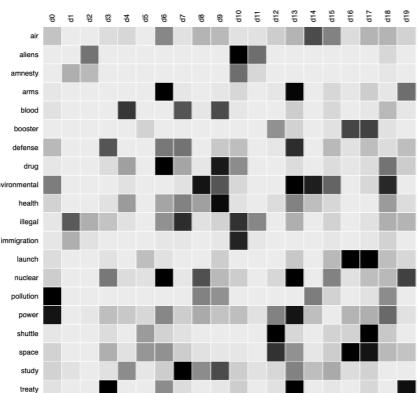
$$W \# V = H$$

$$WH \approx V$$

Probabilistic Latent Semantic Indexing (pLSI)



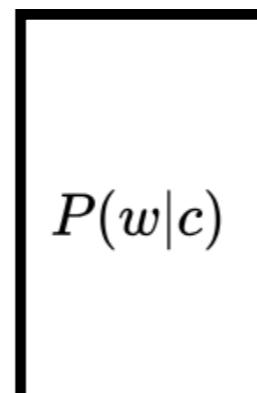
$$P(w, d)$$



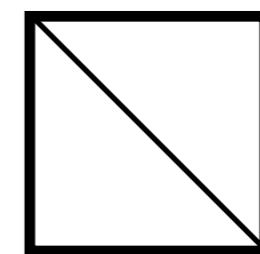
Word x Document
Frequency Matrix

\approx

$$P(c)$$



Word x k

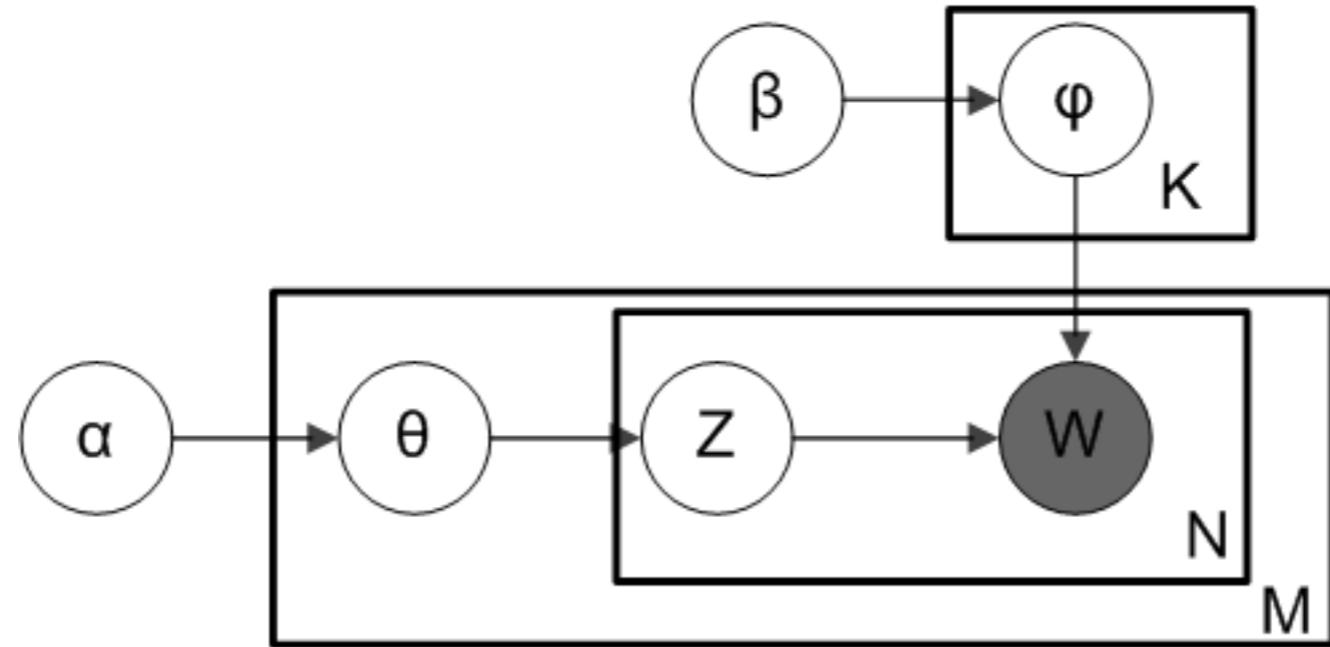


$k \times k$ $k \times \text{Document}$

H

- MLE, no control on topic / word sparsity?

Latent Dirichlet Analysis (LDA)



z: *topics*, **w:** *words*

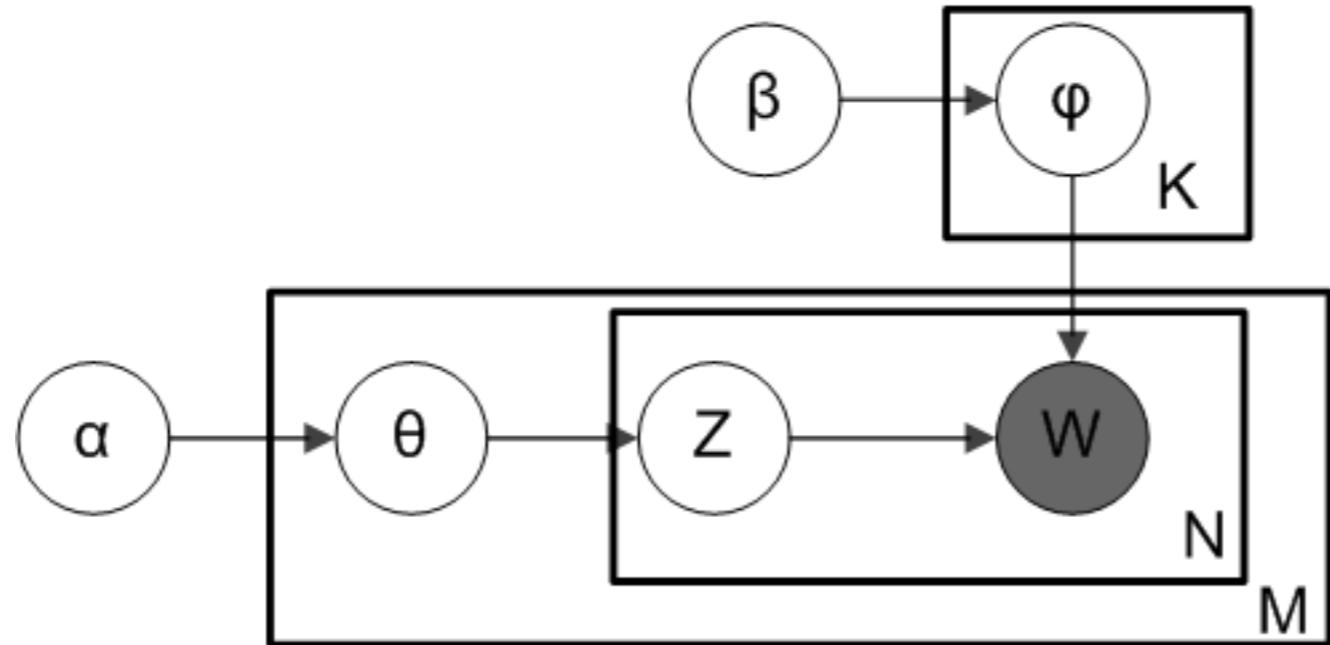
θ : *topic distributions for documents*

φ : *word distributions for topics*

α : the parameter of the Dirichlet prior
on the per-document topic distributions

β : *the parameter of the Dirichlet prior
on the per-topic word distribution*

Latent Dirichlet Analysis (LDA)



Generative process:

$$\varphi_{k=1\dots K} \sim \text{Dirichlet}_V(\beta)$$

$$\theta_{d=1\dots M} \sim \text{Dirichlet}_K(\alpha)$$

$$z_{d=1\dots M, w=1\dots N_d} \sim \text{Categorical}_K(\theta_d)$$

$$w_{d=1\dots M, w=1\dots N_d} \sim \text{Categorical}_V(\varphi_{z_{dw}})$$

Latent Dirichlet Analysis (LDA)

α , β control doc-topic, topic-word sparsity

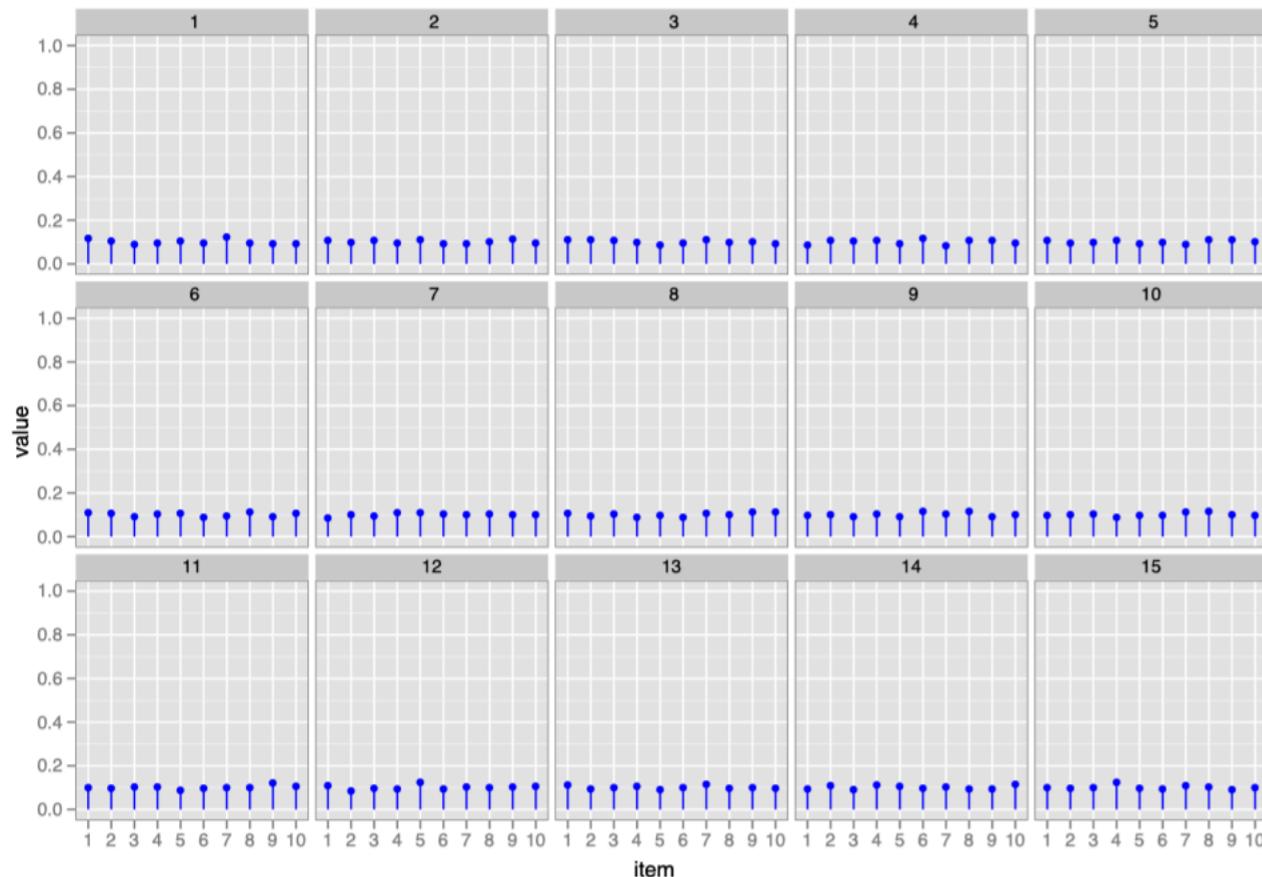
$$P(\theta_j | \alpha) = \frac{\Gamma\left(\sum_{i=1}^K \alpha_i\right)}{\prod_{i=1}^K \Gamma(\alpha_i)} \prod_{i=1}^K \theta_{j,i}^{\alpha_i - 1}$$

It is **conjugate** to the **multinomial**. Given a multinomial observation, the posterior distribution of θ is a **Dirichlet**.

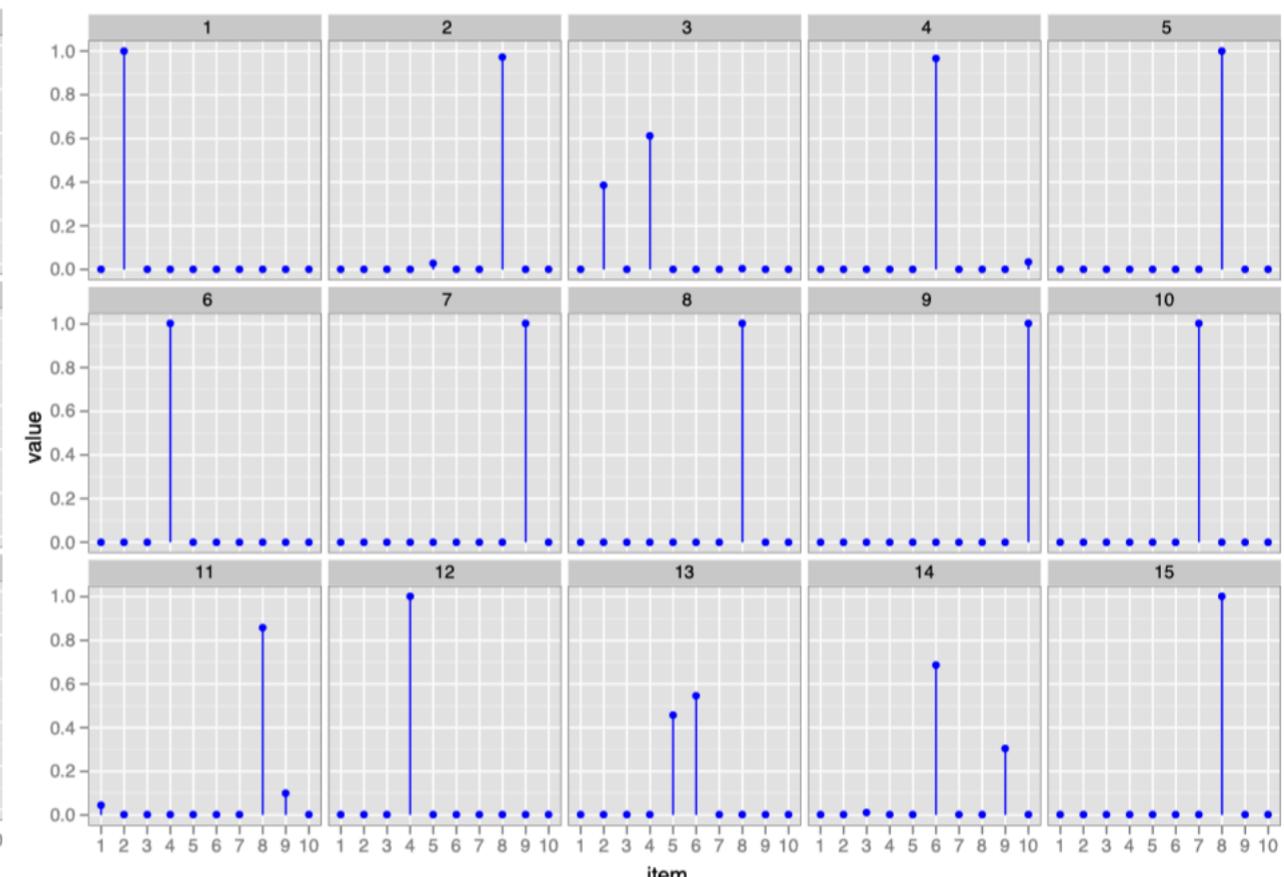
Latent Dirichlet Analysis (LDA)

