



UNIVERSITÀ
DEGLI STUDI
DI PADOVA

Network Science

A.Y. 23/24

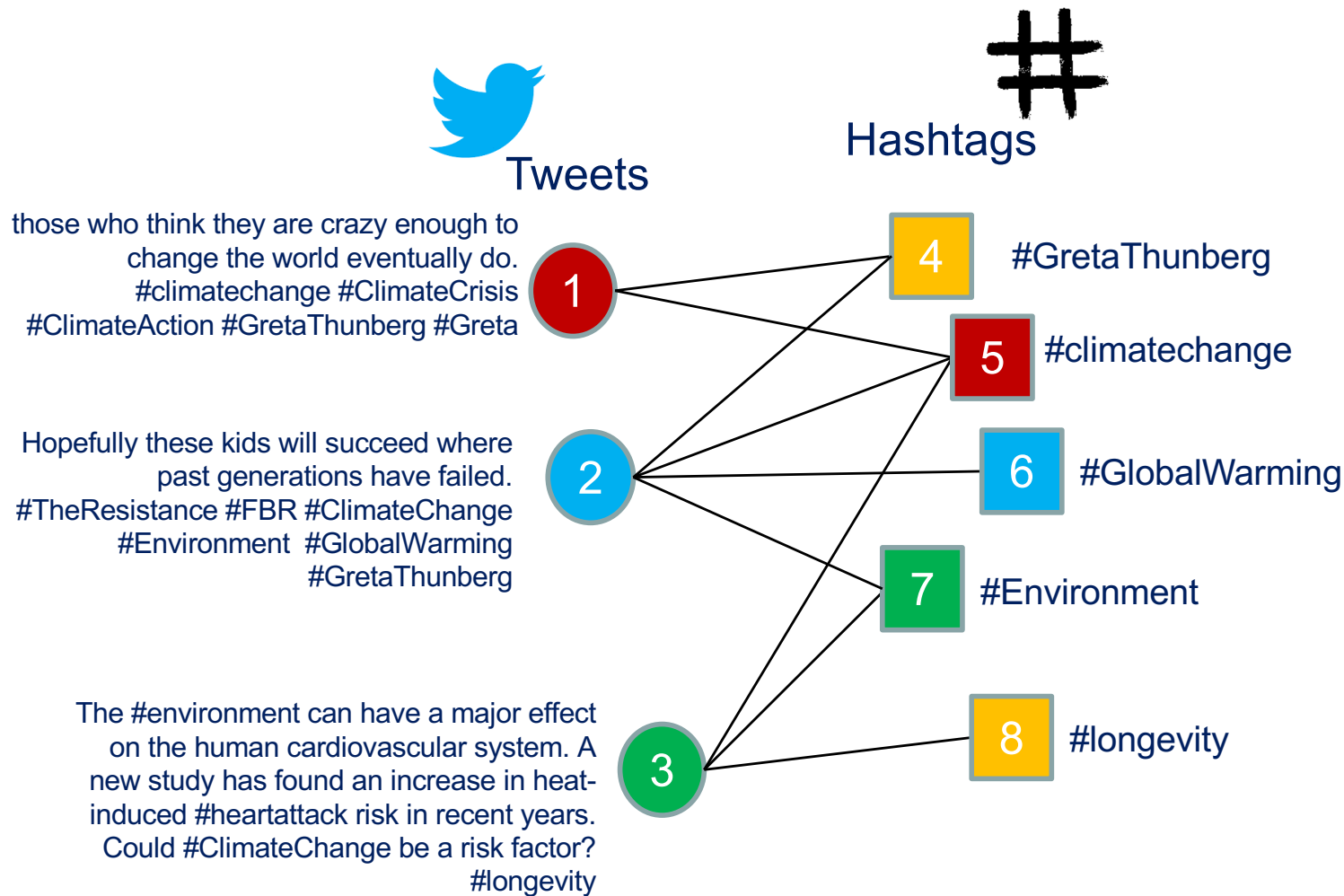
ICT for Internet & multimedia, Data science, Physics of data

Semantic networks

recap



Conceptual picture of a semantic network on Twitter





Probability matrices linking words to documents

number of occurrences
of words in documents

$N_{wd} =$

0	1	1	1	#globalwarming
1	1	1	1	#climatechange
1	0	1	0	#climateaction
0	0	1	1	#gretathunberg
1	0	1	1	#environment

probability of words
given a documents

$P_{w|d} =$

0	$\frac{1}{2}$	$\frac{1}{5}$	$\frac{1}{4}$
$\frac{1}{3}$	$\frac{1}{2}$	$\frac{1}{5}$	$\frac{1}{4}$
$\frac{1}{3}$	0	$\frac{1}{5}$	0
0	0	$\frac{1}{5}$	$\frac{1}{4}$
$\frac{1}{3}$	0	$\frac{1}{5}$	$\frac{1}{4}$

we identify a
document
probability

$p_d = \begin{cases} \frac{1}{D} & \text{equally likely} \\ \frac{n_d}{\sum_d n_d} & \text{custom} \end{cases}$

we capture the **statistical**
properties by
normalizing by columns



Probability matrices projecting to words or documents

bipartite
network

joint probability of words
and documents

$$\mathbf{P}_{wd} = \mathbf{P}_{w|d} \text{diag}(\mathbf{p}_d)$$

0	1/8	1/20	1/16
1/12	1/8	1/20	1/16
1/12	0	1/20	0
0	0	1/20	1/16
1/12	0	1/20	1/16



marginal probabilities

$$\mathbf{p}_w = \mathbf{P}_{wd} \mathbf{1} \quad \mathbf{p}_d = \mathbf{P}_{wd}^T \mathbf{1}$$

$$p_{w_1, w_2} = \sum_d p_{w_1|d, \cancel{w_2}} p_{w_2, d}$$

$$\mathbf{P}_{ww} = \mathbf{P}_{wd} \text{diag}(\mathbf{p}_d)^{-1} \mathbf{P}_{wd}^T$$

$$\mathbf{p}_w = \mathbf{P}_{ww} \mathbf{1}$$

projection on words

projection on documents

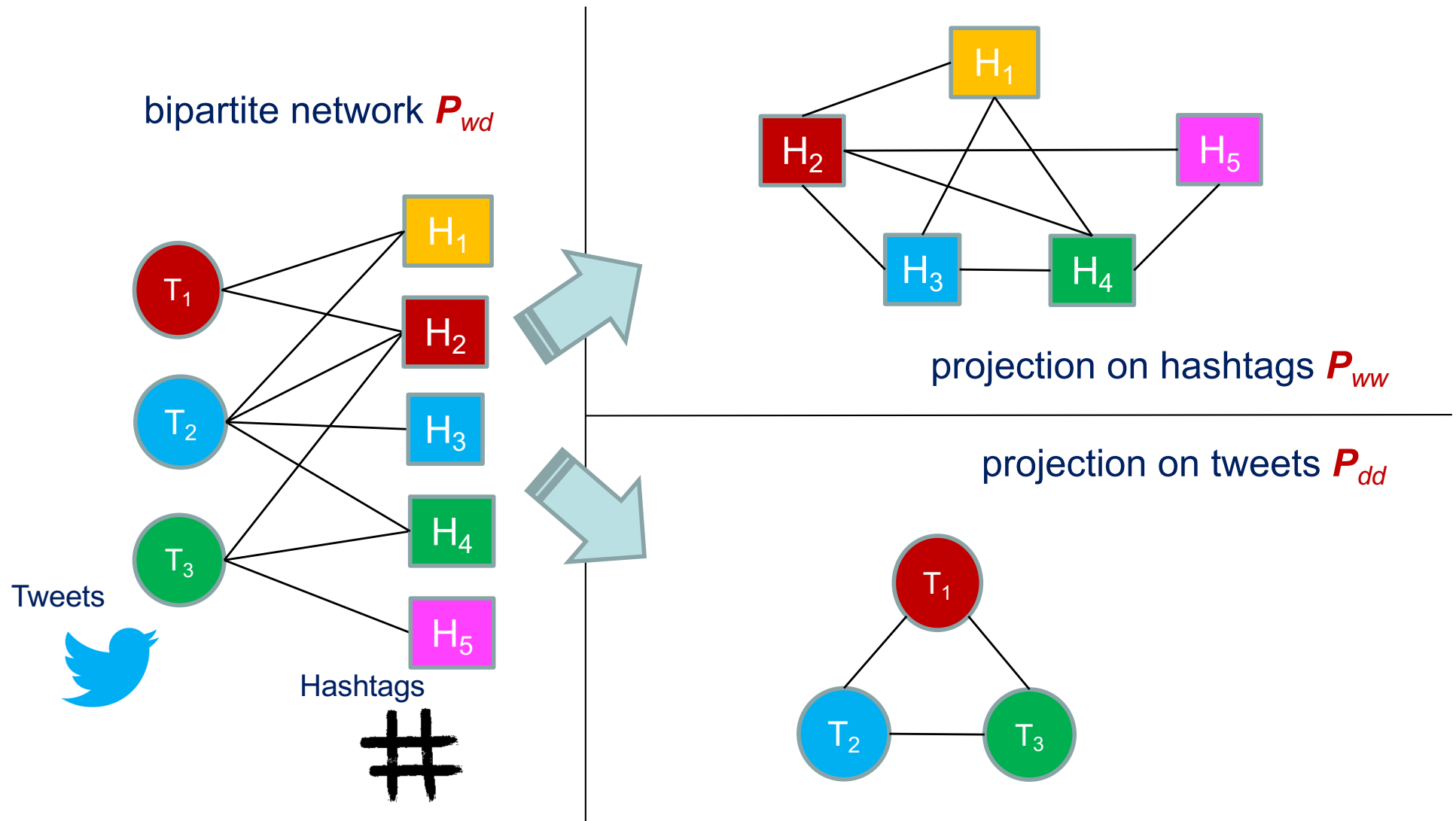
$$p_{d_1, d_2} = \sum_w p_{d_1|w, \cancel{d_2}} p_{d_2, w}$$

$$\mathbf{P}_{dd} = \mathbf{P}_{wd}^T \text{diag}(\mathbf{p}_w)^{-1} \mathbf{P}_{wd}$$

$$\mathbf{p}_d = \mathbf{P}_{dd} \mathbf{1}$$

Bipartite and projected networks

a comparison





The role of TF-IDF

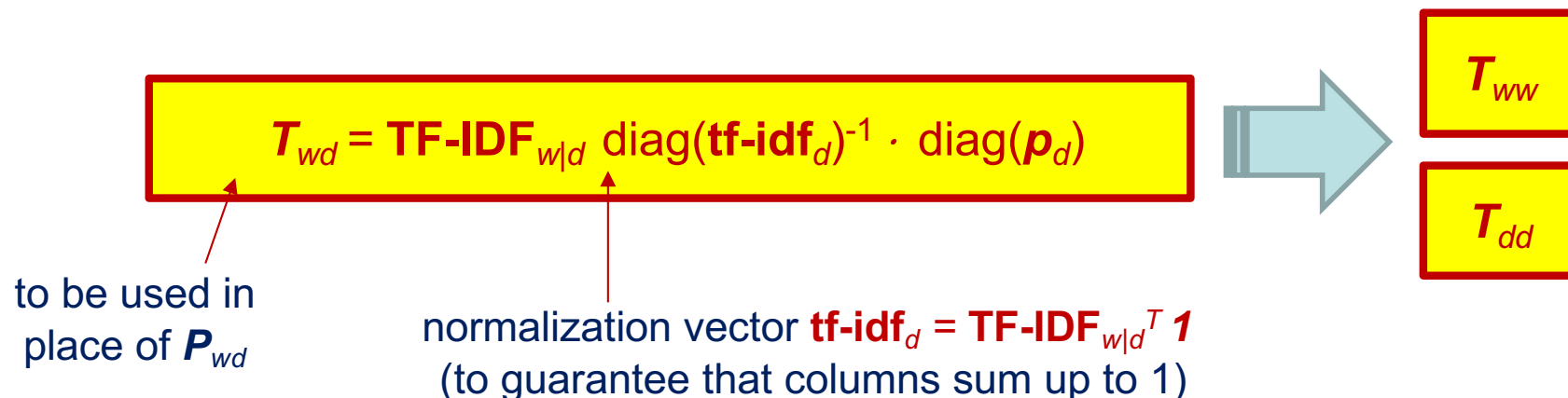
term frequency – inverse document frequency

term frequency =
frequency (probability) of
the word in the document

inverse document frequency =
(log) fraction of documents
that contain the word

$$\text{TF-IDF}_{w|d} = p_{w|d} \cdot -\log \left(\frac{\sum_d (n_{wd} > 0)}{D} \right)$$