α, β controls doc-topic, topic-word sparsity

$$P(\theta_j | \alpha) = \frac{\Gamma\left(\sum_{i=1}^K \alpha_i\right)}{\prod_{i=1}^K \Gamma(\alpha_i)} \prod_{i=1}^K \theta_{j,i}^{\alpha_i - 1}$$

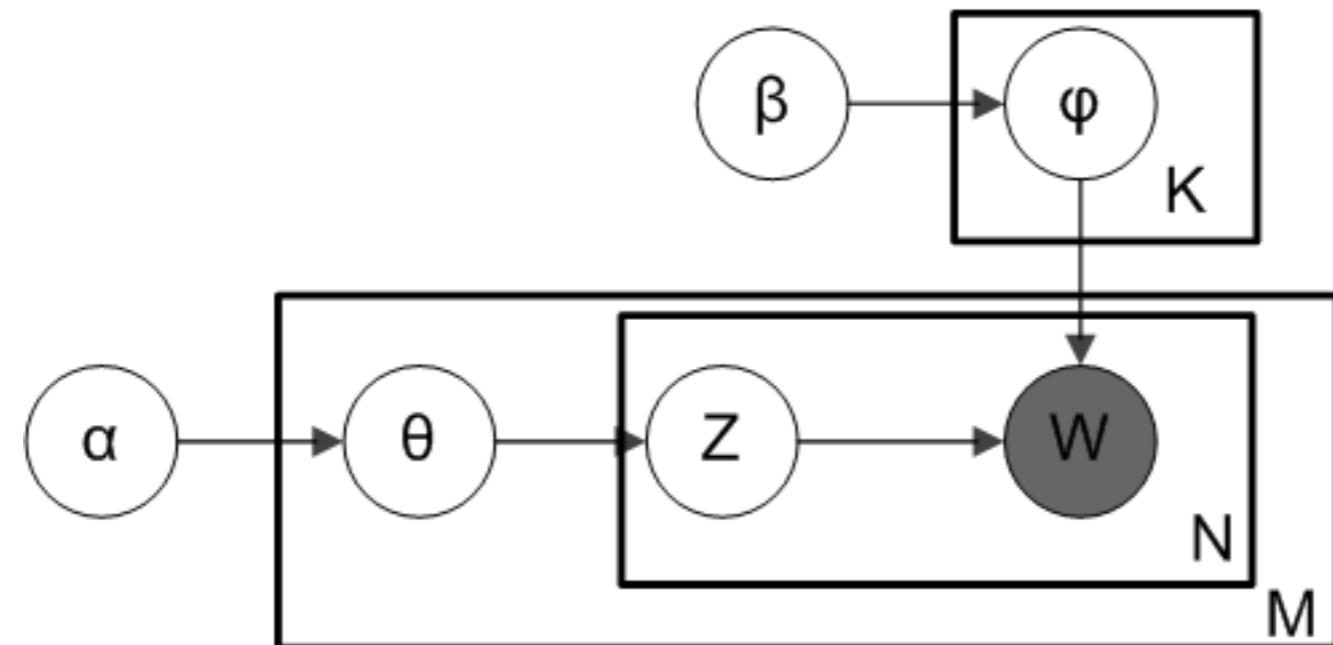


$$\alpha = 100$$



$$\alpha = 0.01$$

Latent Dirichlet Analysis (LDA)

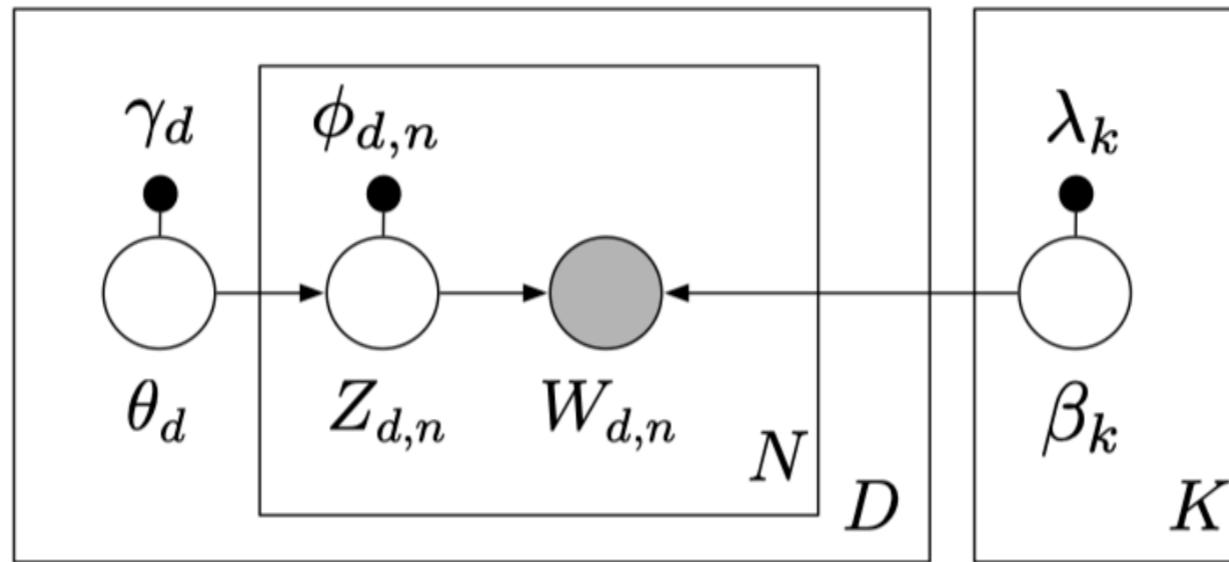


Posterior inference for LDA: $P(\varphi, \theta, z | w)$

$$\frac{P(\varphi, \theta, z, w)}{\int_{\varphi} \int_{\theta} \sum_z P(\varphi, \theta, z, w)}$$

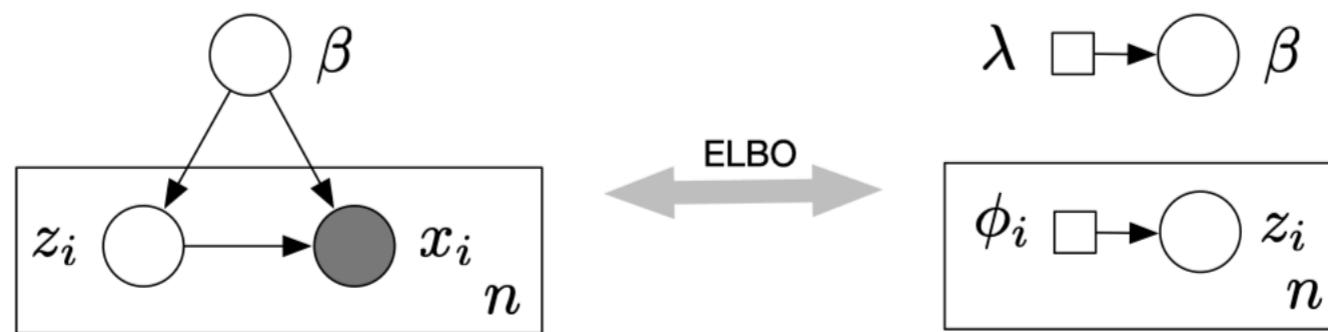
Intractable

Latent Dirichlet Analysis (LDA)

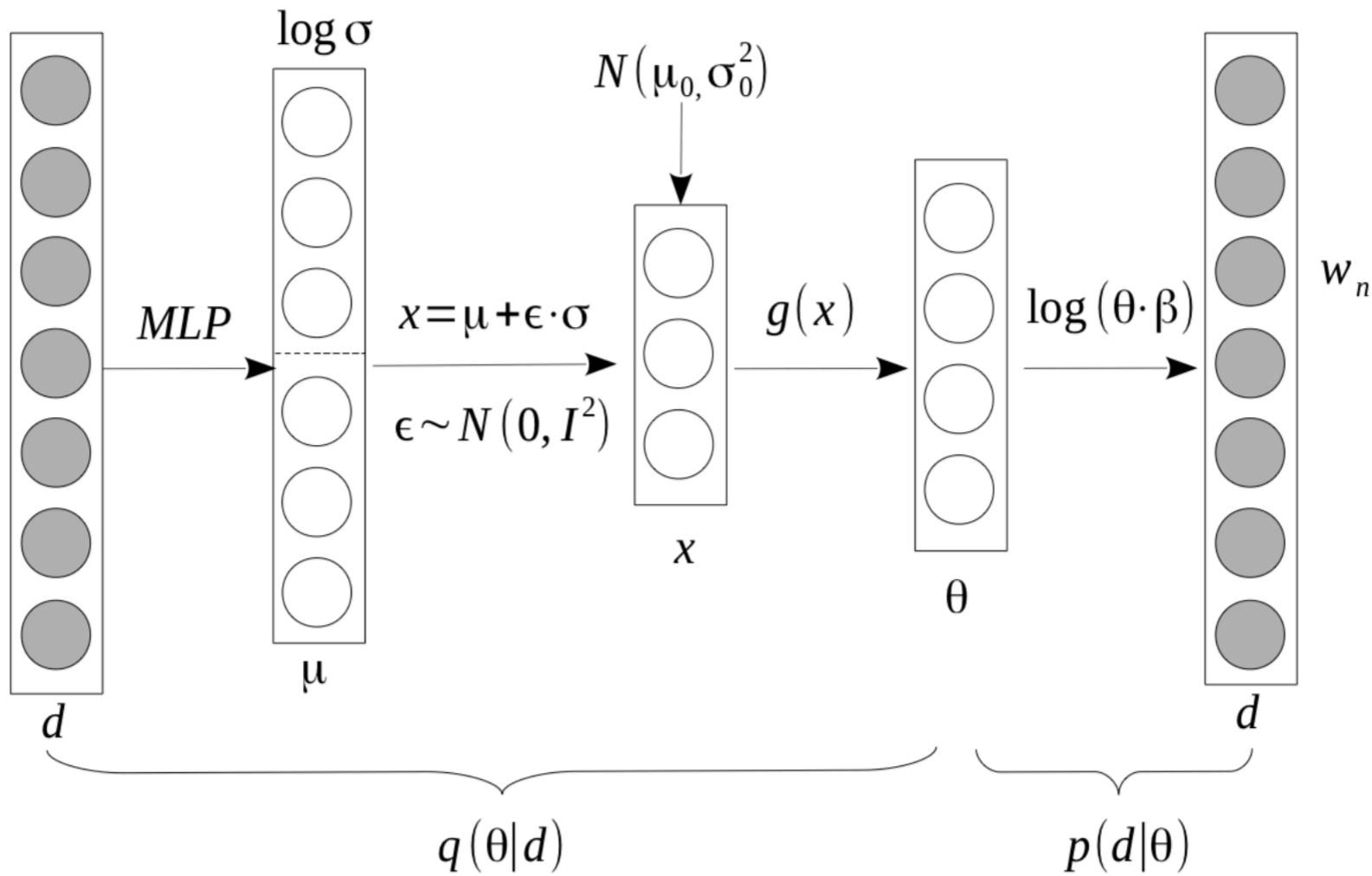


Variational distribution \Rightarrow Optimization problem

$$q(\beta, \theta, z) = \prod_{k=1}^K q(\beta_k | \lambda_k) \prod_{d=1}^D q(\theta_d | \gamma_d) \prod_{n=1}^N q(z_{d,n} | \phi_{d,n}) \approx P(\varphi, \theta, z | \mathbf{w})$$

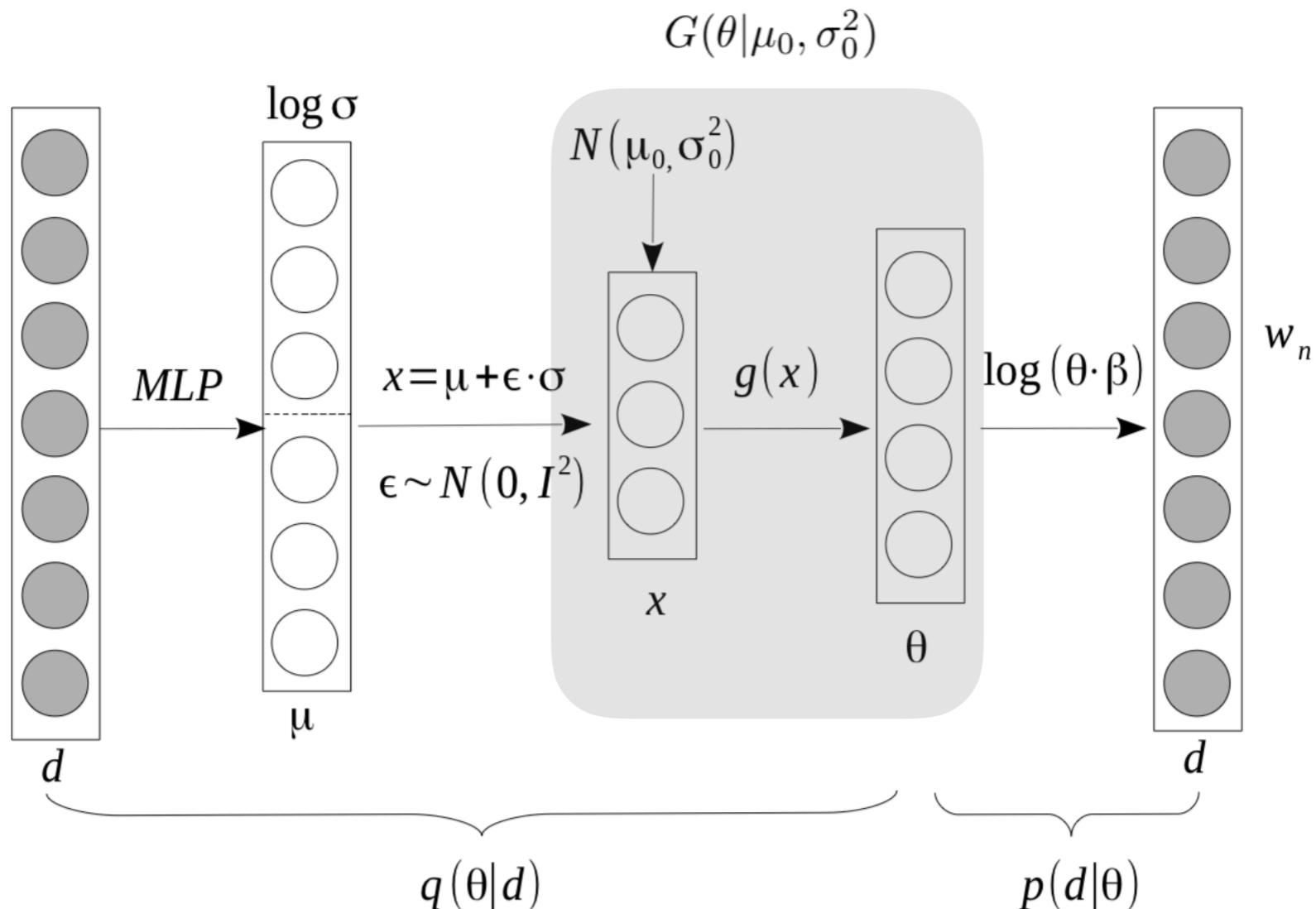


Neural Topic Modeling (Gaussian prior)