- ❑ An heuristic
- ❑ **Punishes** words that appear in many documents
- ❑ **Enhances** words that are document specific



Topic detection

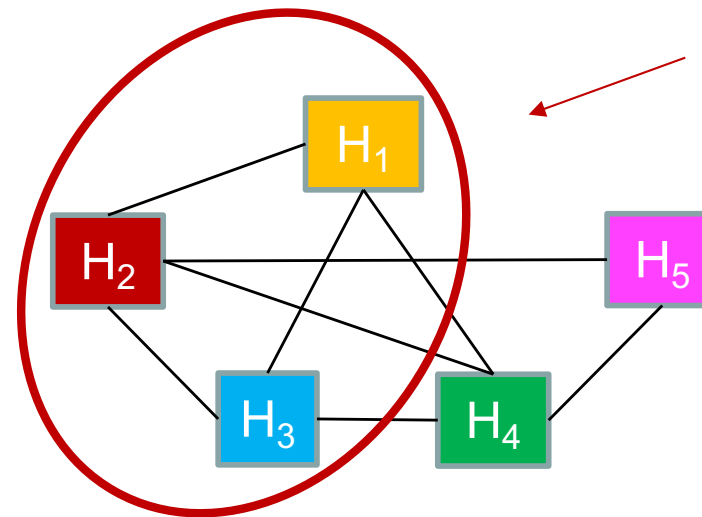
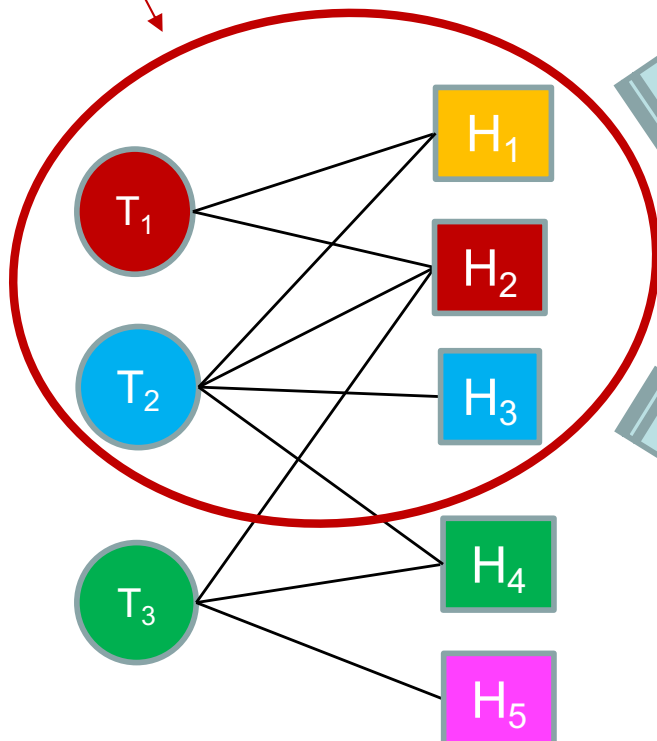
i.e., community detection in semantic networks

Topic detection

in bipartite and projection networks

bipartite network P_{wd} or T_{wd}

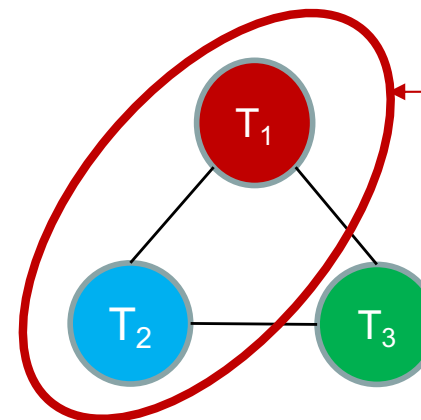
a community identifies both
documents and words



a community
identifies words

projection on words P_{ww} or T_{ww}

projection on documents P_{dd} or T_{dd}

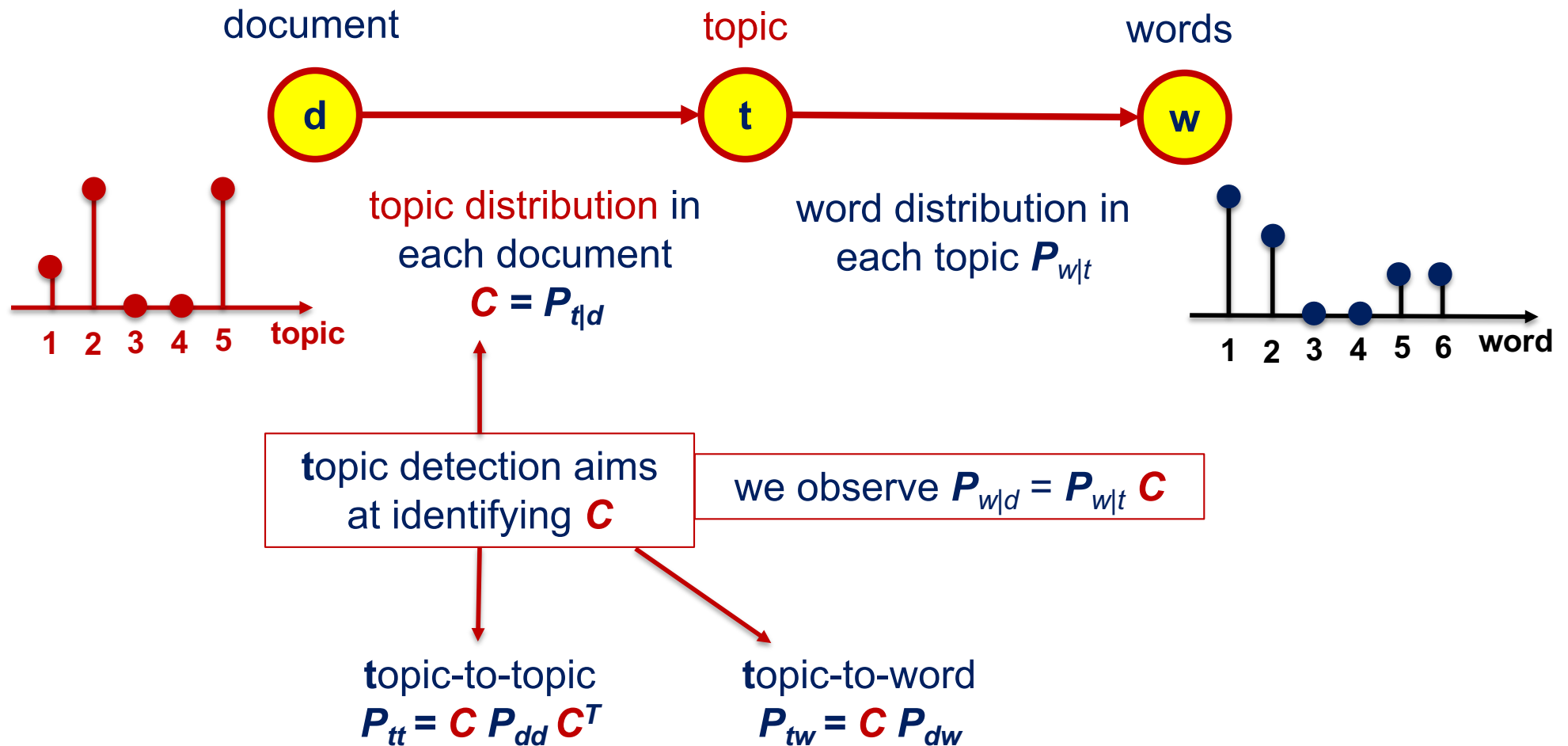


a community
identifies documents



The reference model

Under the presence of topics





Modularity and normalized cut

a wrap-up in topic detection

C topic assignment to be assessed for quality

$$P_{tt} = C P_{dd} C^T$$

can be interpreted as a probability matrix linking topics, its entries are the sum of the links of **A** from topic *i* to topic *j*

P_{11}	P_{12}	P_{13}
P_{21}	P_{22}	P_{23}
P_{31}	P_{32}	P_{33}

$$p_t = P_{tt} \mathbf{1}$$

can be interpreted as the probability vector of topics

modularity

$$Q = \sum_t (P_{tt} - p_t^2) < 1$$

to be maximized

normalized cut

normalized version

$$Ncut = 1 - \frac{\sum_t P_{tt}/p_t}{\sum_t 1} > 0$$

to be minimized



PageRank vector (ranking of documents)

$$\mathbf{r} = (1-c) \mathbf{P}_{d|d} \mathbf{r} + c \mathbf{1}/N$$

Here \mathbf{c}_i
is the i th
row of \mathbf{C}

$$\mathbf{P}_{d|d} = \mathbf{P}_{dd} \text{diag}^{-1}(\mathbf{p}_d)$$

$$\mathbf{q}_i = \left(1 - (1-c) \frac{\mathbf{c}_i \mathbf{1}}{N}\right) \mathbf{z}_i \mathbf{1} - c \mathbf{c}_i \mathbf{P}_{d|d} \mathbf{z}_i^T$$

$$\mathbf{z}_i = \mathbf{c}_i \text{diag}(\mathbf{r})$$

$$\text{InfoMap} = f(\mathbf{q}) + \sum_i f([q_i, \mathbf{z}_i])$$

normalized version

$$\frac{\text{InfoMap}}{f(\mathbf{r})} - 1$$

to be minimized

entropy function

$$f(\mathbf{x}) = - \sum_i x_i \log \left(\frac{x_i}{\sum_j x_j} \right)$$



Normalized mutual information

a wrap-up in topic detection

statistical dependencies about
words and topics

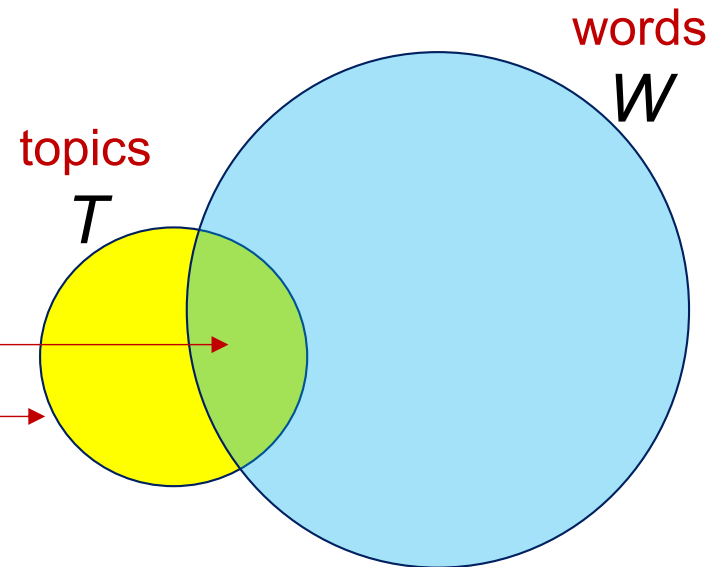
$$\mathbf{P}_{wt} = \mathbf{P}_{wd} \mathbf{C}^T$$

probability of a
topic

$$\mathbf{p}_t = \mathbf{P}_{wt}^T \mathbf{1}$$

fraction of knowledge related to
the topic that is explained by
words (equal to 1 if topics use
different words)

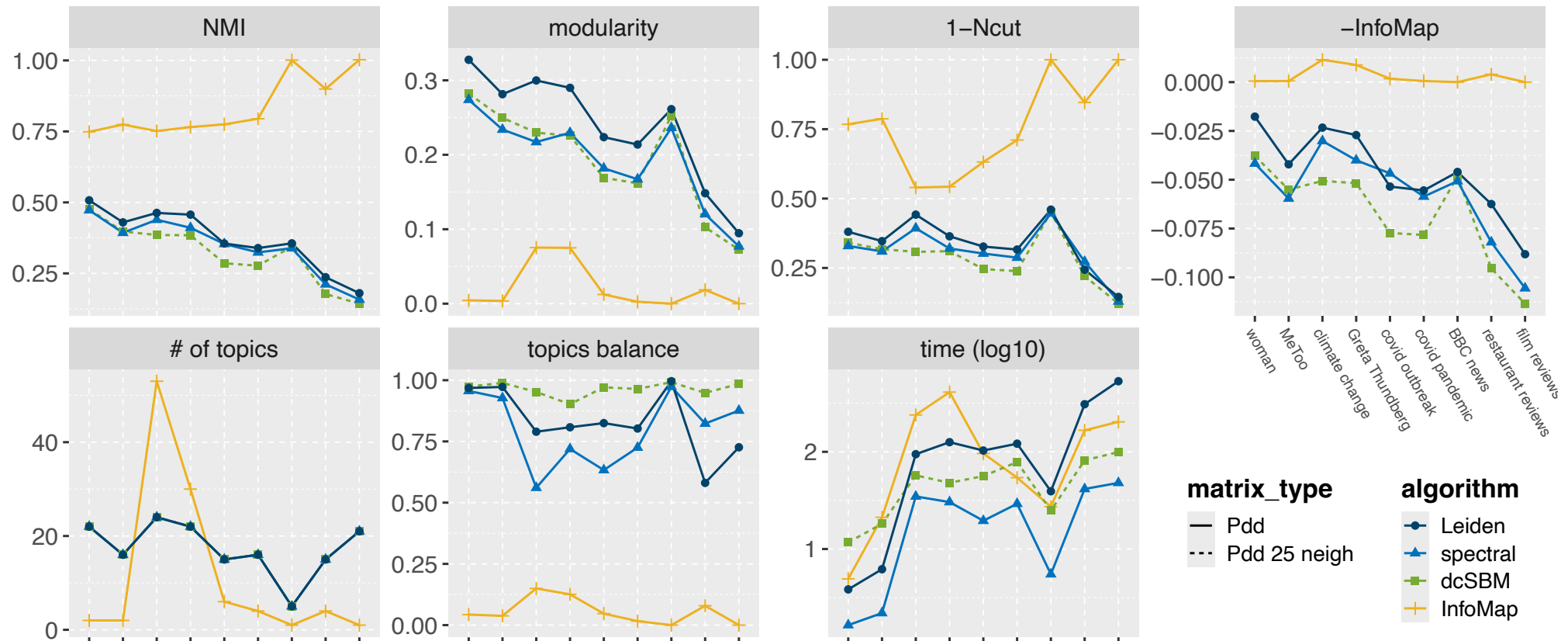
$$\text{NMI} = \frac{I(W;T)}{H(T)}$$





Louvain, InfoMap and Spectral clust.

On a semantic network

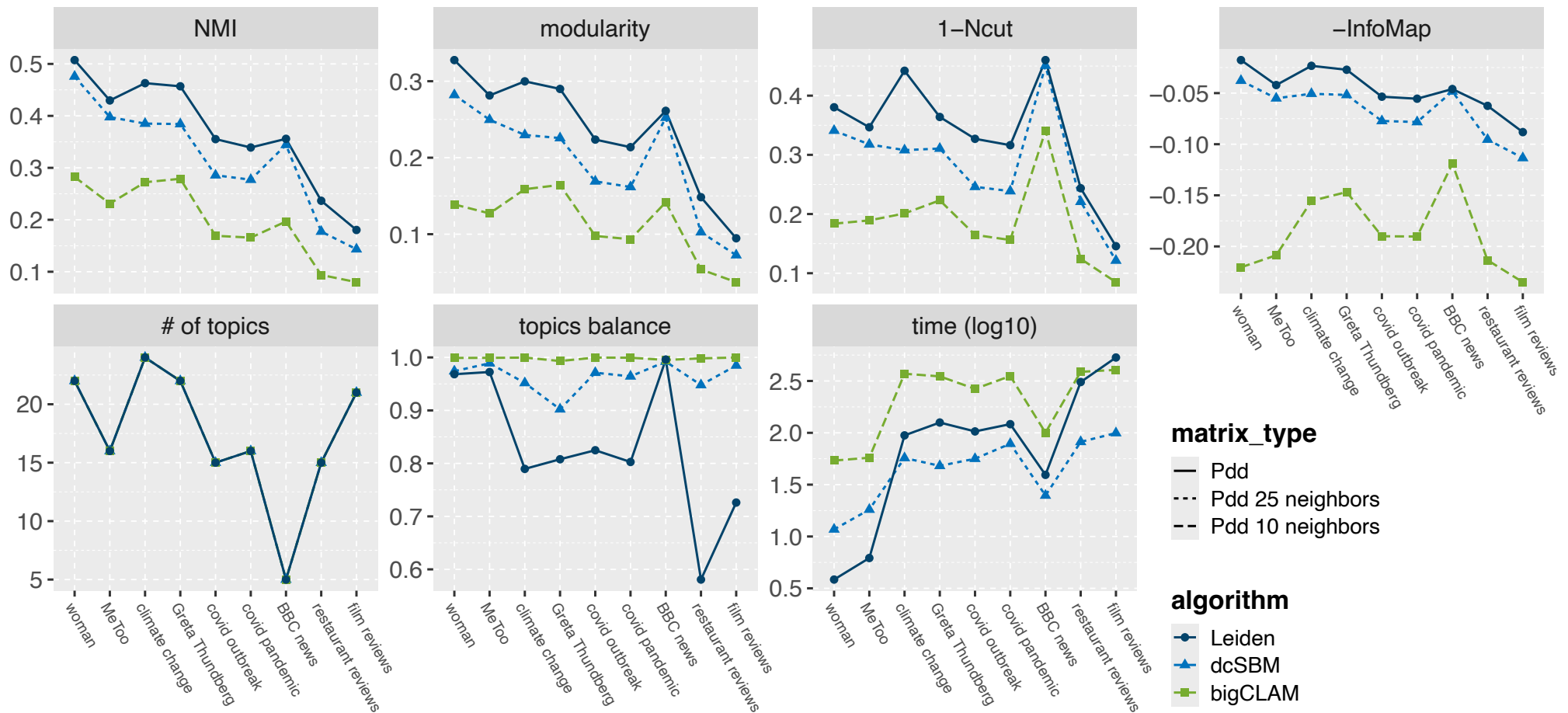




BigClam and SBMs at work

On a semantic network

SBM approaches





- ❑ Non-negative matrix factorization (**NMF**)
- ❑ Latent Dirichlet allocation (**LDA**)
- ❑ Variational auto-encoders (**VAE**)
- ❑ Embeddings and **BERTopic**

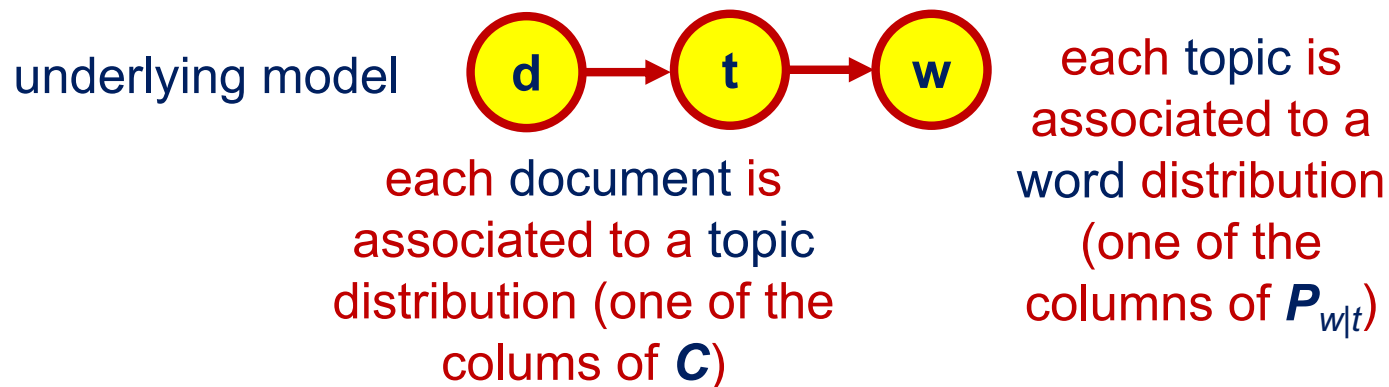
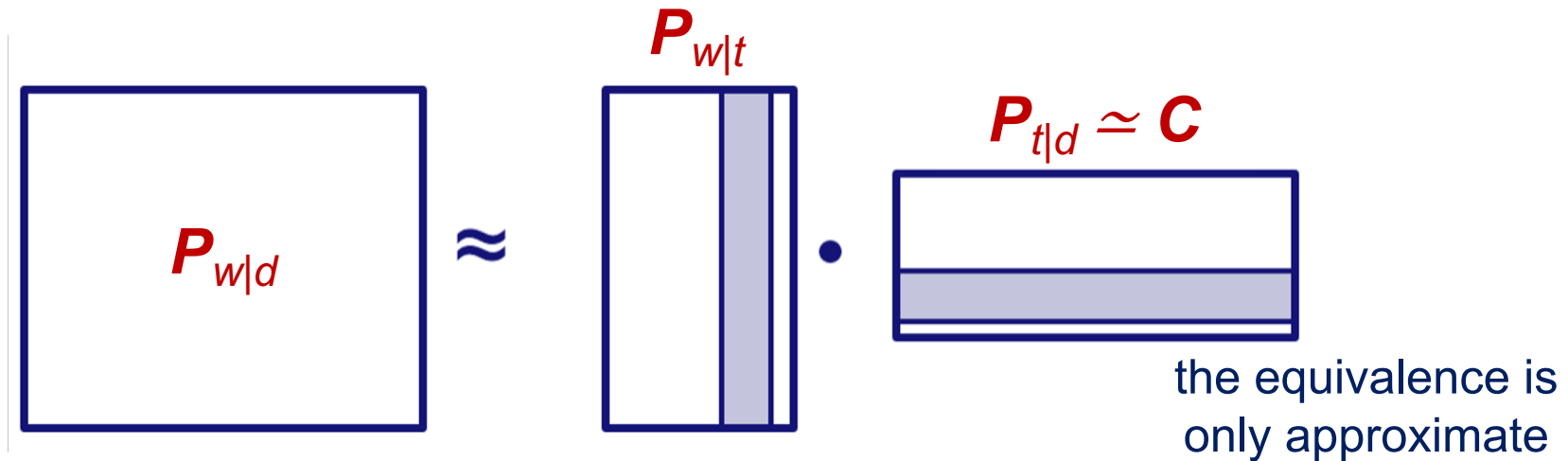
Non-negative Matrix Factorization

and its application to topic detection



NMF = nonnegative matrix factorization

rationale





$\mathbf{A} = \mathbf{P}_{w|d}$ is column stochastic

$$\operatorname{argmin}_{\mathbf{W} \geq 0, \mathbf{H} \geq 0} \sum_{ij} |A_{ij} - [\mathbf{WH}]_{ij}|^2$$

minimizing the Frobenius norm does not ensure a column stochastic product \mathbf{WH}

$$\operatorname{argmin}_{\mathbf{W} \geq 0, \mathbf{H} \geq 0} \sum_{ij} A_{ij} \log \left(\frac{A_{ij}}{[\mathbf{WH}]_{ij}} \right) - A_{ij} + [\mathbf{WH}]_{ij}$$

$$f(y) = x \log \left(\frac{x}{y} \right) - x + y$$

$$f'(y) = -\frac{x}{y} + 1 = 0 \rightarrow y = x$$

minimizing the generalized Kullback-Leibler divergence ensures a column stochastic product \mathbf{WH}