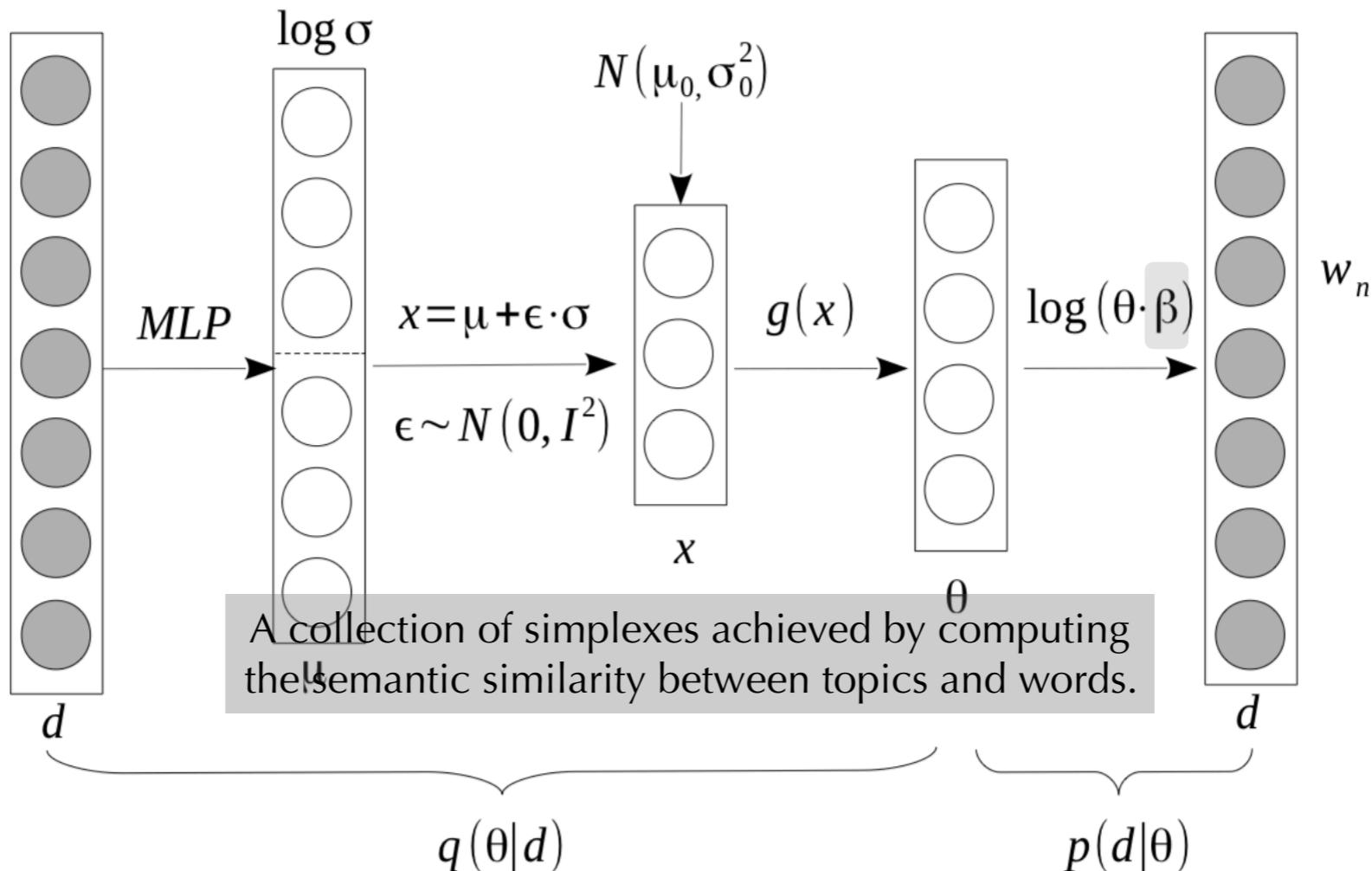
**Network structure of the inference model $q(\theta | d)$,
and of the generative model $p(d | \theta)$.**

Neural Topic Modeling (Gaussian prior)



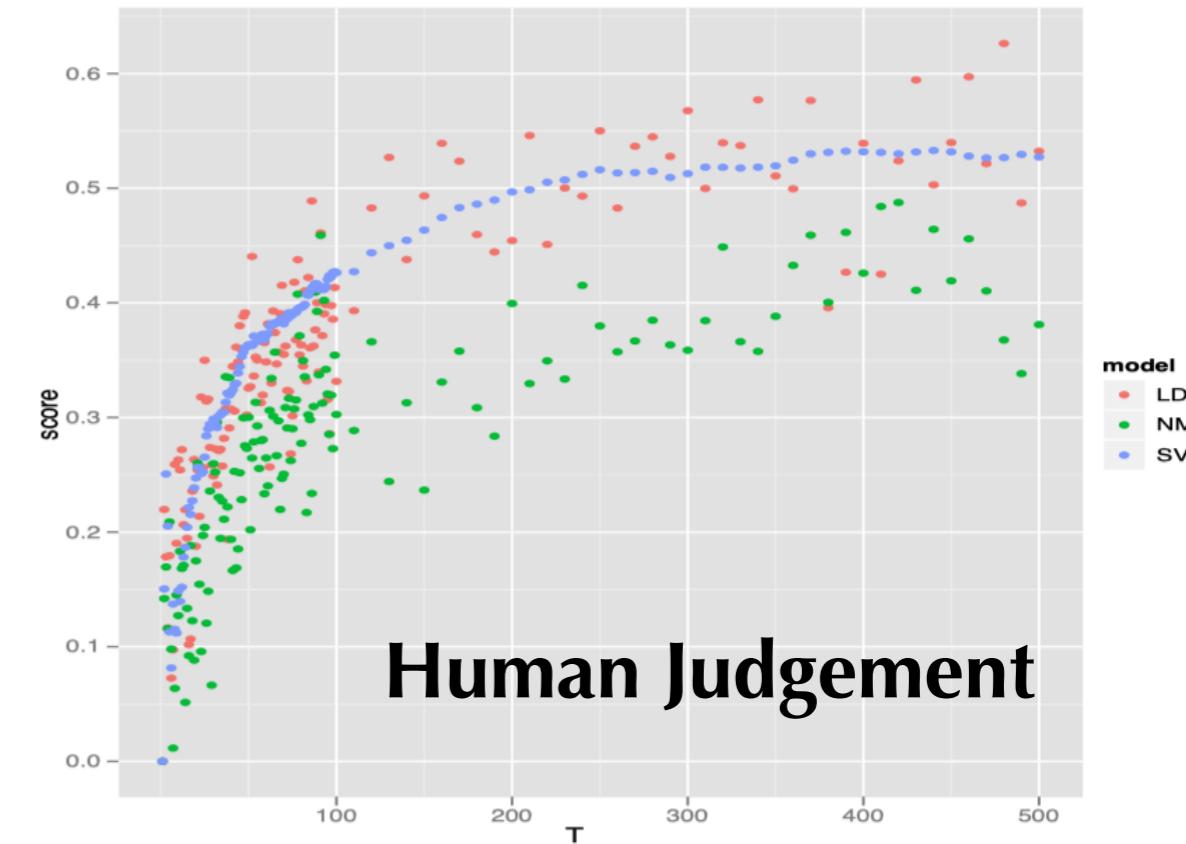
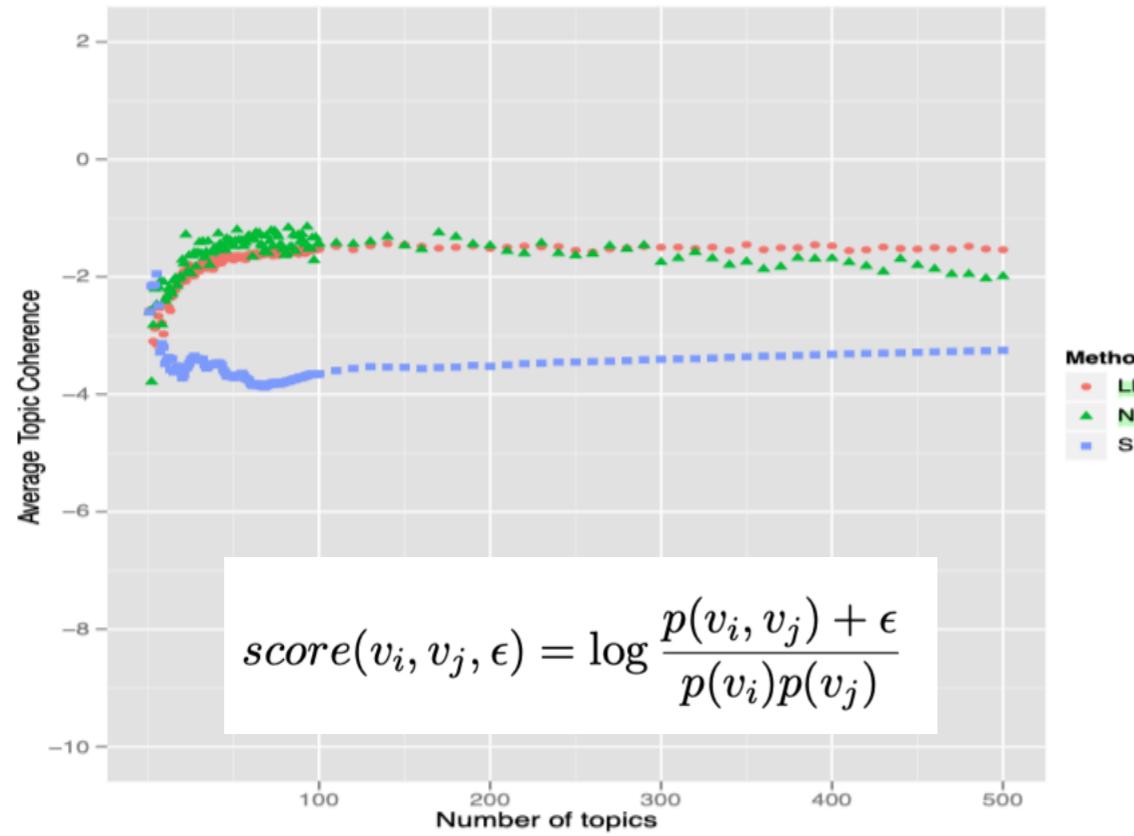
**Network structure of the inference model $q(\theta | d)$,
and of the generative model $p(d | \theta)$.**

Neural Topic Modeling (Gaussian prior)



**Network structure of the inference model $q(\theta | d)$,
and of the generative model $p(d | \theta)$.**

Comparison - SVD, NMF, LDA

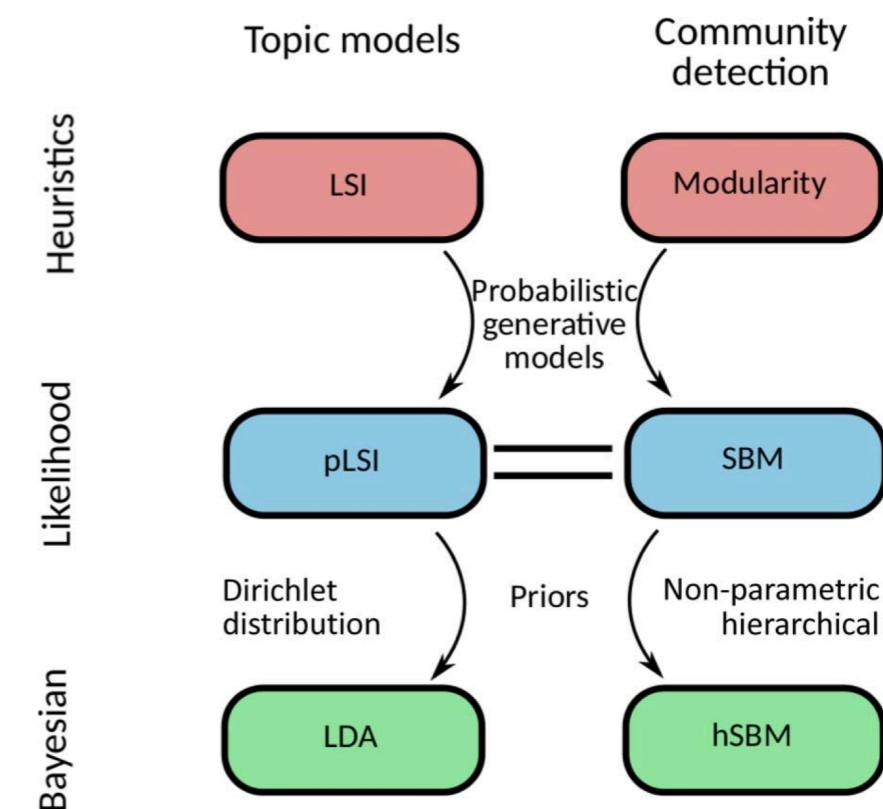
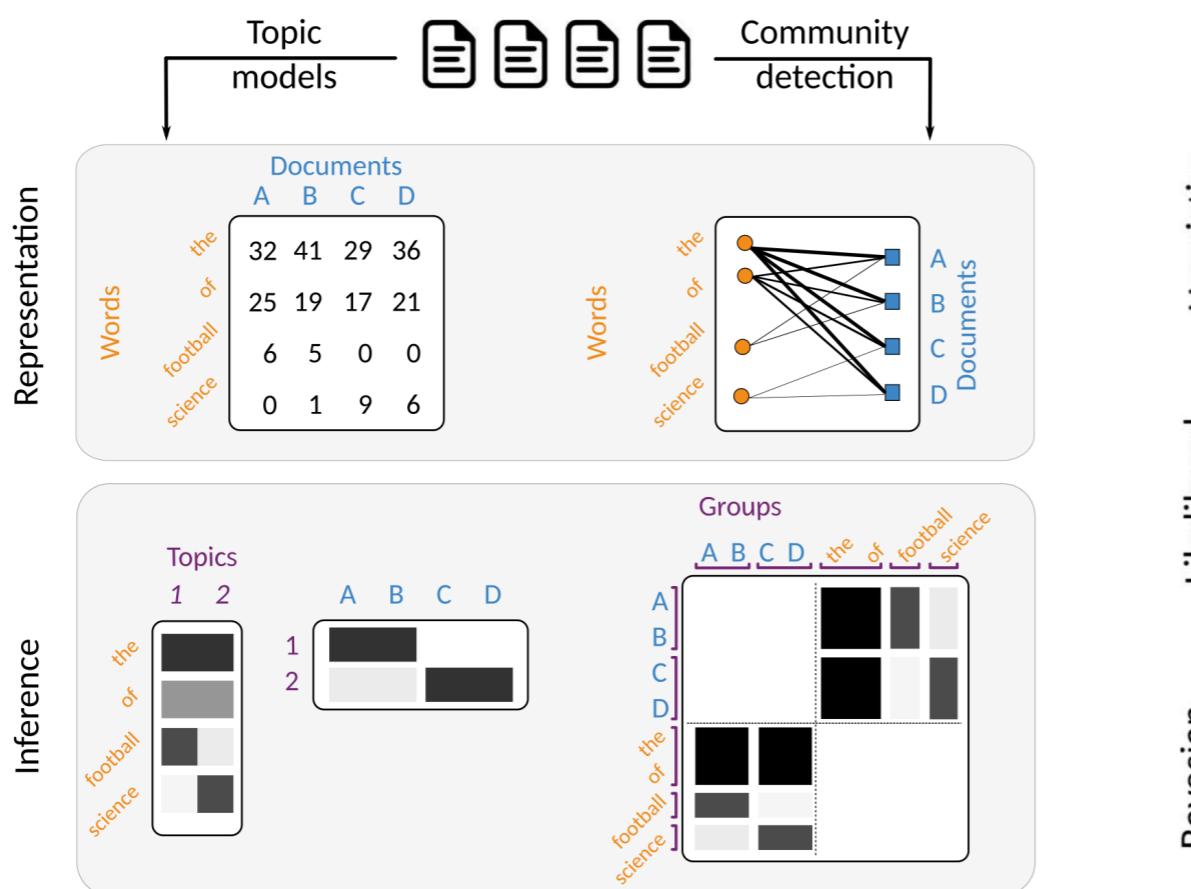