Ho & Van Dooren. "Non-negative matrix factorization with fixed row and column sums." (2008)



```
from sklearn.decomposition import NMF  
Pwgd = Pwd/Pwd.sum(axis=0).flatten()
```

run on different number of topics, then choose
the best fit, e.g., according to modularity

```
# fit nmf model  $X = W \cdot H$   
model = NMF(n_components=i, init='nndsvd',  
            solver='mu', beta_loss='kullback-leibler')  
W = model.fit_transform(Pwgd)  
H = sps.csr_matrix(model.components_)  
# column normalized versions  
H = sps.diags(W.sum(axis=0).flatten())*H # Ptgd  
W = W/W.sum(axis=0).flatten() # Pwgt  
# community assignment C  
C = sps.csr_matrix(np.transpose(H/H.sum(axis=0).flatten()))
```

wisely initialize
for best
performance

choose generalized
Kullback-Leibler
divergence, and the
related solver

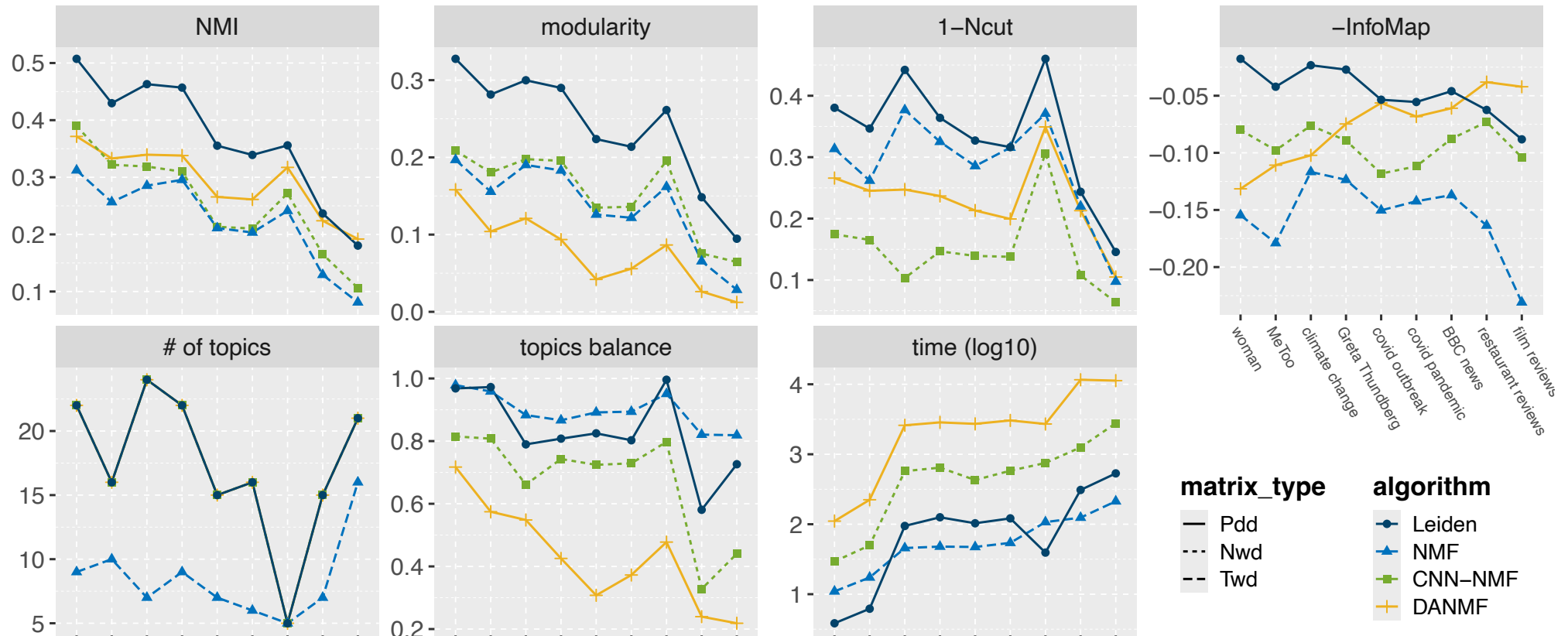
need to make W
column stochastic,
to have H column
stochastic too

force column stochasticity in H (not needed though)



NMF at work

On a semantic network

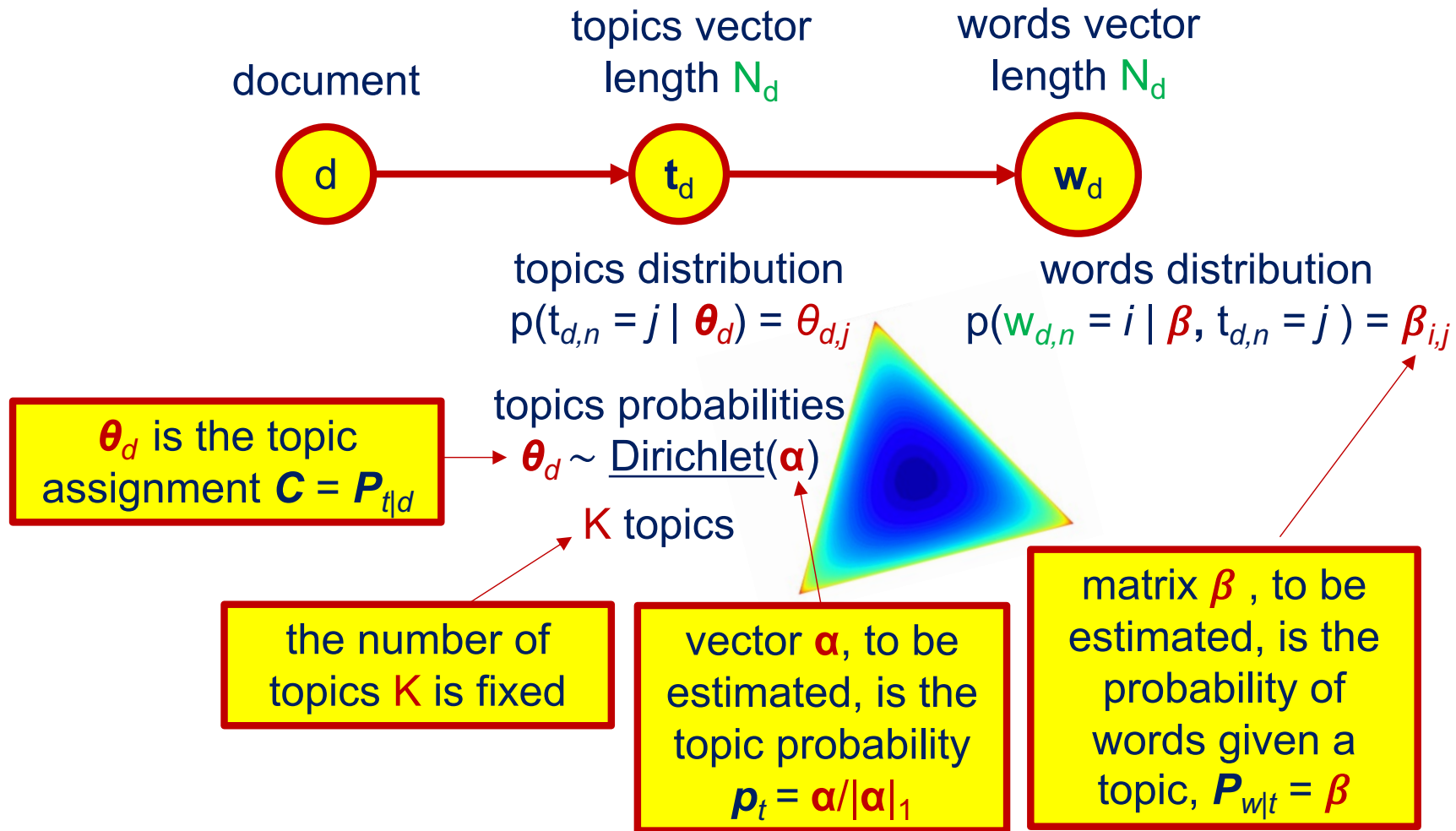




- ❑ Naturally provides a soft topic assignment
- ❑ NMF – not strikingly good
 - probably due to the fact that we want to express a sparse matrix through an eigenvector-like product with few eigenvectors (the fit is far from ideal)
- ❑ Comparison – with Louvain
 - much weaker
- ❑ Complexity – generally slow
 - need to test it for different numbers of topics ☹
 - fast for fixed topic number

Latent Dirichlet allocation

LDA = a stochastic model for topic detection





topics assignment probability (Dirichlet)

$$p(\boldsymbol{\theta}_d | \boldsymbol{\alpha}) = \frac{\Gamma(\sum_{k=1}^K \alpha_k)}{\sum_{k=1}^K \Gamma(\alpha_k)} \prod_{k=1}^K [\theta_{d,k}]^{\alpha_k - 1}$$

words probability

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

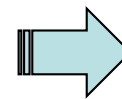
this dependence
between $\boldsymbol{\beta}$ and $\boldsymbol{\theta}$
is the trickiest part

overall probability

$$p(\text{corpus} | \boldsymbol{\alpha}, \boldsymbol{\beta}) = \prod_d \int p(\mathbf{w}_d | \boldsymbol{\beta}, \boldsymbol{\theta}_d) p(\boldsymbol{\theta}_d | \boldsymbol{\alpha}) d\boldsymbol{\theta}_d$$

target optimization

$$\operatorname{argmax}_{\boldsymbol{\alpha}, \boldsymbol{\beta}} p(\text{corpus} | \boldsymbol{\alpha}, \boldsymbol{\beta})$$



$$\mathbf{C} = \mathbf{P}_{t|d} = \boldsymbol{\theta}$$
$$\mathbf{P}_{wt} = \boldsymbol{\beta} \operatorname{diag}(\boldsymbol{\alpha} / |\boldsymbol{\alpha}|_1)$$

this is what we get



```
from sklearn.decomposition import LatentDirichletAllocation
```

```
# fit lda model
```

```
lda = LatentDirichletAllocation(n_components=i,  
                                learning_method="batch")
```

```
lda.fit(Mwd.T)
```

initialise and fit model

```
# community assignment C = Ptgd'
```

```
C = sps.csr_matrix(lda.transform(Mwd.T))
```