- **SVD fails to form individual topics that aggregate similar words.**
- Conversely, both NMF and LDA learn concise and coherent topics and achieved similar performance.
- **NMF learns more incoherent topics than LDA and SVD.**

Beyond Probabilistic Generative Models

Main Drawbacks of LDA:

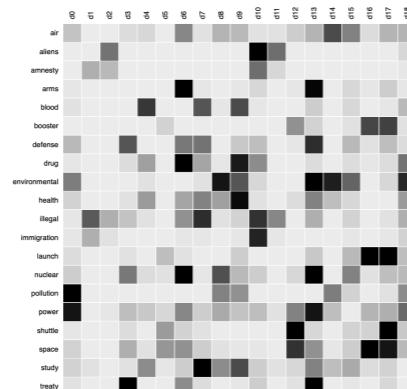
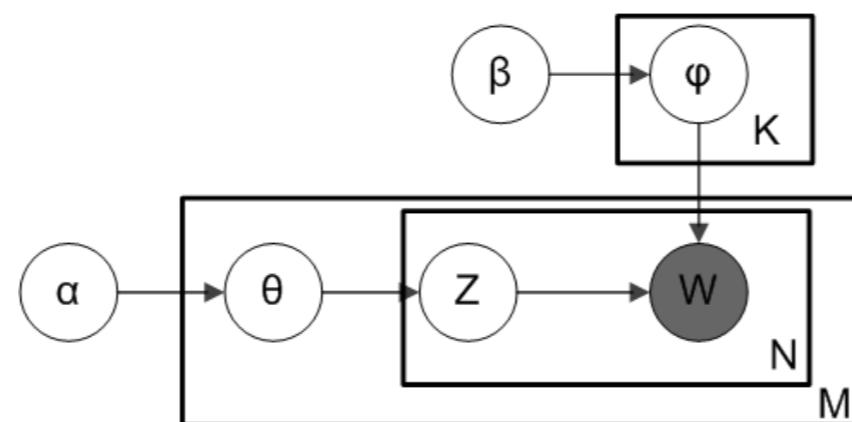
- lacks an intrinsic methodology to choose the **number of topics** and contains a **large number of free parameters** that can cause overfitting.
- no **justification for the use of the Dirichlet prior** in the model formulation besides mathematical convenience.



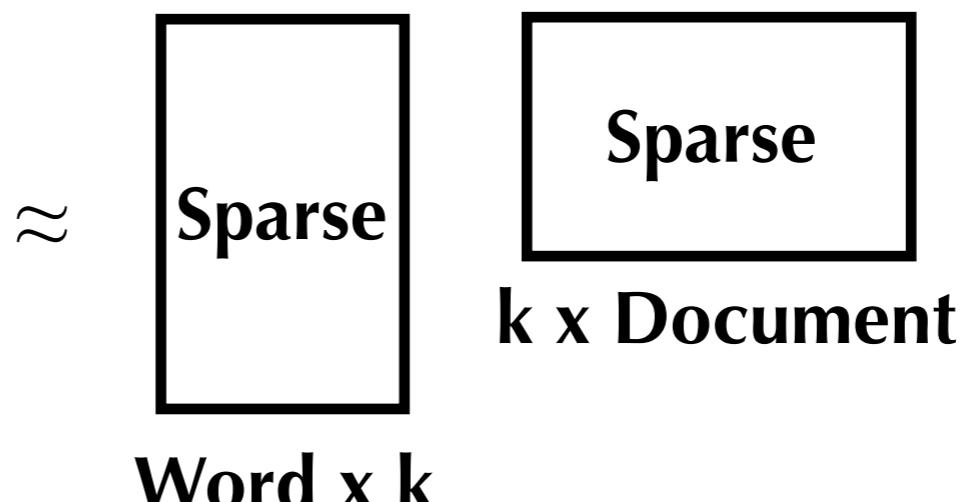
Beyond Probabilistic Generative Models

However, LDA is charming because

- We have α, β to control doc-topic, topic-word sparsity

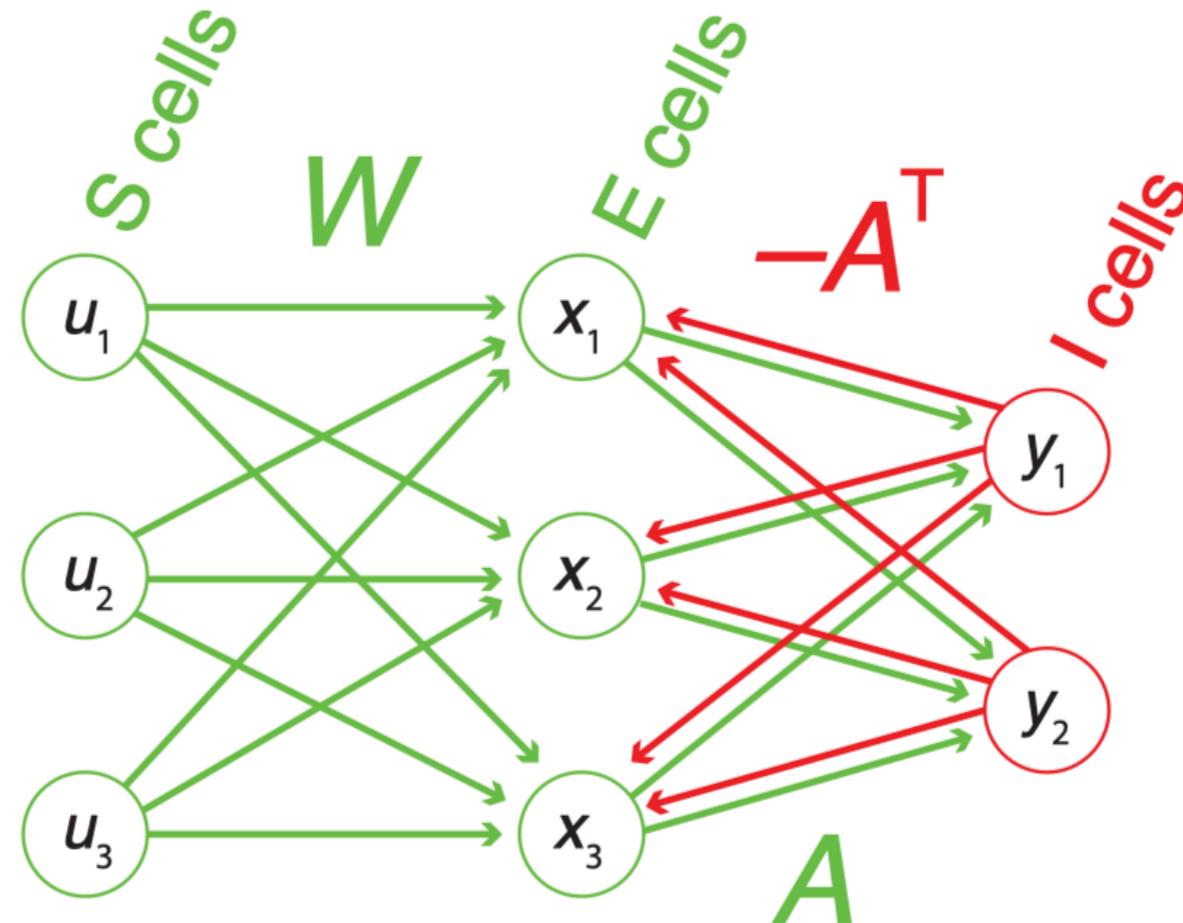


Word x Document
Frequency Matrix



Less sparse - more general topics
More sparse - more representative topics

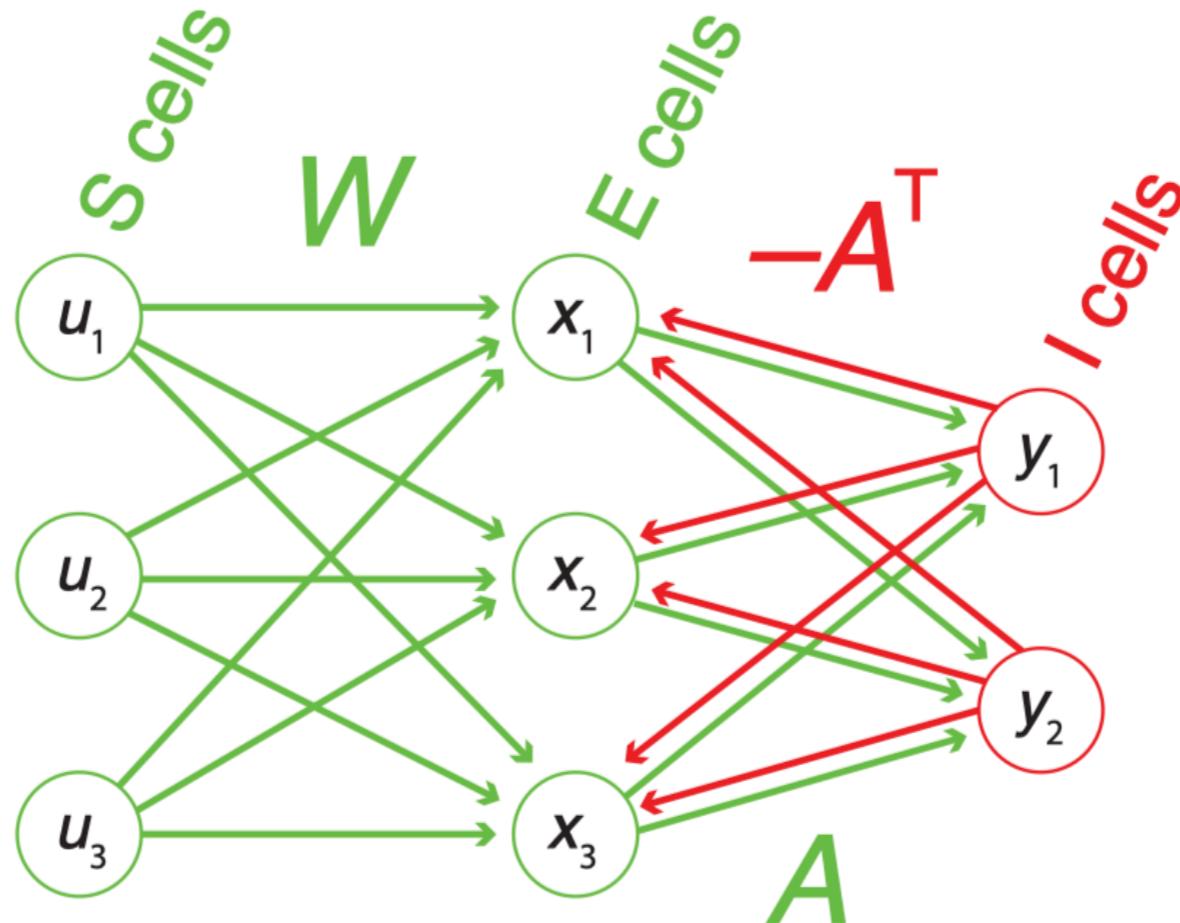
Disynaptic Neural Networks



Excitatory-inhibitory net respecting Dale's Law.

- The $S \rightarrow E$ connection from u_a to x_i has strength W_{ia} .
- The E neurons x_i are reciprocally coupled with the I neurons y_α .
- The $E \rightarrow I$ connection from x_i to y_α has strength $A_{\alpha i}$
- The $I \rightarrow E$ connection from y_α to x_i has strength $-A_{\alpha i}$.

Disynaptic Neural Networks

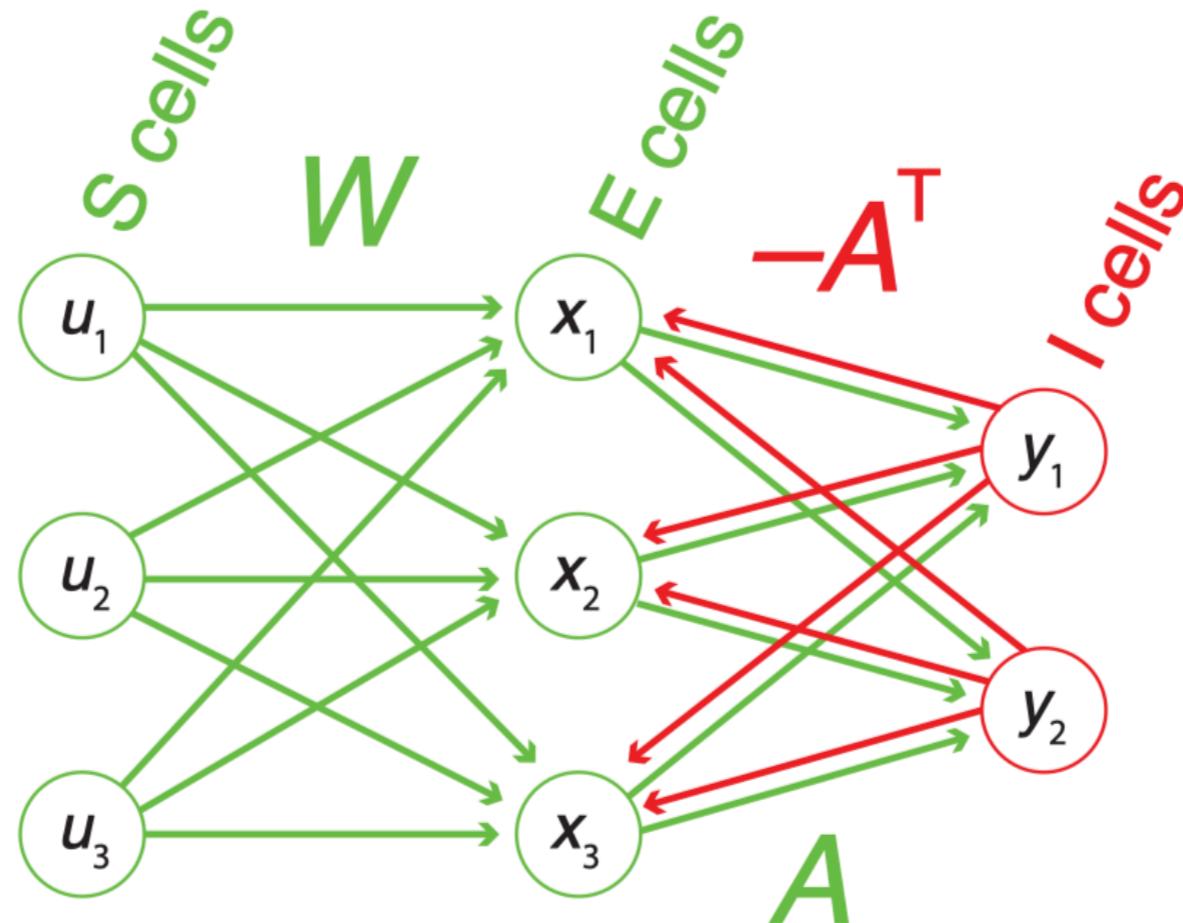


Activity dynamics:

$$x_i := \left[(1 - dt)x_i + dt \lambda_i^{-1} \left(\sum_{a=1}^n W_{ia} u_a - \sum_{\alpha=1}^r y_\alpha A_{\alpha i} \right) \right]^+$$

$$y_\alpha = \sum_{i=1}^m A_{\alpha i} x_i$$

Disynaptic Neural Networks



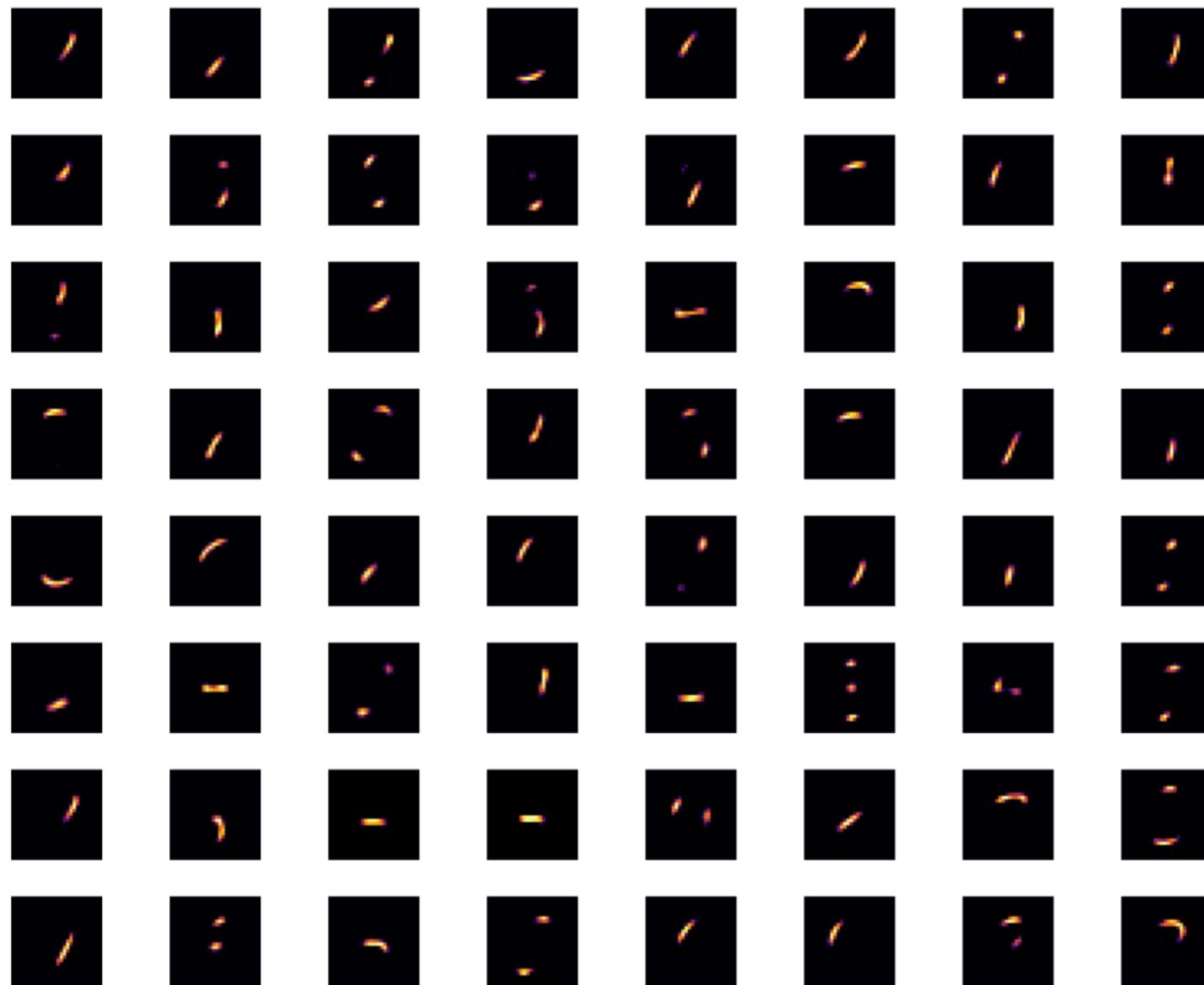
Learning rules (Hebbian / Anti-hebbian learning):

$$\Delta W_{ia} \propto x_i u_a - \gamma W_{ia} - \kappa \sum_b W_{ib}$$

$$\Delta A_{\alpha j} \propto y_\alpha x_j - (q^2 - p^2) A_{\alpha j} - p^2 \sum_i A_{\alpha i}$$

$$\Delta \lambda_i \propto x_i^2 - q^2$$

Disynaptic Neural Networks



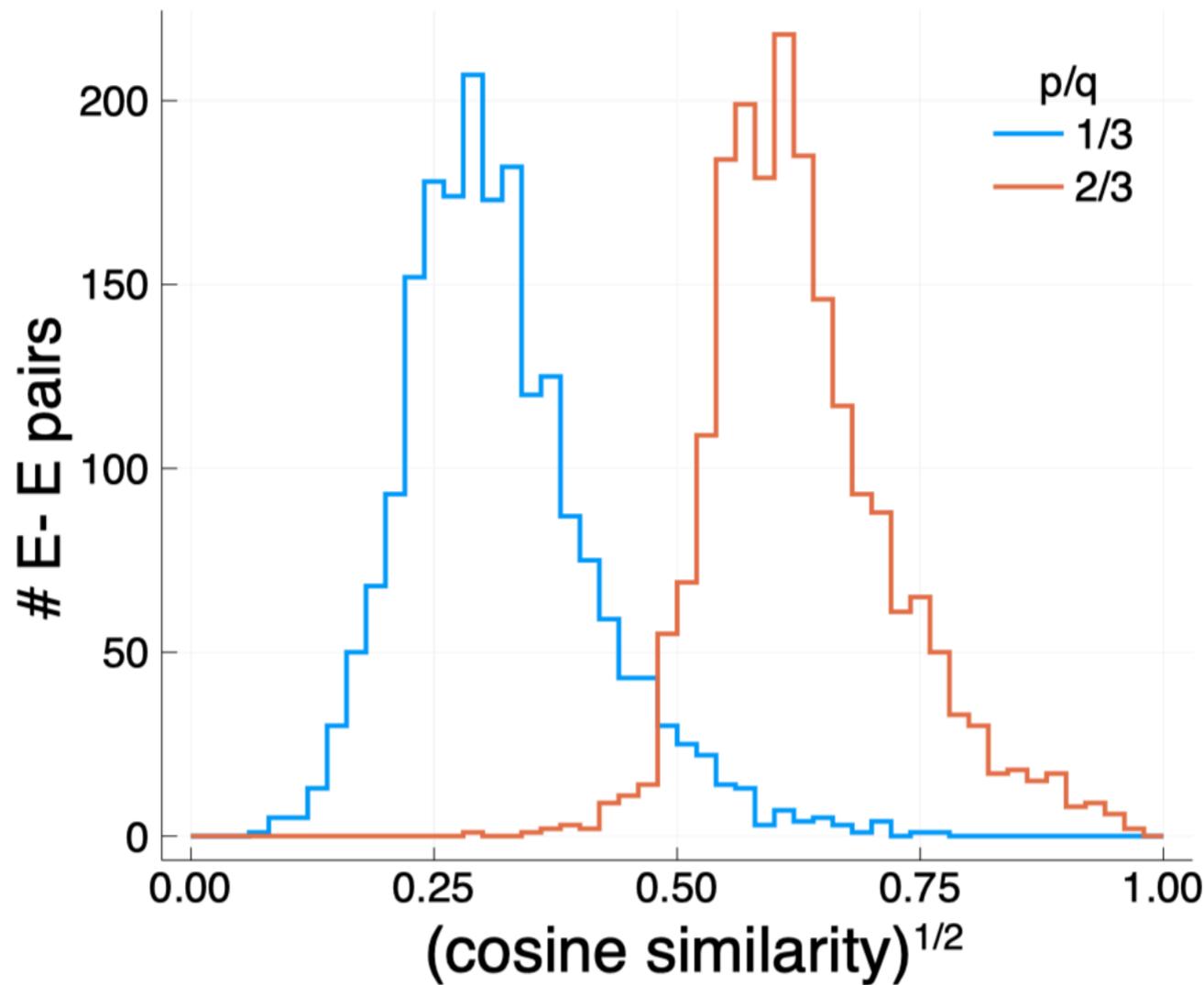
**For each E neuron, the convergent S → E connections
constitute a sensory feature learned from the stimuli.
($\gamma/\kappa = 5$)**

Disynaptic Neural Networks



For each E neuron, the convergent S → E connections
constitute a sensory feature learned from the stimuli.
 $(\gamma/\kappa = 50)$

Disynaptic Neural Networks



The ratio p/q controls the degree of decorrelation.

Correlation Game

Define the goal of unsupervised learning
as the constrained optimization:

$$\max_{X \geq 0} \Phi^* \left(\frac{XU^\top}{T} \right) \text{ subject to copositivity of } D - \frac{XX^\top}{T}$$

Correlation Game

Define the goal of unsupervised learning
as the constrained optimization:

$$\max_{X \geq 0} \Phi^* \left(\frac{XU^\top}{T} \right) \text{ subject to copositivity of } D - \frac{XX^\top}{T}$$

where $\Phi^*(C) = \max_{W \geq 0} \left\{ \sum_{ia} W_{ia} C_{ia} - \Phi(W) \right\}$, $D_{ij} = \begin{cases} q^2, & i = j, \\ p^2, & i \neq j \end{cases}$.

Correlation Game

Define the goal of unsupervised learning
as the constrained optimization:

$$\max_{X \geq 0} \Phi^* \left(\frac{XU^\top}{T} \right) \text{ subject to copositivity of } D - \frac{XX^\top}{T}$$

where $\Phi^*(C) = \max_{W \geq 0} \left\{ \sum_{ia} W_{ia} C_{ia} - \Phi(W) \right\}, \quad D_{ij} = \begin{cases} q^2, & i = j, \\ p^2, & i \neq j \end{cases}.$

convex conjugate of regularizer

$$\Phi(W) = \frac{\gamma}{2} \sum_{ia} W_{ia}^2 + \frac{\kappa}{2} \sum_i \left(\sum_a W_{ia} \right)^2$$

Correlation Game

Define the goal of unsupervised learning
as the constrained optimization:

$$\max_{X \geq 0} \Phi^* \left(\frac{XU^\top}{T} \right) \text{ subject to copositivity of } D - \frac{XX^\top}{T}$$

where $\Phi^*(C) = \max_{W \geq 0} \left\{ \sum_{ia} W_{ia} C_{ia} - \Phi(W) \right\}$, $D_{ij} = \begin{cases} q^2, & i = j, \\ p^2, & i \neq j \end{cases}$.

$$\langle x_i^2 \rangle \leq D_{ii} \quad \langle x_i x_j \rangle \leq D_{ij}$$

constraints on correlation

Correlation Game

Define the goal of unsupervised learning
as the constrained optimization:

$$\max_{X \geq 0} \Phi^* \left(\frac{XU^\top}{T} \right) \text{ subject to copositivity of } D - \frac{XX^\top}{T}$$



Introducing Lagrange multipliers \mathbf{A} and Λ :

$$\min_{A, \Lambda \geq 0} \max_{X \geq 0} \left\{ \Phi^* \left(\frac{XU^\top}{T} \right) + \frac{1}{2} \text{Tr} \left(D - \frac{XX^\top}{T} \right) (A^\top A + \Lambda) \right\}$$

Correlation Game

- Original Problem(s): (1) \equiv (2) \equiv (3)

$$\max_{\substack{X \geq 0 \\ L \geq 0}} \min_{W \geq 0} \left\{ \frac{1}{T} \operatorname{Tr} W^\top X U^\top - \Phi(W) + \frac{1}{2} \left(D - \frac{XX^\top}{T} \right) L \right\} \quad (1)$$

$$\max_{\substack{X \geq 0 \\ W \geq 0 \\ L \geq 0}} \min_{W \geq 0} \left\{ \frac{1}{T} \operatorname{Tr} W^\top X U^\top - \Phi(W) + \frac{1}{2} \left(D - \frac{XX^\top}{T} \right) L \right\} \quad (2)$$

$$\max_{\substack{W \geq 0 \\ X \geq 0 \\ L \geq 0}} \min_{W \geq 0} \left\{ \frac{1}{T} \operatorname{Tr} W^\top X U^\top - \Phi(W) + \frac{1}{2} \left(D - \frac{XX^\top}{T} \right) L \right\} \quad (3)$$