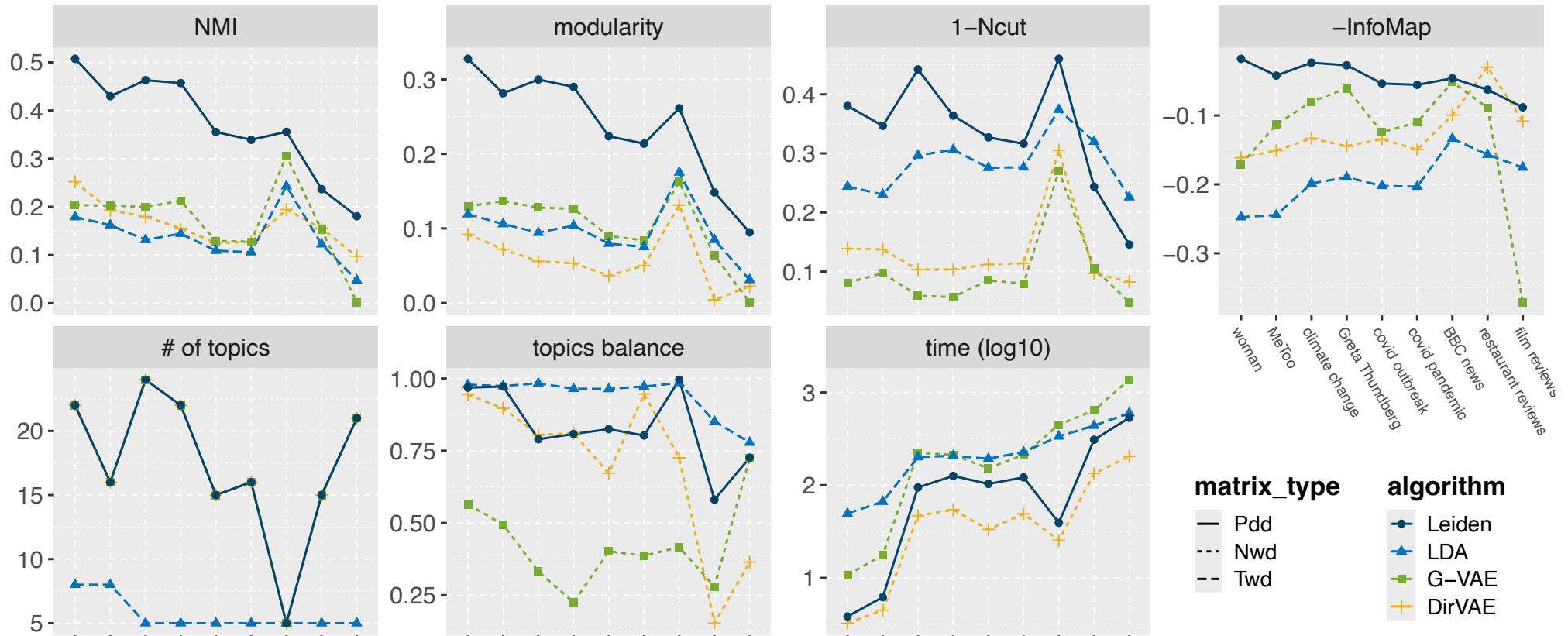
extract topic assignment

LDA



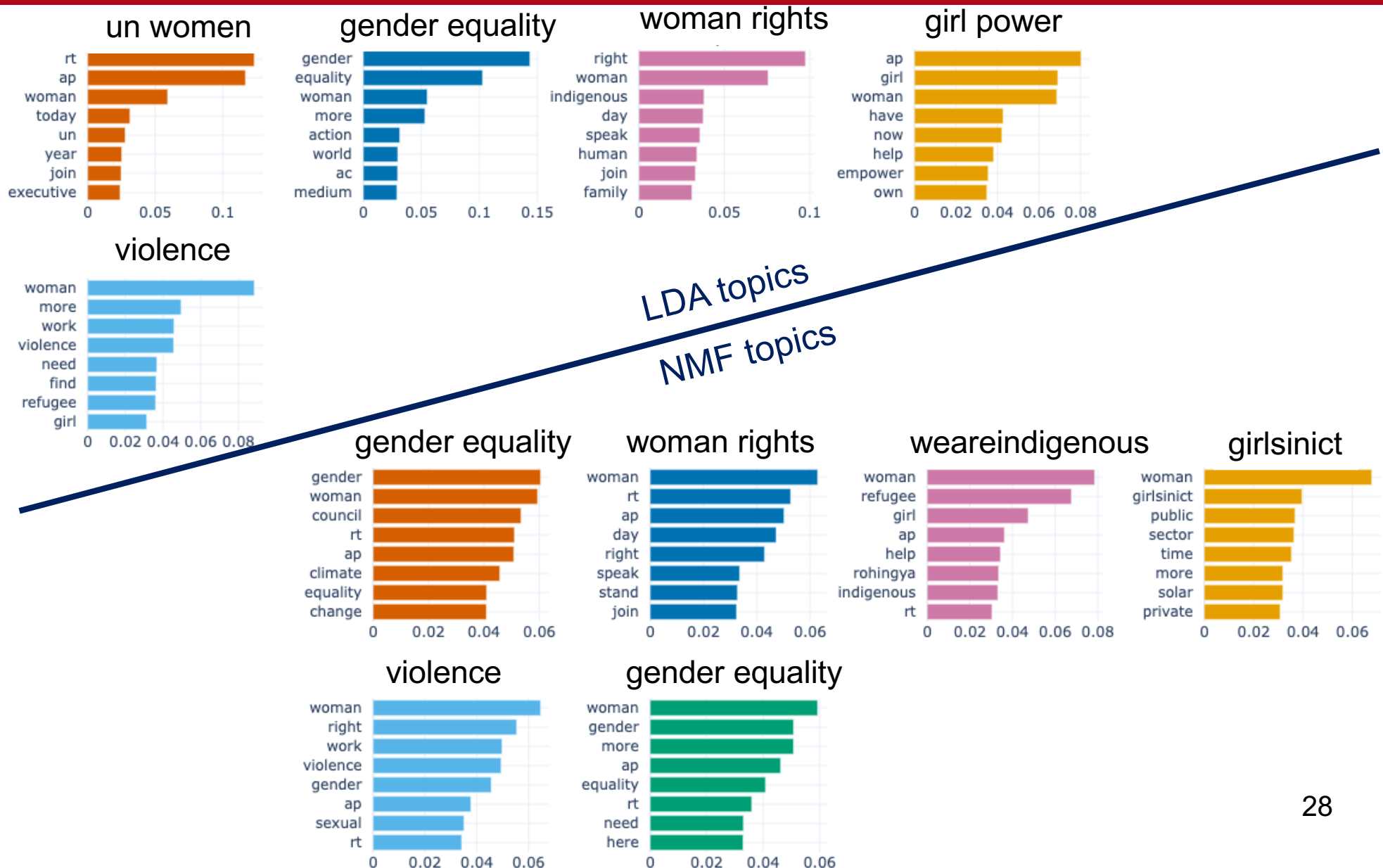
LDA at work

On a semantic network





A comparison NMF versus LDA topics





- ❑ Naturally provides a soft topic assignment
- ❑ LDA – not strikingly good
 - same eigenvector-like product as NMF
 - worse than NMF ... known issue ☹
 - probably due to the **Dirichlet** assumption (questionable)
 - and the **variational inference** (suboptimum approach)
- ❑ Comparison – with Louvain
 - much weaker
- ❑ Complexity – generally slow
 - need to test it for different numbers of topics ☹
 - fast for fixed topic number

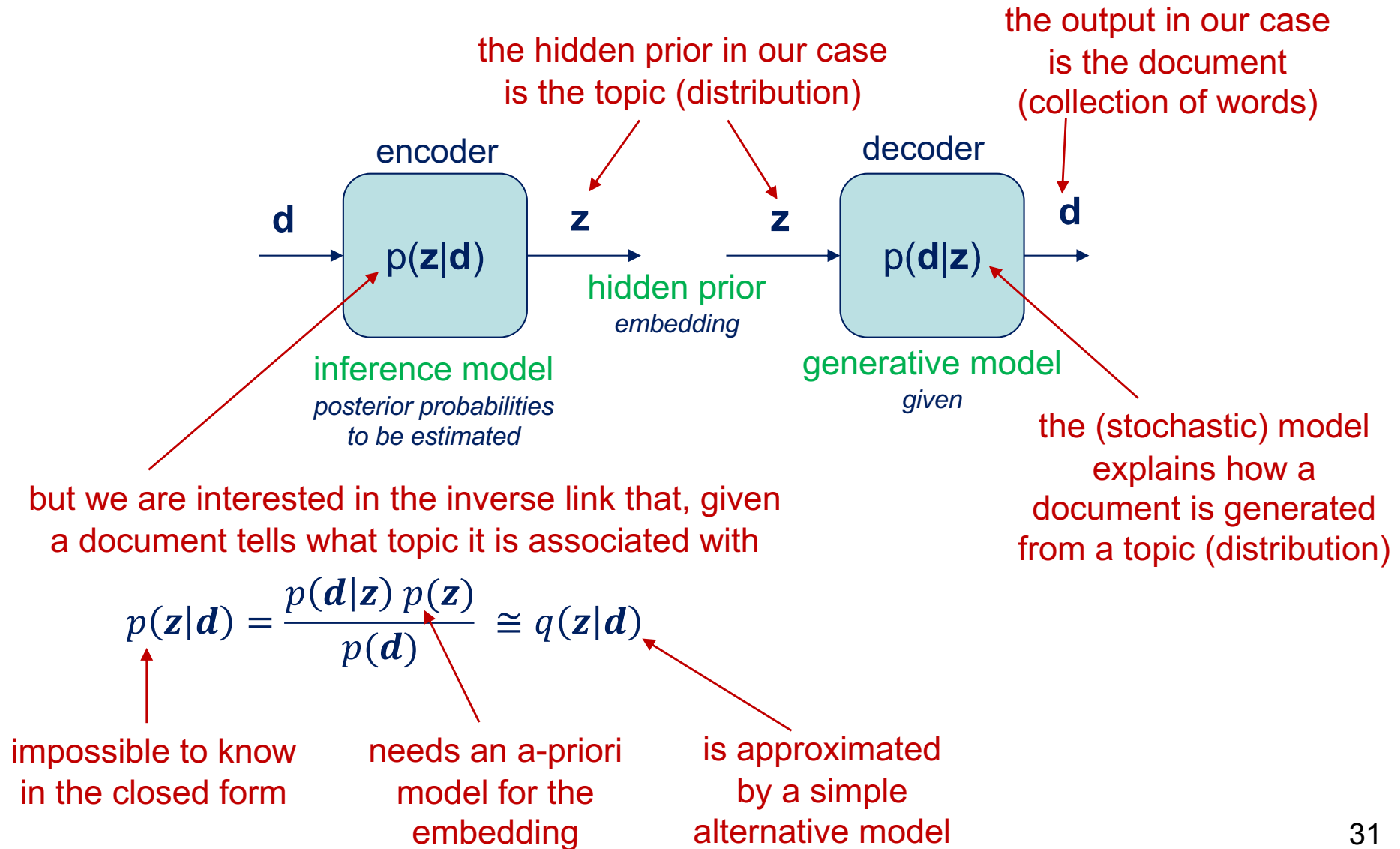
Variational Auto Encoders

an application to topic analysis



Variational Auto-Encoders

Kingma, Welling, "Auto-encoding variational Bayes," (2013)



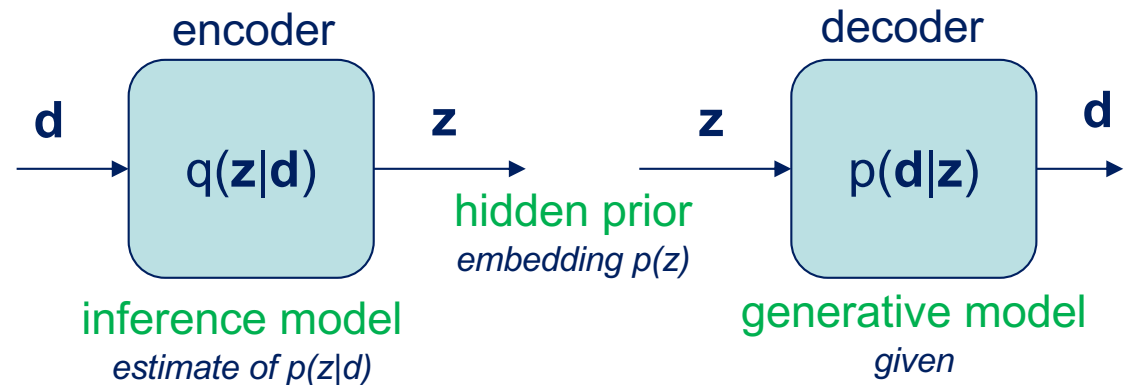


VAE optimization rationale

ELBO = evidence lower bound

ELBO

$$\mathcal{L}_{\theta,\phi}(\mathbf{d}) \leq \log p_{\theta}(\mathbf{d})$$



$$\begin{aligned}\mathcal{L}_{\theta,\phi}(\mathbf{d}) &= \log p_{\theta}(\mathbf{d}) - D_{\text{KL}}\left(q_{\phi}(z|\mathbf{d}) \parallel p_{\theta}(z|\mathbf{d})\right) \\ &= \int dz q_{\phi}(z|\mathbf{d}) \log \left(\frac{p_{\theta}(z, \mathbf{d})}{q_{\phi}(z|\mathbf{d})} \right) \\ &= \underbrace{\int dz q_{\phi}(z|\mathbf{d}) \log \left(p_{\theta}(\mathbf{d}|z) \right)}_{\mathcal{L}_1} + \underbrace{\int dz q_{\phi}(z|\mathbf{d}) \log \left(\frac{p_{\theta}(z)}{q_{\phi}(z|\mathbf{d})} \right)}_{\mathcal{L}_2}\end{aligned}$$

to be maximized wrt parameters θ and ϕ provides fitting on $p(\mathbf{z})$, $p(\mathbf{d}|\mathbf{z})$, and $q(\mathbf{z}|\mathbf{d})$

a-priori model (given)

inference model (approximate) \mathcal{L}_1 generative model (given) \mathcal{L}_2



L2 ELBO function

usually has a compact expression

$$\underbrace{\int dz \, q_{\phi}(z|\mathbf{d}) \log \left(\frac{p_{\theta}(z)}{q_{\phi}(z|\mathbf{d})} \right)}_{\mathcal{L}_2}$$

both should have a simple parametrization on θ and ϕ

e.g., the Gaussian case

$$p_{\theta}(\mathbf{z}) = \frac{1}{\sqrt{\det(2\pi \operatorname{diag}(\boldsymbol{\sigma}_{\theta}^2))}} \exp \left(-\frac{1}{2}(\mathbf{z} - \boldsymbol{\mu}_{\theta})^T \operatorname{diag}^{-1}(\boldsymbol{\sigma}_{\theta}^2)(\mathbf{z} - \boldsymbol{\mu}_{\theta}) \right)$$

$$q_{\phi}(\mathbf{z}|\mathbf{d}) = \frac{1}{\sqrt{\det(2\pi \operatorname{diag}(\boldsymbol{\sigma}_{\phi}^2(\mathbf{d})))}} \exp \left(-\frac{1}{2}(\mathbf{z} - \boldsymbol{\mu}_{\phi}(\mathbf{d}))^T \operatorname{diag}^{-1}(\boldsymbol{\sigma}_{\phi}^2(\mathbf{d}))(\mathbf{z} - \boldsymbol{\mu}_{\phi}(\mathbf{d})) \right)$$

$$\mathcal{L}_2(\theta, \phi) = \frac{1}{2} \sum_i \left(1 + \log \left(\frac{\sigma_{\phi,i}^2(\mathbf{d})}{\sigma_{\theta,i}^2} \right) - \frac{\sigma_{\phi,i}^2(\mathbf{d})}{\sigma_{\theta,i}^2} - \frac{(\mu_{\phi,i}(\mathbf{d}) - \mu_{\theta,i})^2}{\sigma_{\theta,i}^2} \right)$$



L1 ELBO function

approximated through Monte Carlo estimation

$$\underbrace{\int dz q_{\phi}(z|\mathbf{d}) \log(p_{\theta}(\mathbf{d}|z))}_{\mathcal{L}_1}$$

mostly too complex to be
written in the closed form

solution: Monte Carlo approximation

$$\mathcal{L}_1(\theta, \phi) = \frac{1}{L} \sum_{\ell=1}^L \log(p_{\theta}(\mathbf{d}|\mathbf{z}_{\ell}))$$

samples generated according
to the correct distribution

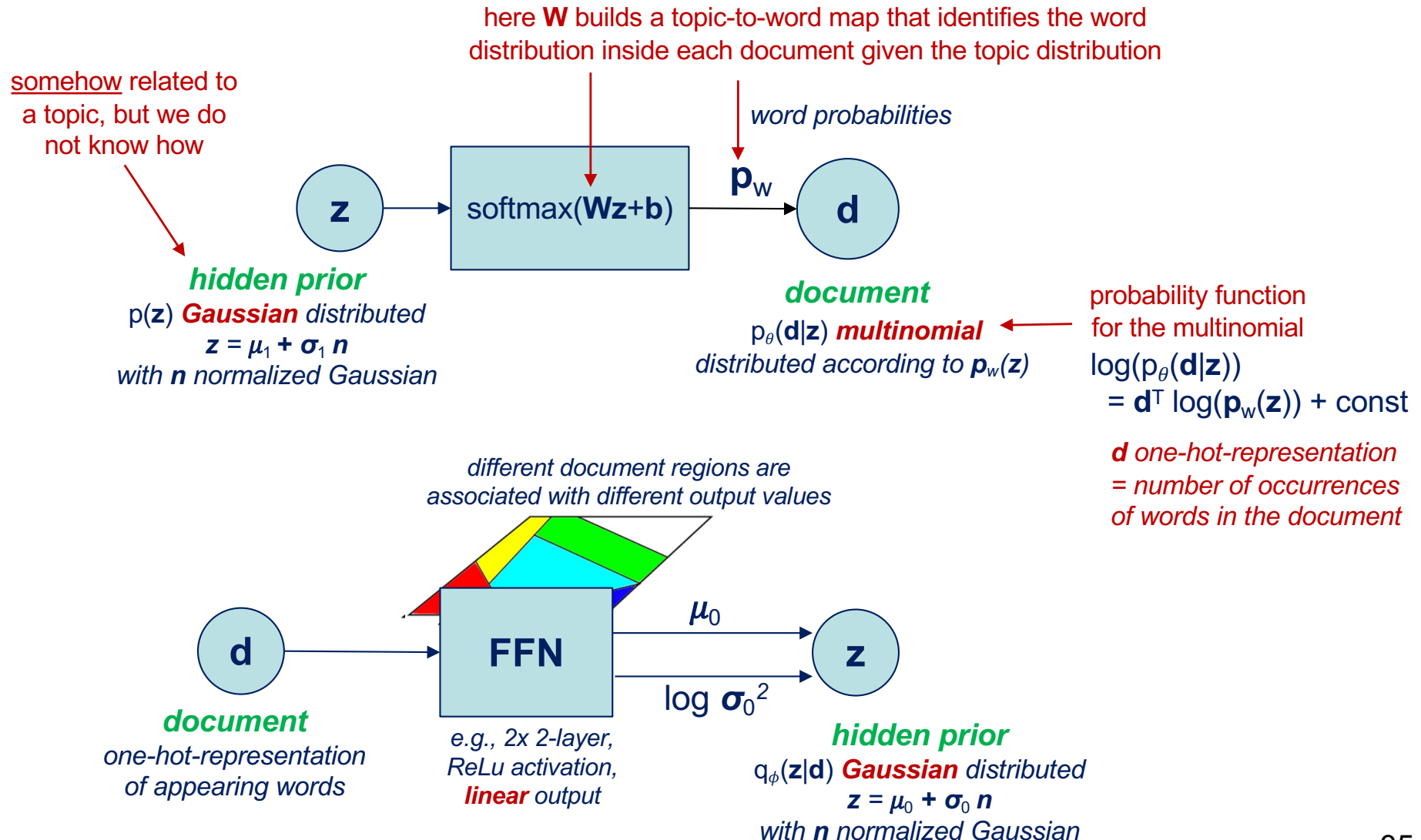
$$\mathbf{z}_{\ell} \sim q_{\phi}(\mathbf{z}|\mathbf{d})$$

e.g., the Gaussian case

$$\mathbf{z}_{\ell} = \boldsymbol{\mu}_{\phi}(\mathbf{d}) + \boldsymbol{\sigma}_{\phi}(\mathbf{d}) \mathbf{n}_{\ell}$$

need to generate these
once, then use them
throughout the process

$$\mathbf{n}_{\ell} \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$$





one-hot-representation of a document = number of occurrences of words in the document

$$\mathcal{L}(\theta, \phi) = \frac{1}{L} \sum_{\ell=1}^L \sum_m \mathbf{d}_m^T \log \left(\underbrace{\text{softmax}(\mathbf{b} + \mathbf{W}(\boldsymbol{\mu}_0(\mathbf{d}_m) + \boldsymbol{\sigma}_0(\mathbf{d}_m) \mathbf{n}_{m,\ell}))}_{\text{decoder map}} \right)$$

decoder model encoder model

normalized Gaussian samples

$$+ \frac{1}{2} \sum_m \sum_i 1 + \log \left(\frac{\sigma_{0,i}^2(\mathbf{d}_m)}{\sigma_{1,i}^2} \right) - \frac{\sigma_{0,i}^2(\mathbf{d}_m)}{\sigma_{1,i}^2} - \frac{(\mu_{0,i}(\mathbf{d}_m) - \mu_{1,i})^2}{\sigma_{1,i}^2}$$

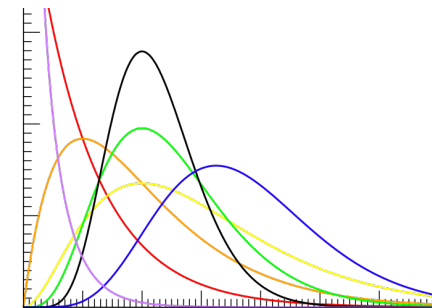
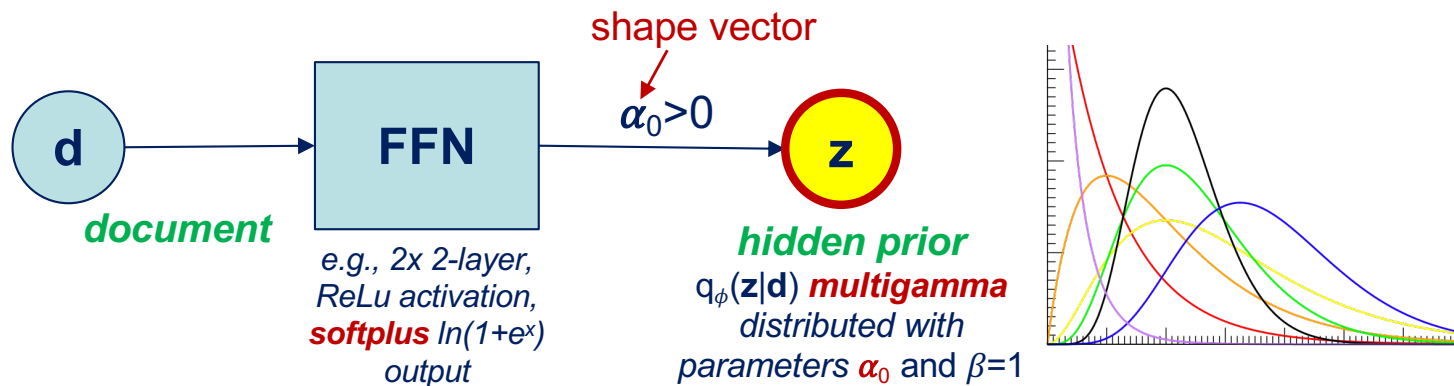
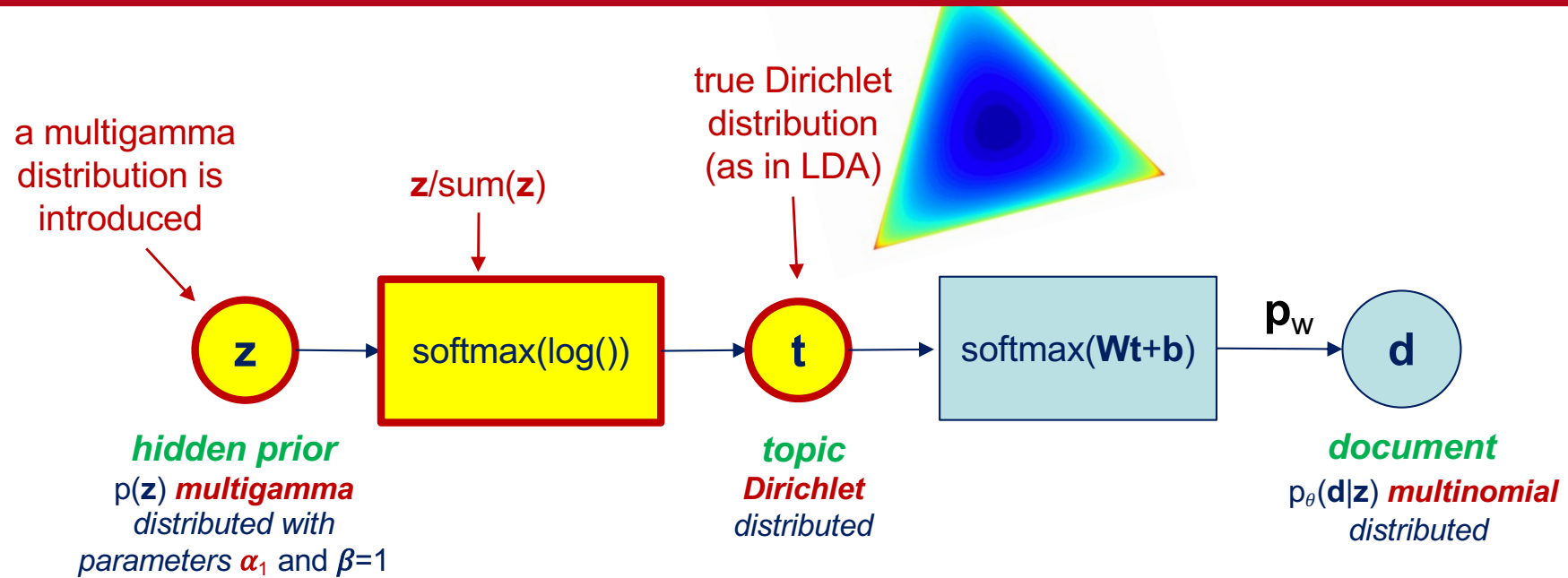
a-priori model

Not very clear where the topic is, though!