- Lagrangian Dual Problem(s): (4) \equiv (5)

$$\min_{\substack{L \geq 0 \\ X \geq 0 \\ W \geq 0}} \max_{\substack{W \geq 0 \\ L \geq 0}} \left\{ \frac{1}{T} \operatorname{Tr} W^\top X U^\top - \Phi(W) + \frac{1}{2} \left(D - \frac{XX^\top}{T} \right) L \right\} \quad (4)$$

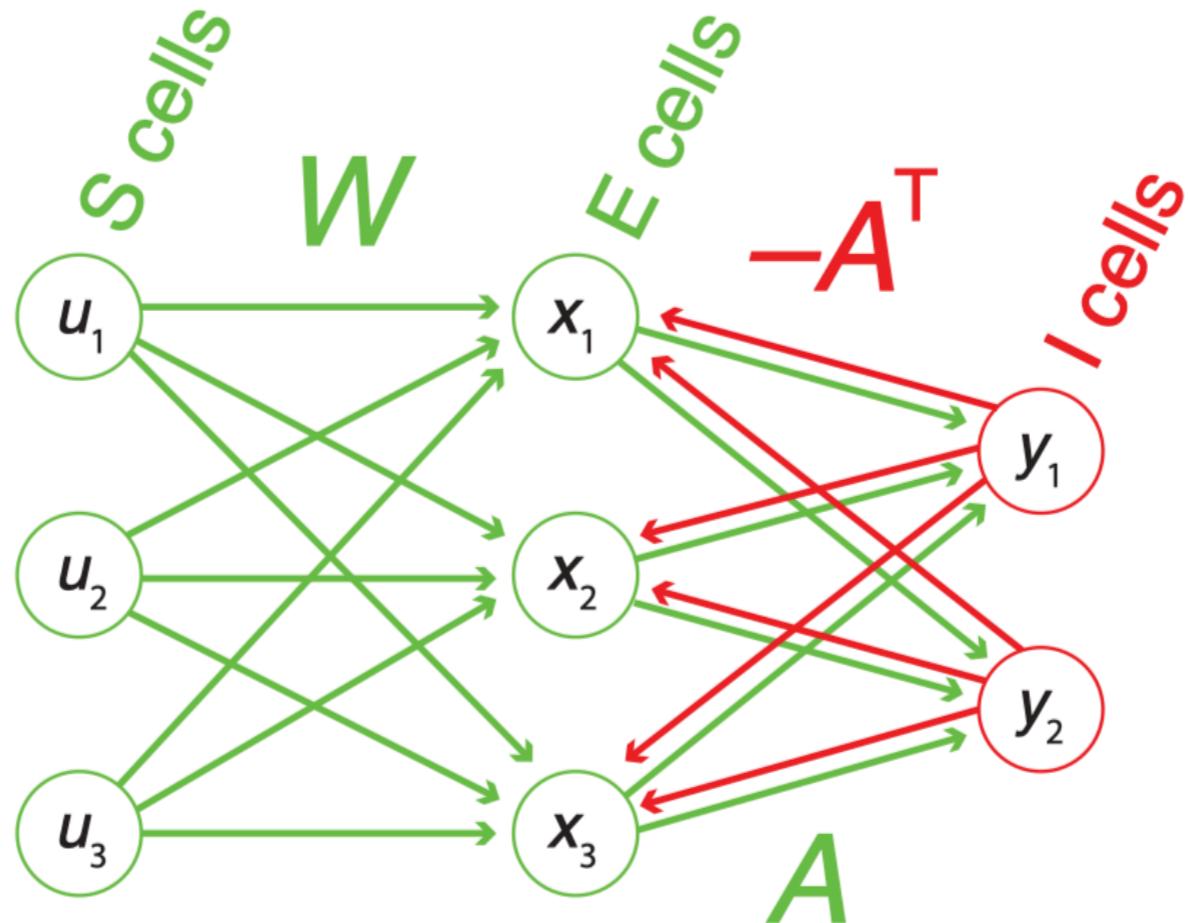
$$\min_{\substack{L \geq 0 \\ W \geq 0 \\ X \geq 0}} \max_{\substack{W \geq 0 \\ L \geq 0}} \left\{ \frac{1}{T} \operatorname{Tr} W^\top X U^\top - \Phi(W) + \frac{1}{2} \left(D - \frac{XX^\top}{T} \right) L \right\} \quad (5)$$

- Neural Network Algorithm:

$$\max_{\substack{W \geq 0 \\ L \geq 0 \\ X \geq 0}} \min_{W \geq 0} \left\{ \frac{1}{T} \operatorname{Tr} W^\top X U^\top - \Phi(W) + \frac{1}{2} \left(D - \frac{XX^\top}{T} \right) L \right\} \quad (6)$$

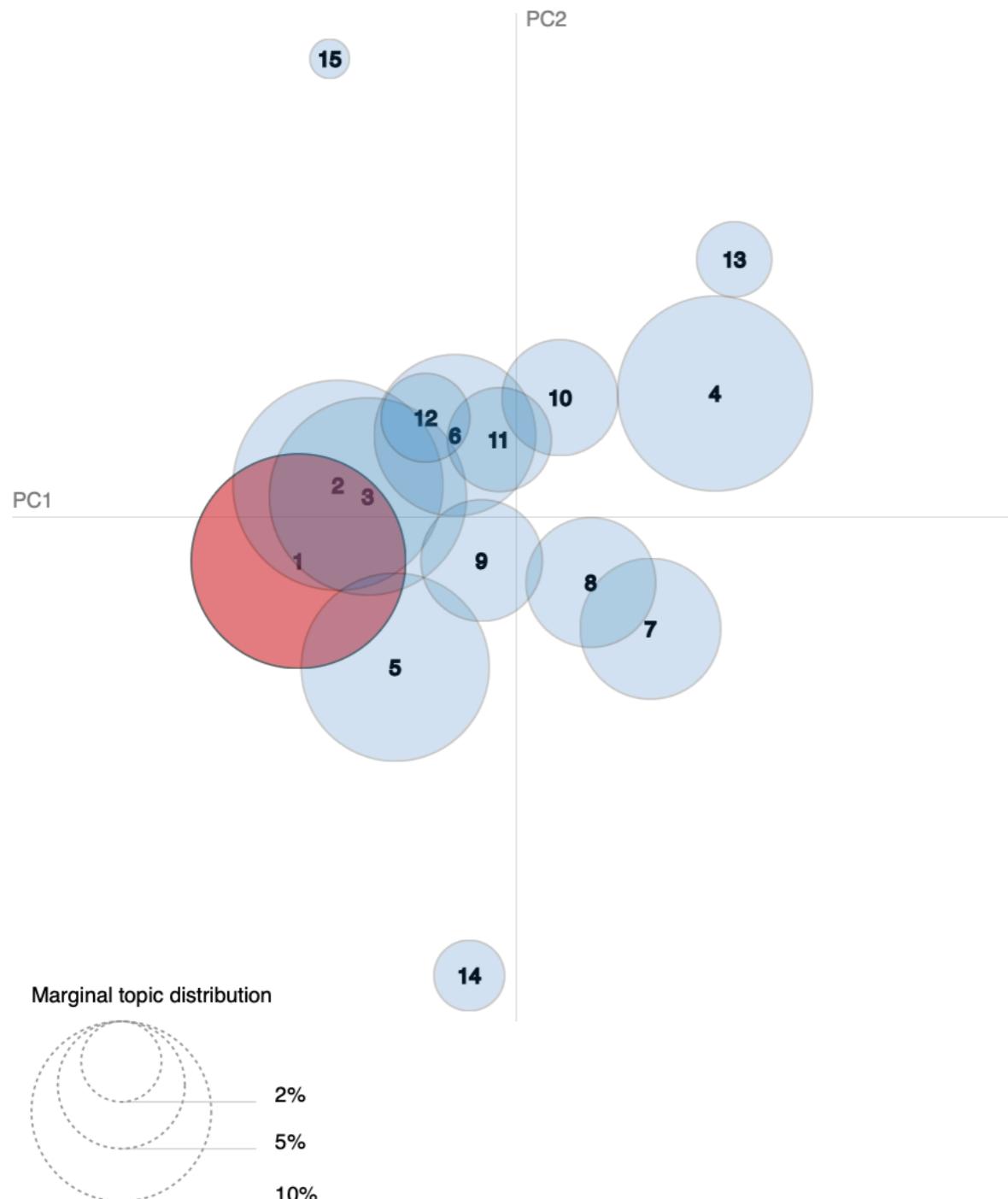
$$(1) \equiv (2) \equiv (3) \leq (6) \leq (4) \equiv (5)$$

Disynaptic Neural Topic Modelling?

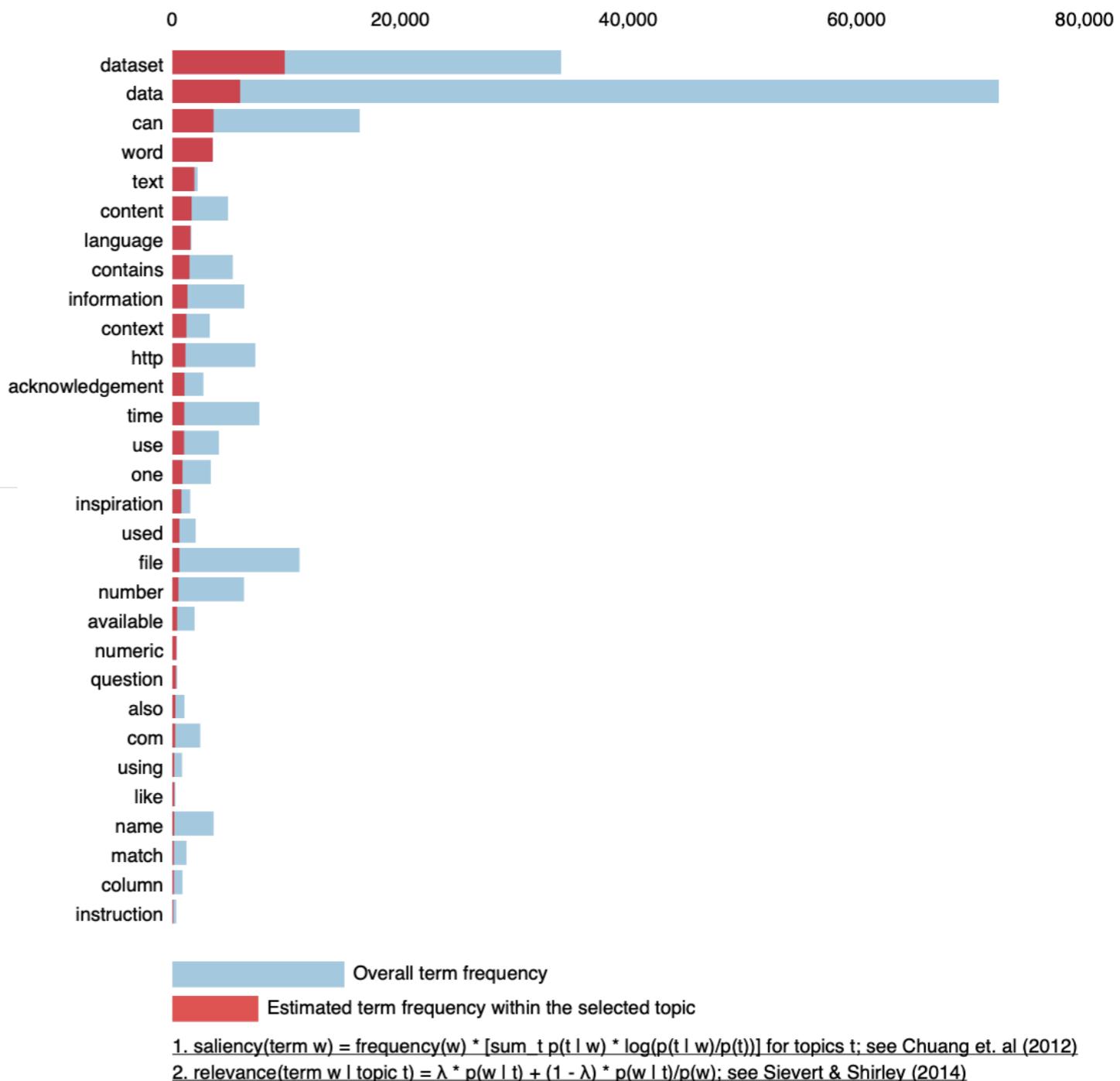


Disynaptic Neural Topic Modelling?

Intertopic Distance Map (via multidimensional scaling)



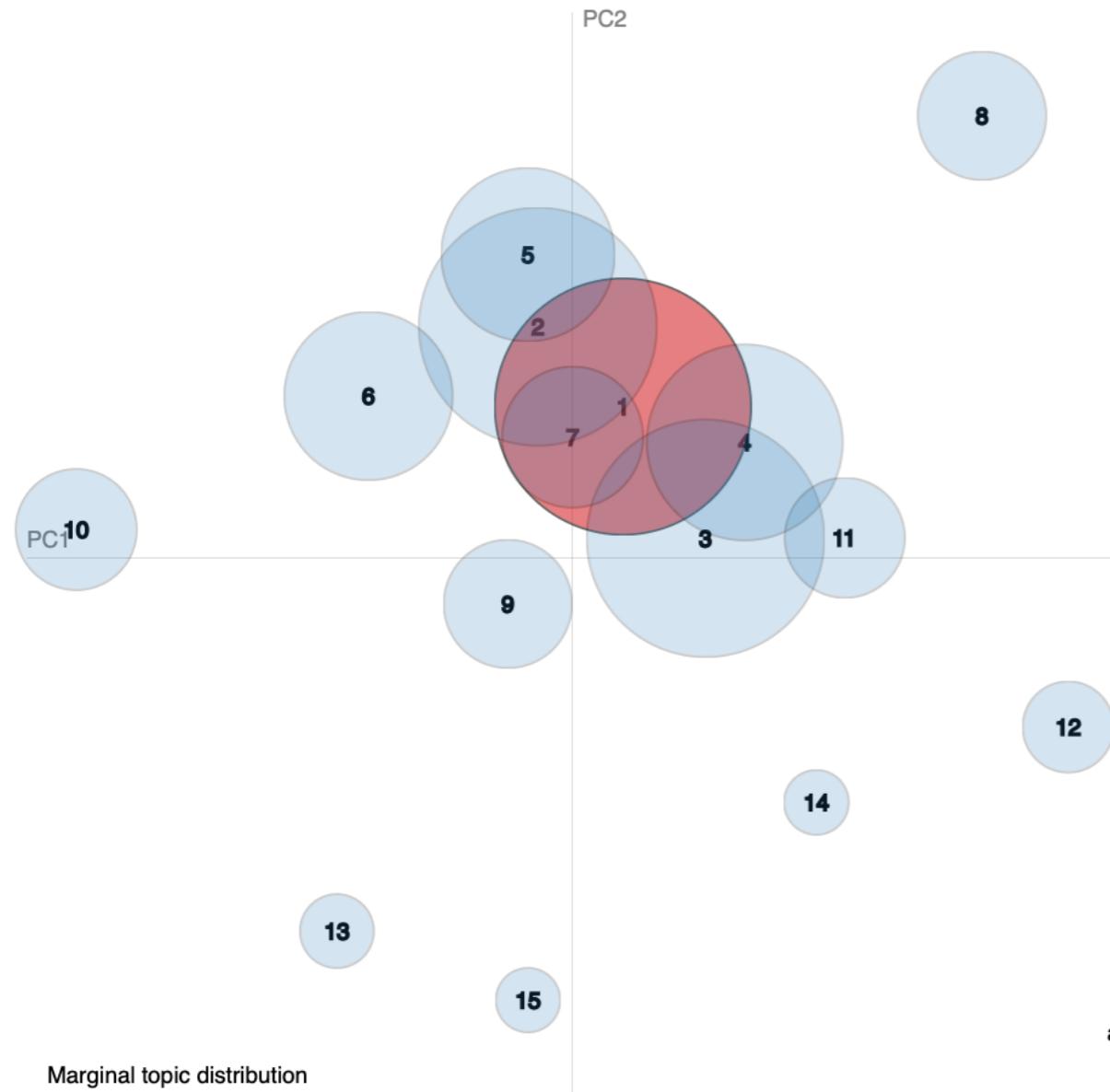
Top-30 Most Relevant Terms for Topic 1 (14.2% of tokens)



$$p/q = 1/3$$

Disynaptic Neural Topic Modelling?

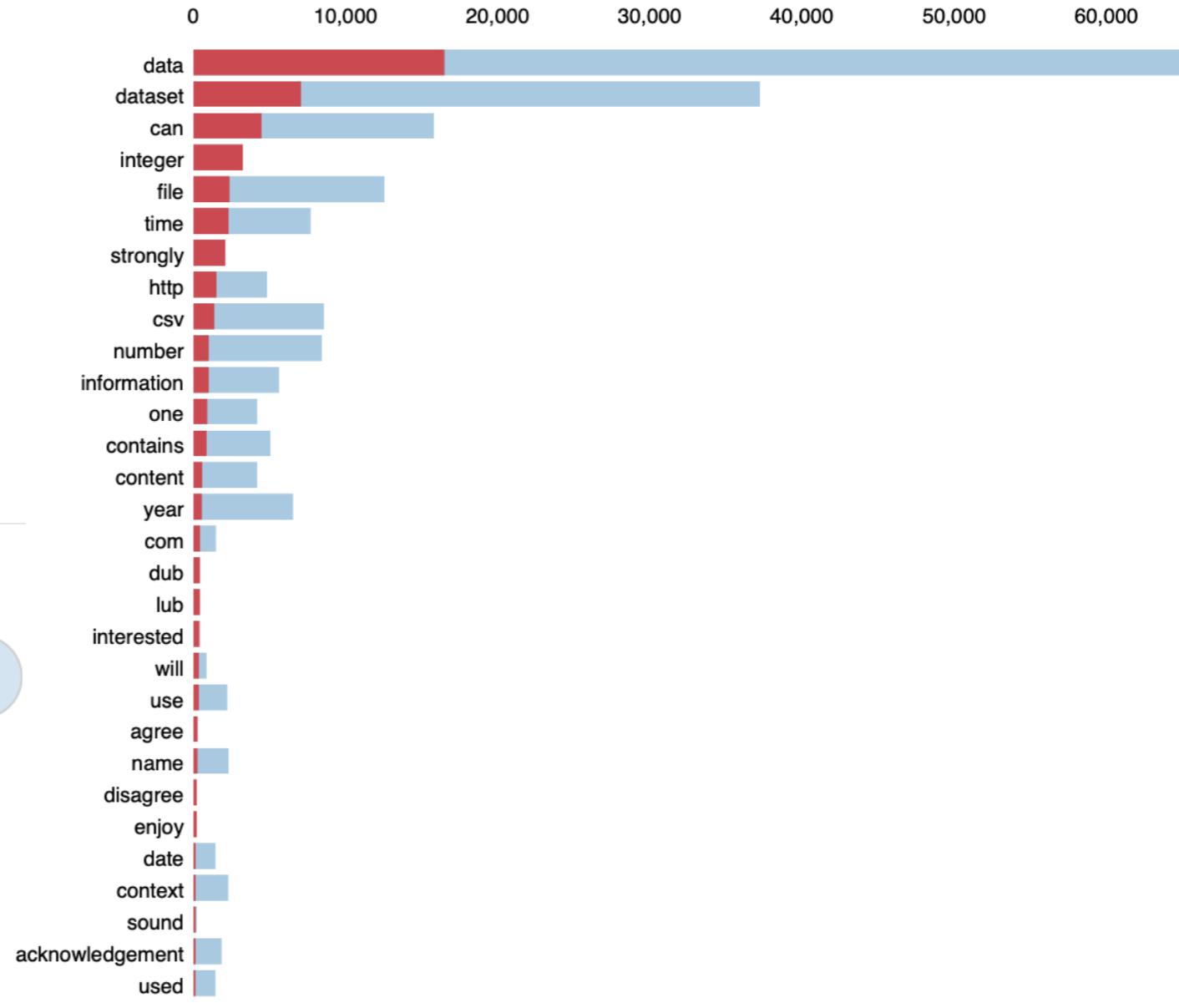
Intertopic Distance Map (via multidimensional scaling)



Marginal topic distribution



Top-30 Most Relevant Terms for Topic 1 (17.3% of tokens)



Overall term frequency

Estimated term frequency within the selected topic

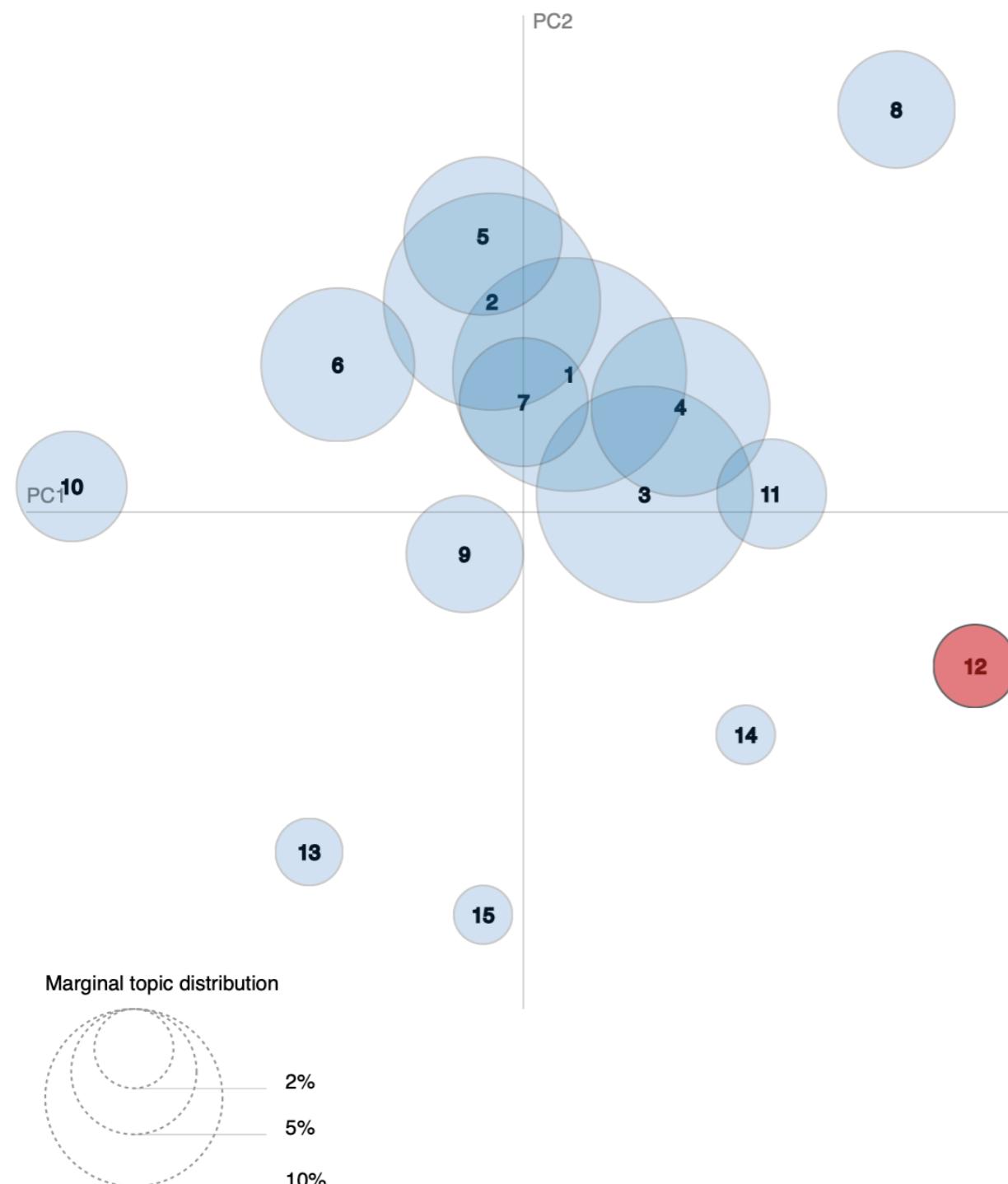
1. saliency(term w) = frequency(w) * [sum_t p(t | w) * log(p(t | w)/p(t))] for topics t; see Chuang et. al (2012)

2. relevance(term w | topic t) = $\lambda * p(w | t) + (1 - \lambda) * p(w | t)/p(w)$; see Sievert & Shirley (2014)

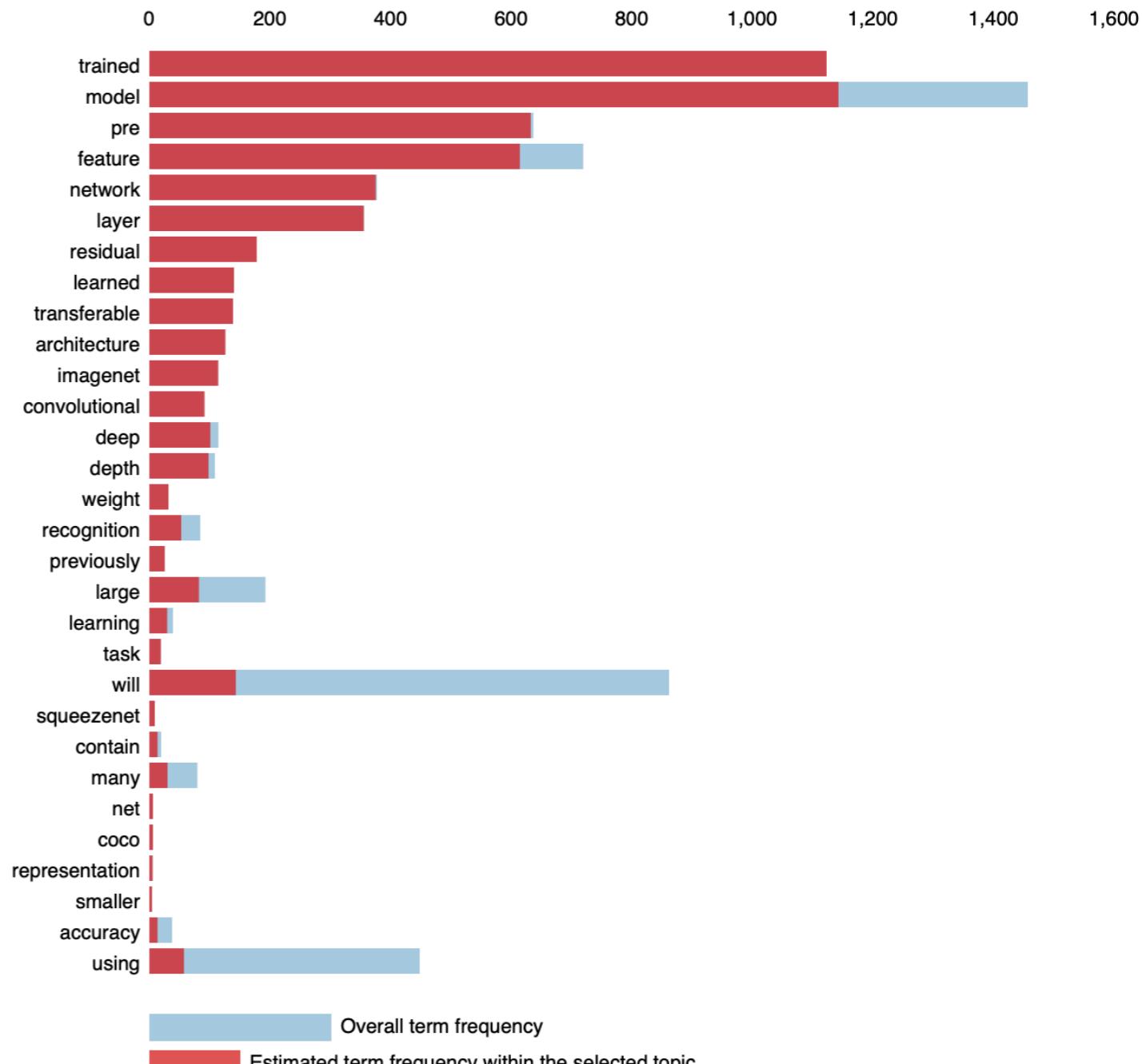
$$p/q = 1/18$$

Disynaptic Neural Topic Modelling?