Take 2 - DirVAE

Joo, Li, Park, Moon, "Dirichlet variational autoencoder," (2020)
<https://www.sciencedirect.com/science/article/pii/S003132032030317>





one-hot-representation of a document = number of occurrences of words in the document

$f(u, \alpha) = (u \alpha \Gamma(\alpha_1))^{1/\alpha}$

approx. uniform to multigamma map

normalized uniform samples

decoder model

encoder model

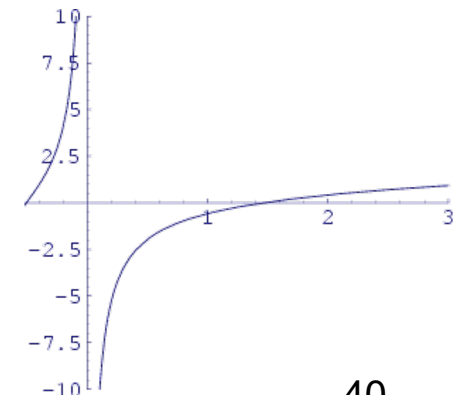
$$\mathcal{L}(\theta, \phi) = \frac{1}{L} \sum_{\ell=1}^L \sum_m \mathbf{d}_m^T \log \left(\underbrace{\text{softmax}(\mathbf{b} + \mathbf{W} \text{softmaxlog}(f(\mathbf{u}_{m,\ell}, \alpha_0(\mathbf{d}_m))))}_{\text{decoder map}} \right)$$

$$+ \sum_m \sum_i \log \left(\frac{\Gamma(\alpha_{0,i}(\mathbf{d}_m))}{\Gamma(\alpha_{1,i})} \right) - (\alpha_{0,i}(\mathbf{d}_m) - \alpha_{1,i}) \psi(\alpha_{0,i}(\mathbf{d}_m))$$

a-priori model

digamma function

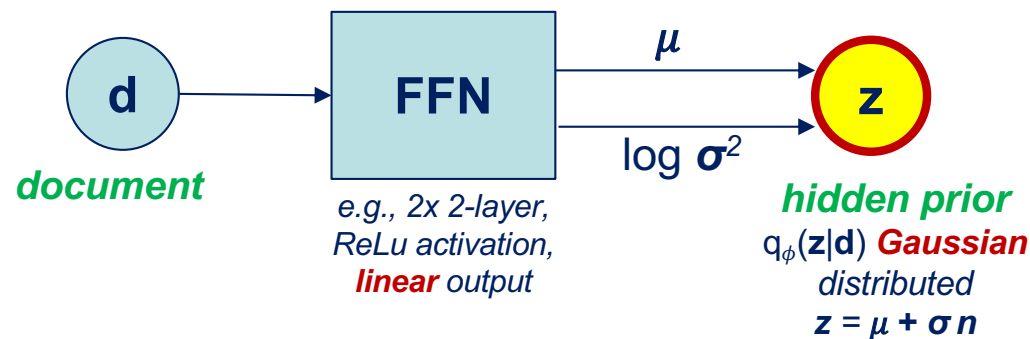
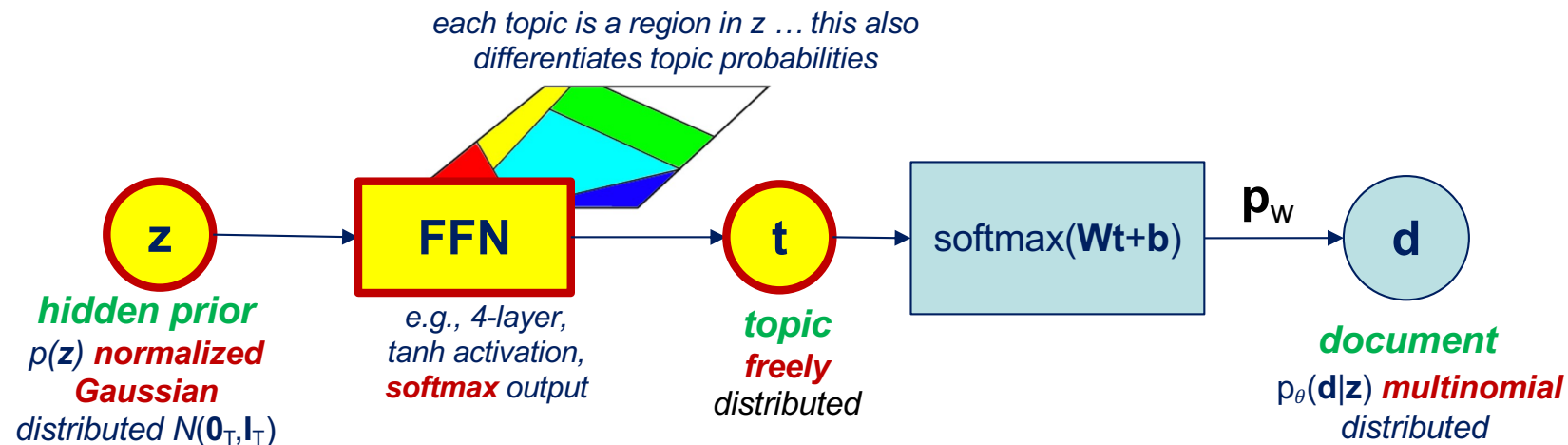
$\psi(x) = \frac{\Gamma'(x)}{\Gamma(x)}$



Now we know where
the topic is!



Take 3 – NFTM





one-hot-representation of a document = number of occurrences of words in the document

$$\mathcal{L}(\theta, \phi) = \frac{1}{L} \sum_{\ell=1}^L \sum_m \mathbf{d}_m^T \log \left(\underbrace{\text{softmax}(\mathbf{b} + \mathbf{W} \text{FFN}_1(\boldsymbol{\mu}_0(\mathbf{d}_m) + \boldsymbol{\sigma}_0(\mathbf{d}_m) \mathbf{n}_{m,\ell}))}_{\text{decoder map}} \right)$$

decoder model encoder model

normalized Gaussian samples

$$+ \frac{1}{2} \sum_m \sum_i 1 + \log(\sigma_{0,i}^2(\mathbf{d}_m)) - \sigma_{0,i}^2(\mathbf{d}_m) - (\mu_{0,i}(\mathbf{d}_m))^2$$

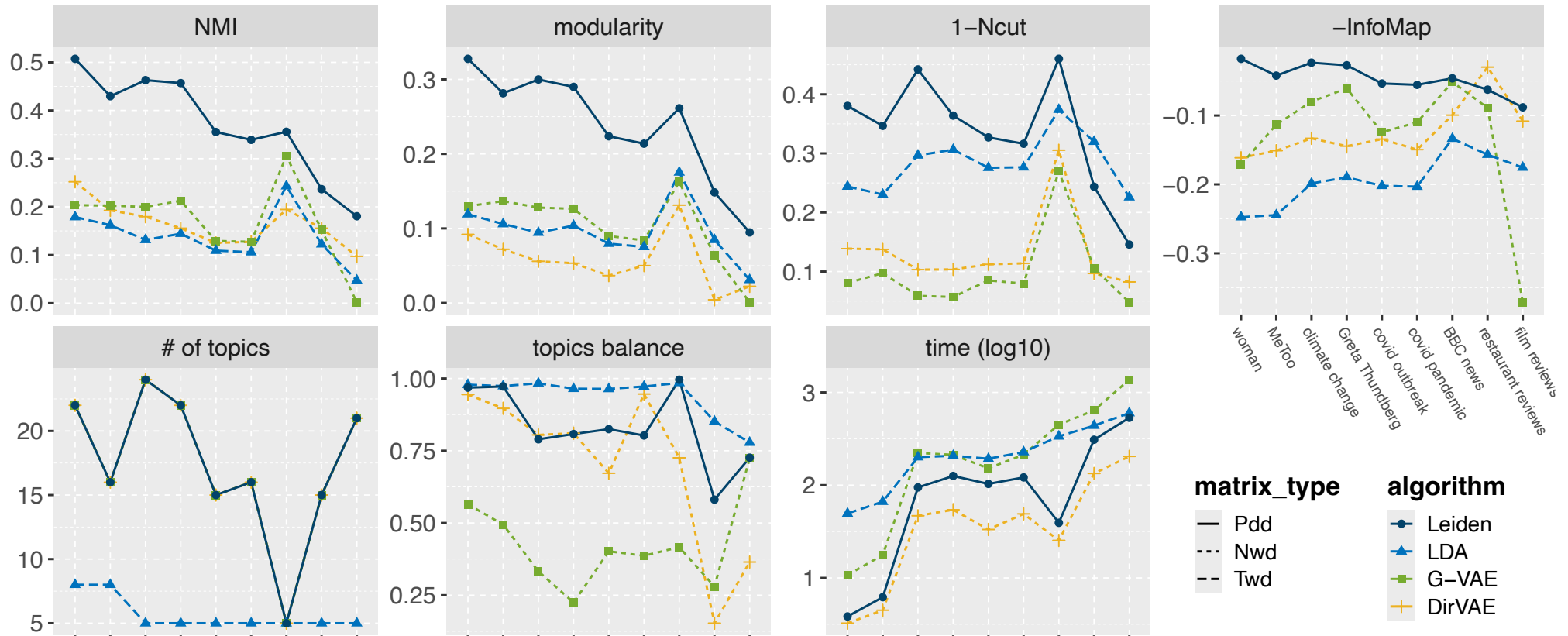
Our estimate of the **topic distribution** for the m th document!

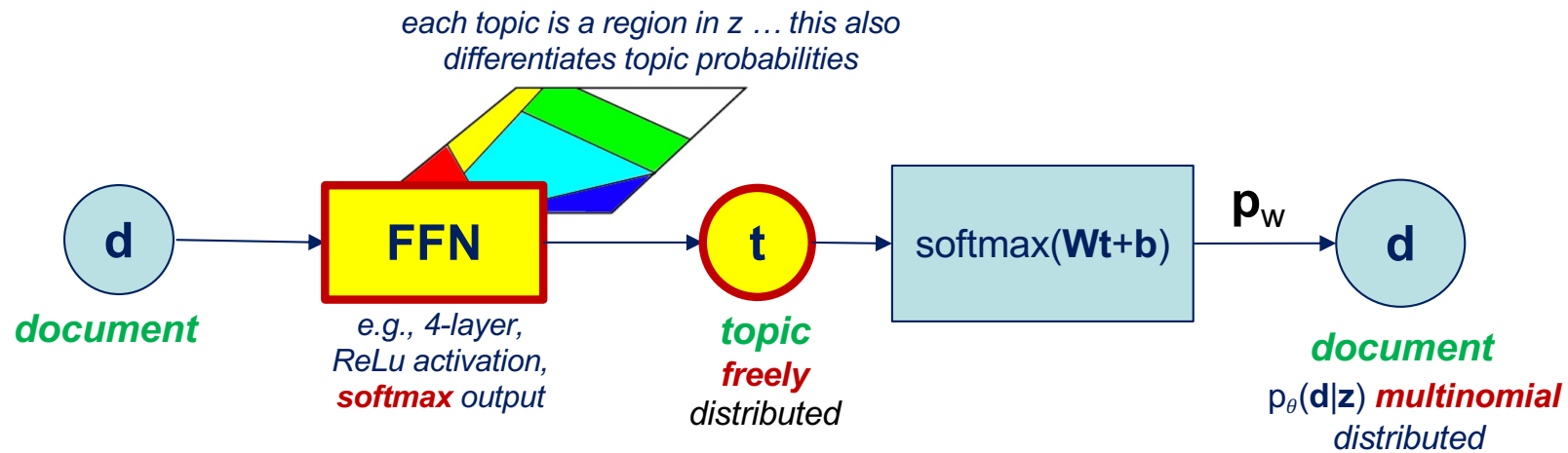
$$\mathbf{c}_m = \frac{1}{L} \sum_{\ell=1}^L \text{FFN}_1(\boldsymbol{\mu}_0(\mathbf{d}_m) + \boldsymbol{\sigma}_0(\mathbf{d}_m) \mathbf{n}_{m,\ell})$$



VAE at work

On a semantic network





one-hot-representation of a
document = number of occurrences
of words in the document