Intertopic Distance Map (via multidimensional scaling)



Top-30 Most Relevant Terms for Topic 12 (2.2% of tokens)



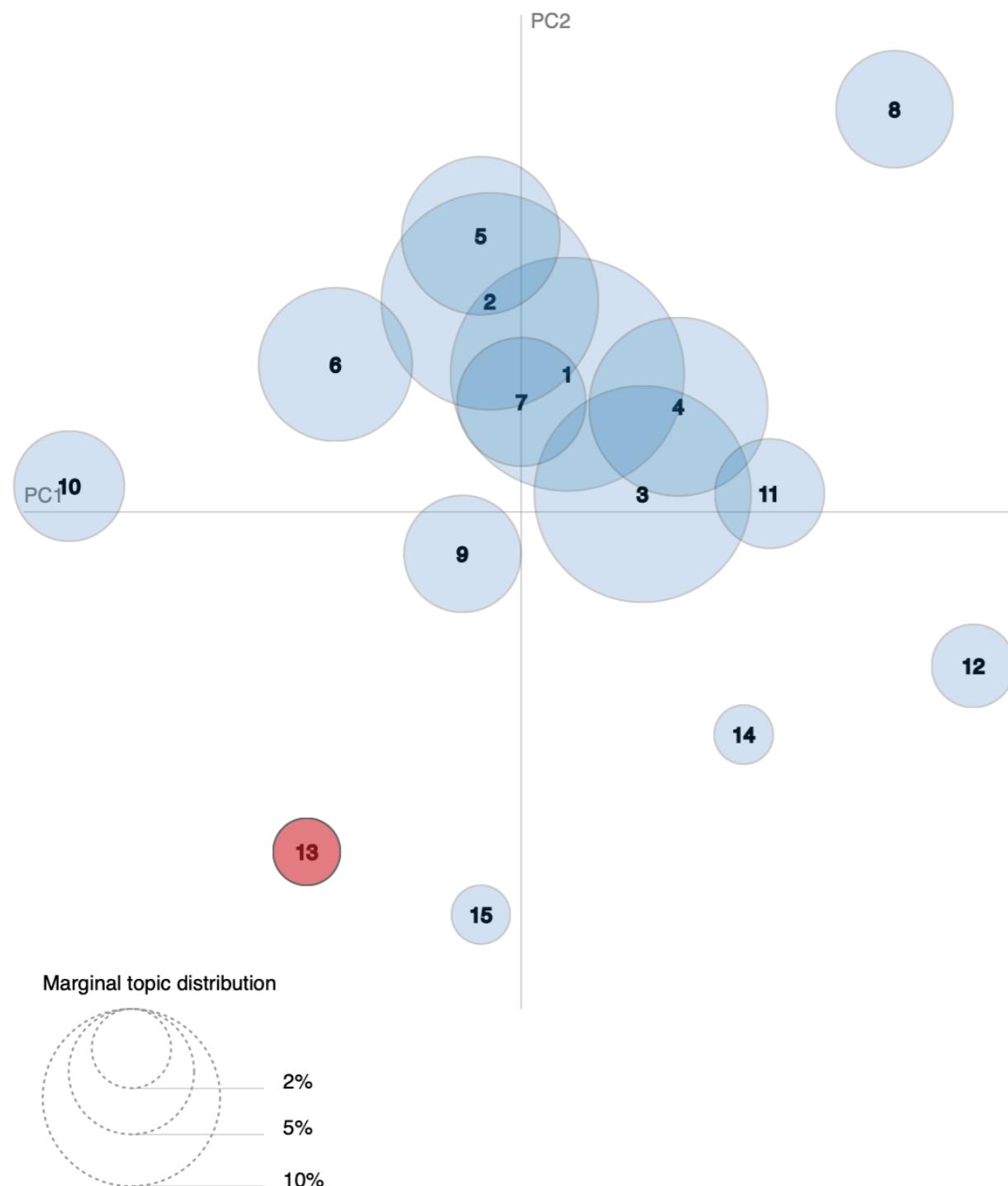
1. saliency(term w) = frequency(w) * [sum_t p(t | w) * log(p(t | w)/p(t))] for topics t; see Chuang et. al (2012)

2. relevance(term w | topic t) = $\lambda * p(w | t) + (1 - \lambda) * p(w | t) / p(w)$; see Sievert & Shirley (2014)

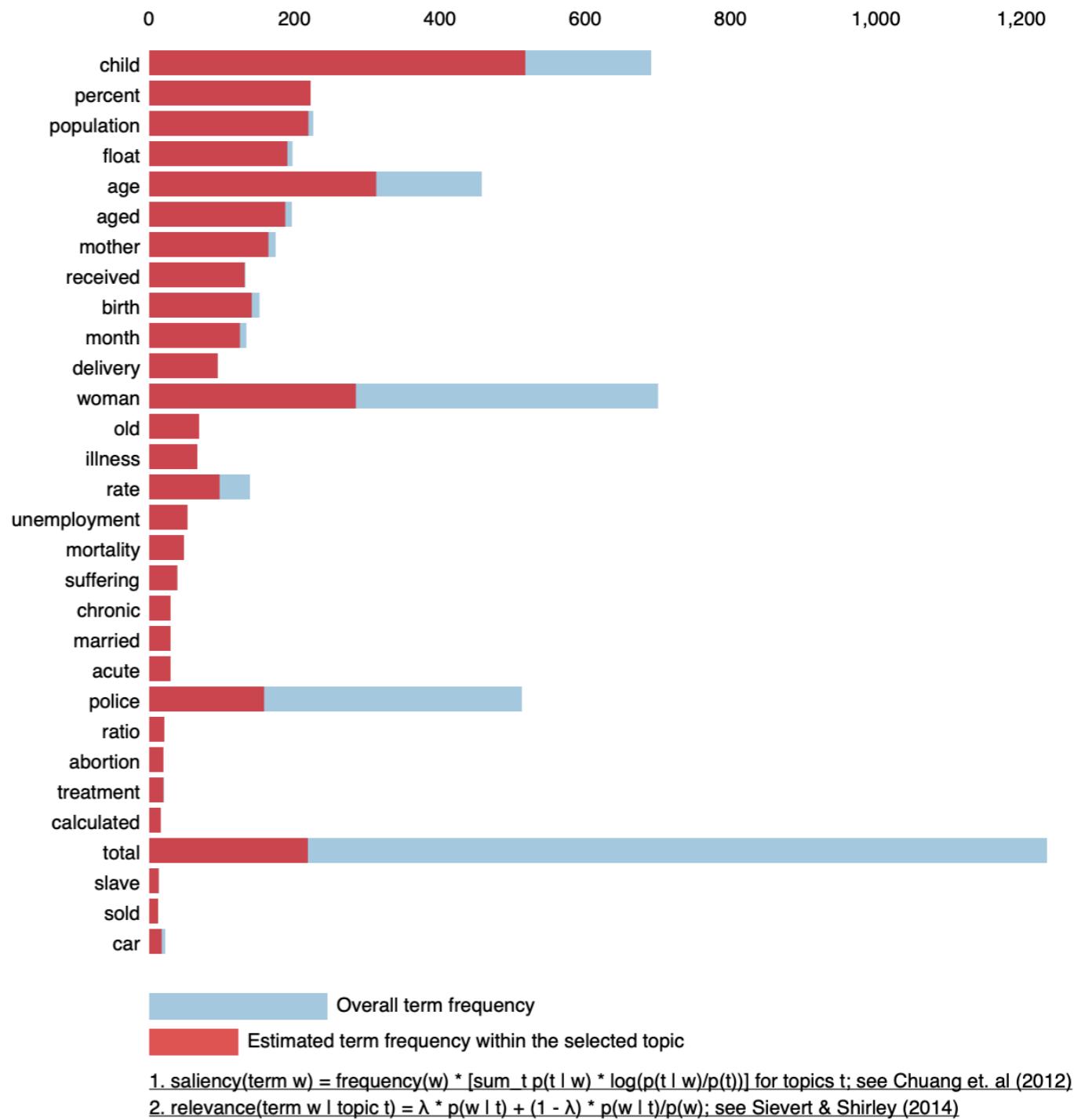
$p/q = 1/18$, topic 12 “training neural nets”

Disynaptic Neural Topic Modelling?

Intertopic Distance Map (via multidimensional scaling)



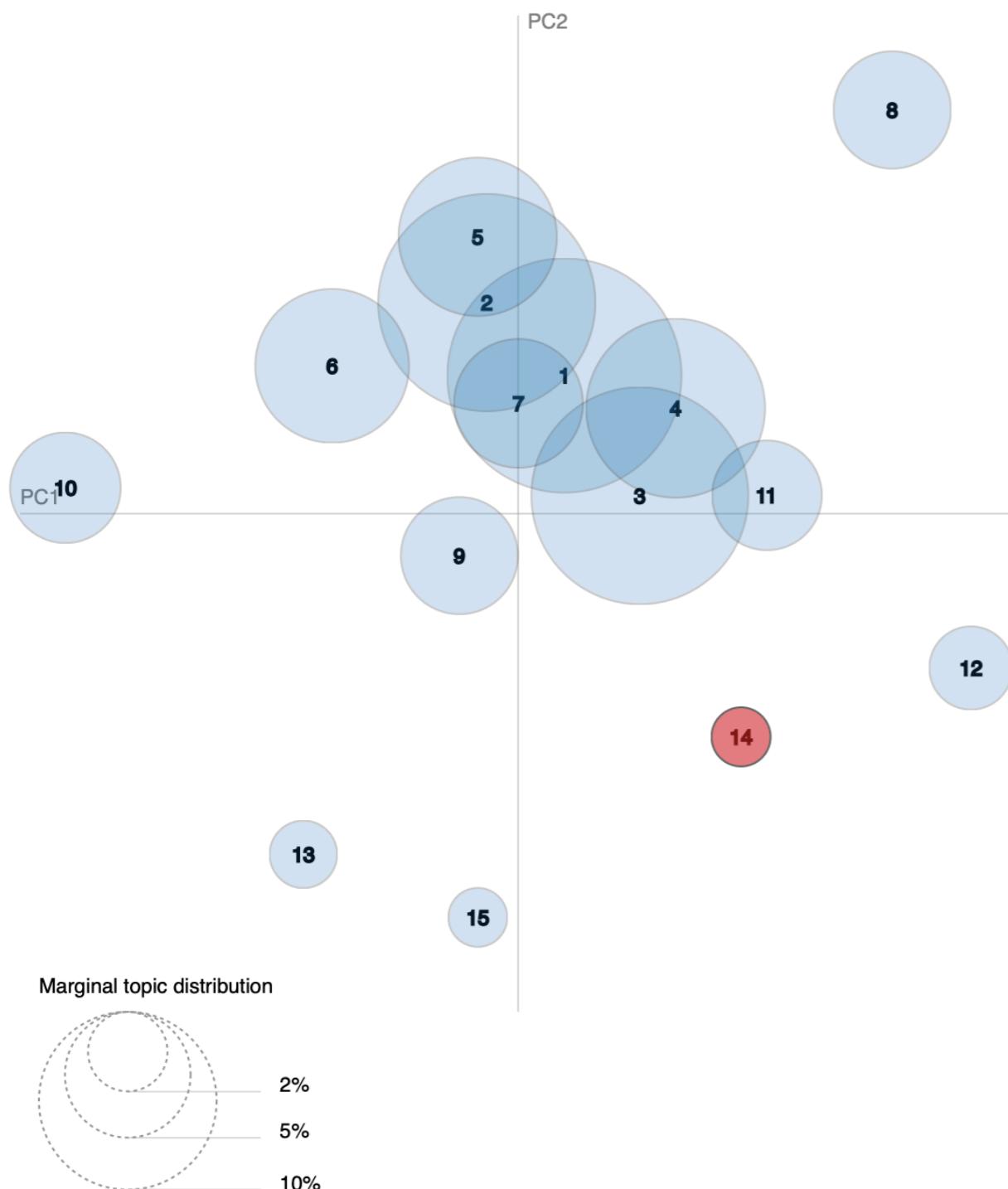
Top-30 Most Relevant Terms for Topic 13 (1.4% of tokens)



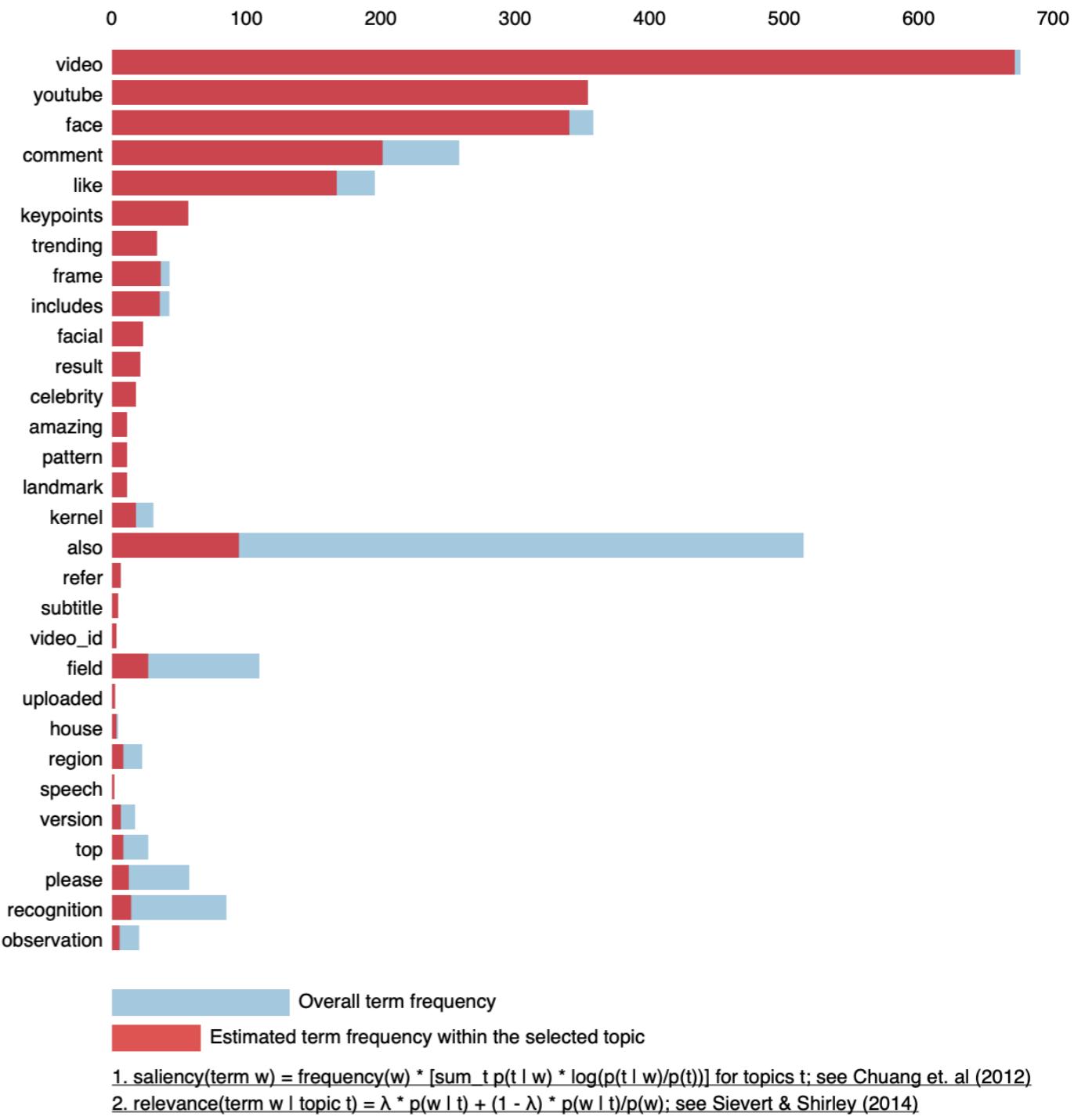
p/q = 1/18, topic 13 “population”

Disynaptic Neural Topic Modelling?

Intertopic Distance Map (via multidimensional scaling)



Top-30 Most Relevant Terms for Topic 14 (1.1% of tokens)



$p/q = 1/18$, topic 14 “video detection”