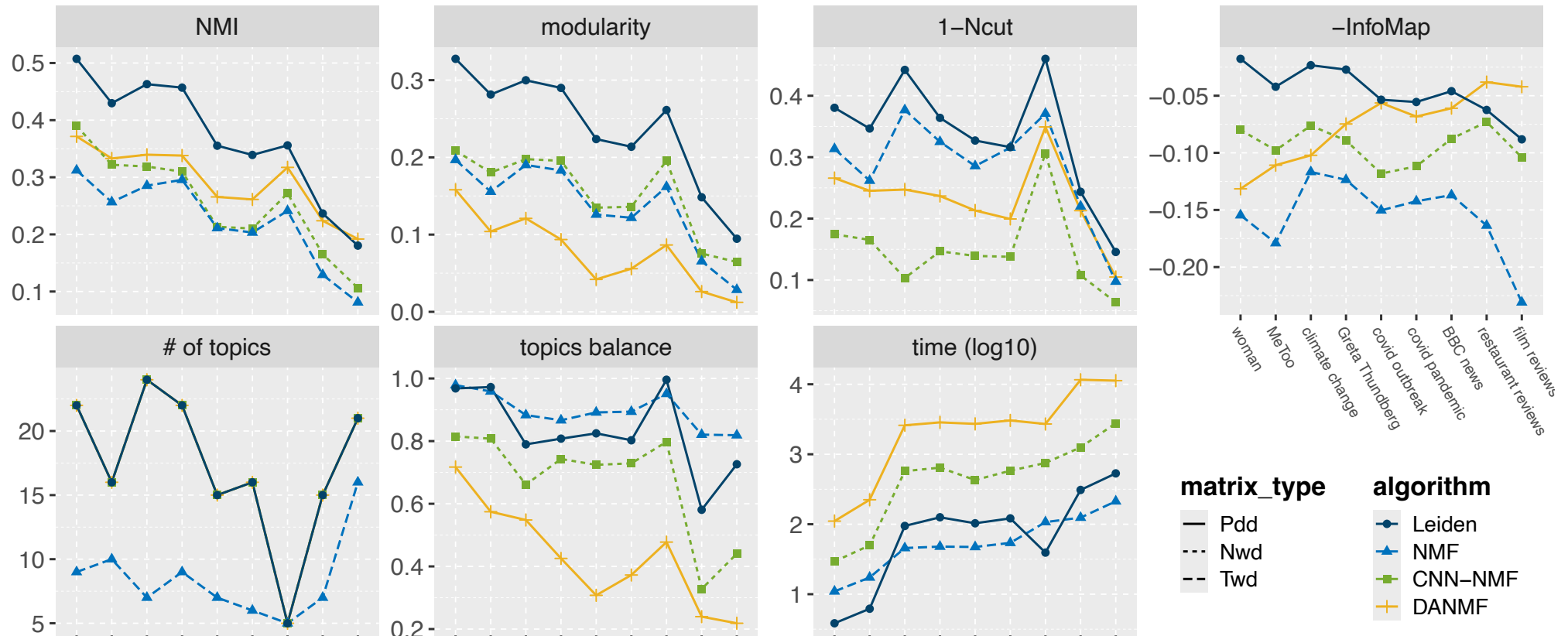
decoder model encoder model

$$\mathcal{L}(\theta, \phi) = \sum_m d_m^T \log(\text{softmax}(\mathbf{b} + \mathbf{W} \text{FFN}(\mathbf{d}_m)))$$



CNN-NMF at work

On a semantic network





- ❑ Naturally provides a soft topic assignment
- ❑ VAE – interesting approach
 - more flexible model than NMF or LDA
 - gives improvements
- ❑ Comparison – with Louvain
 - still far away
 - would be nice to see other **Deep Learning** approaches
 - ... your task! 😊

Transformer Architecture

with application to BERT, RoBERTa, OpenAI GPT



≡ Attention (machine learning)

Article [Talk](#)

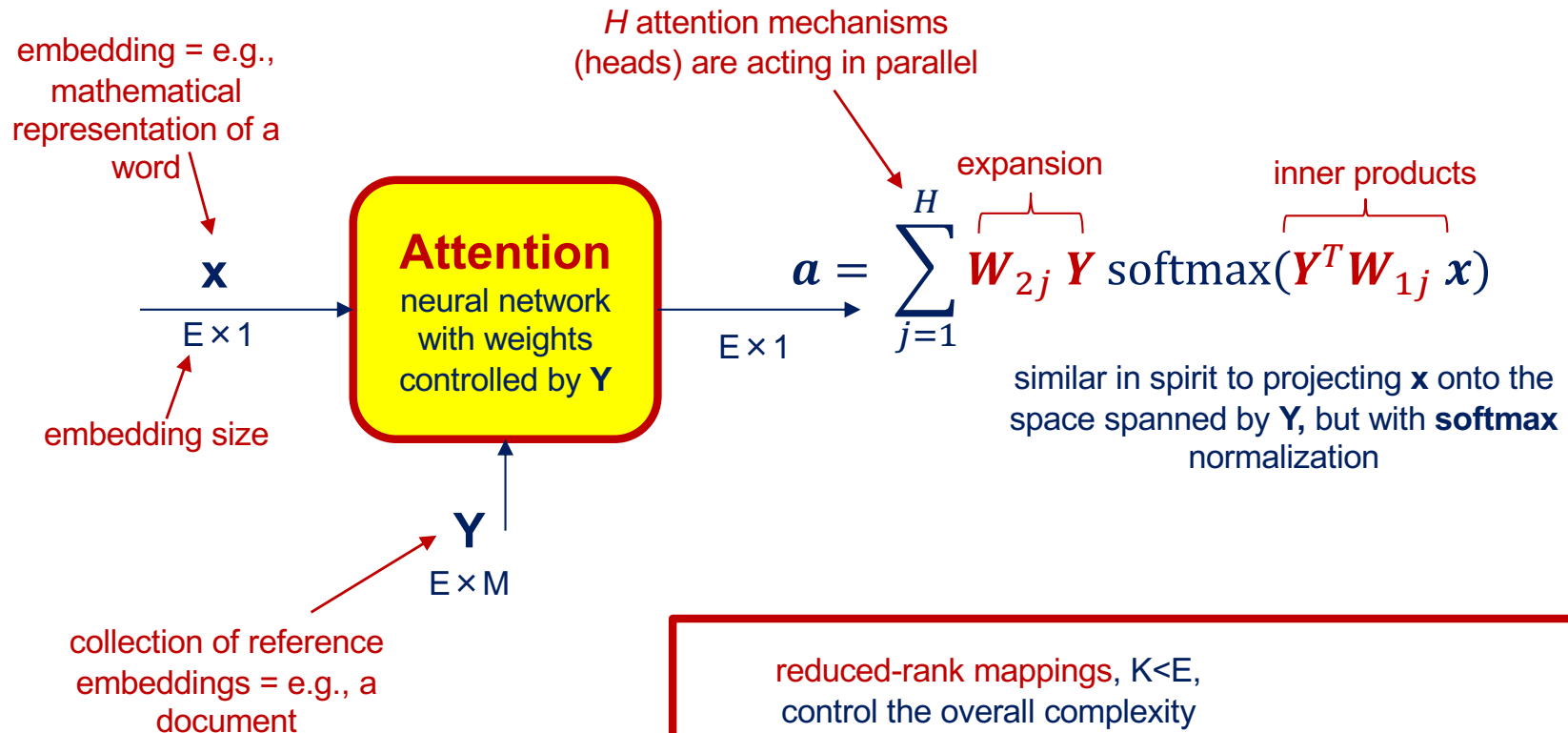
From Wikipedia, the free encyclopedia

In [artificial neural networks](#), **attention** is a technique that is meant to mimic [cognitive attention](#). This effect **enhances** some **parts of the input data while diminishing other parts** — the motivation being that the network should devote more focus to the important parts of the data, even though they may be small portion of an image or sentence. Learning which part of the data is more important than another depends on the context, and this is trained by [gradient descent](#).



The Attention Module

Vaswani, Ashish, et al. "Attention is all you need" (2017)



reduced-rank mappings, $K < E$,
control the overall complexity

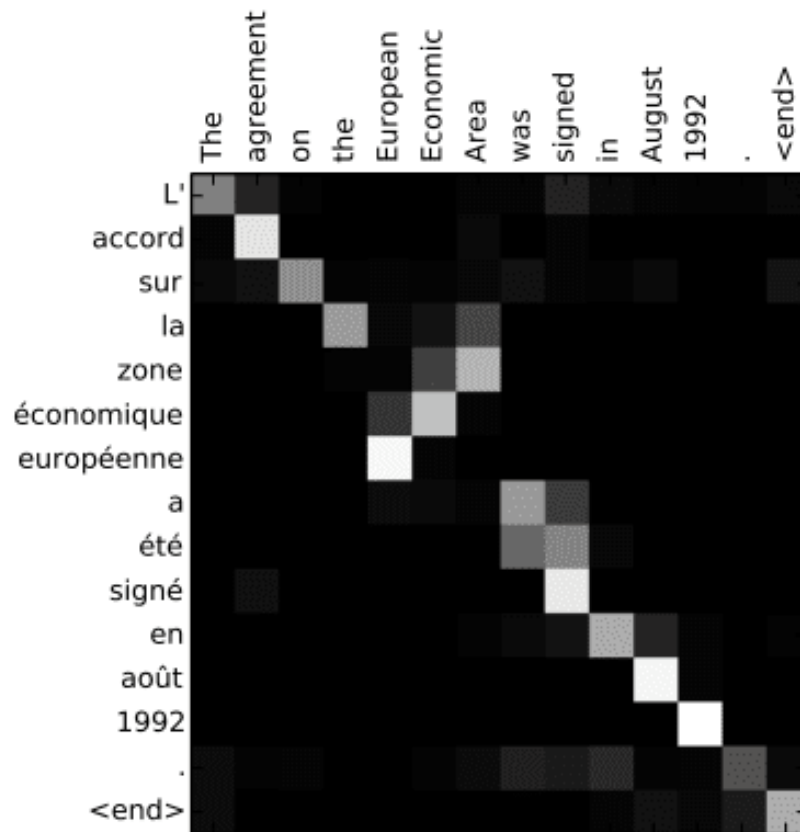
$$\begin{matrix} \boxed{W_{ij}} & = & \boxed{V_{ij}} & \boxed{U_{ij}} \\ E \times E & & E \times K & K \times E \end{matrix}$$

overall
complexity
 $4EKH$

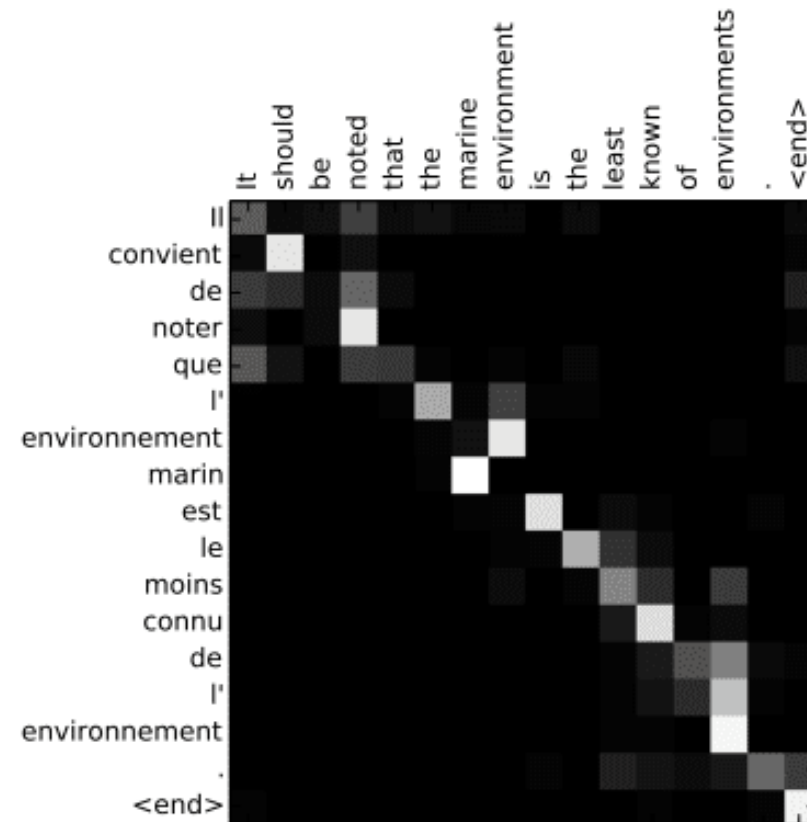


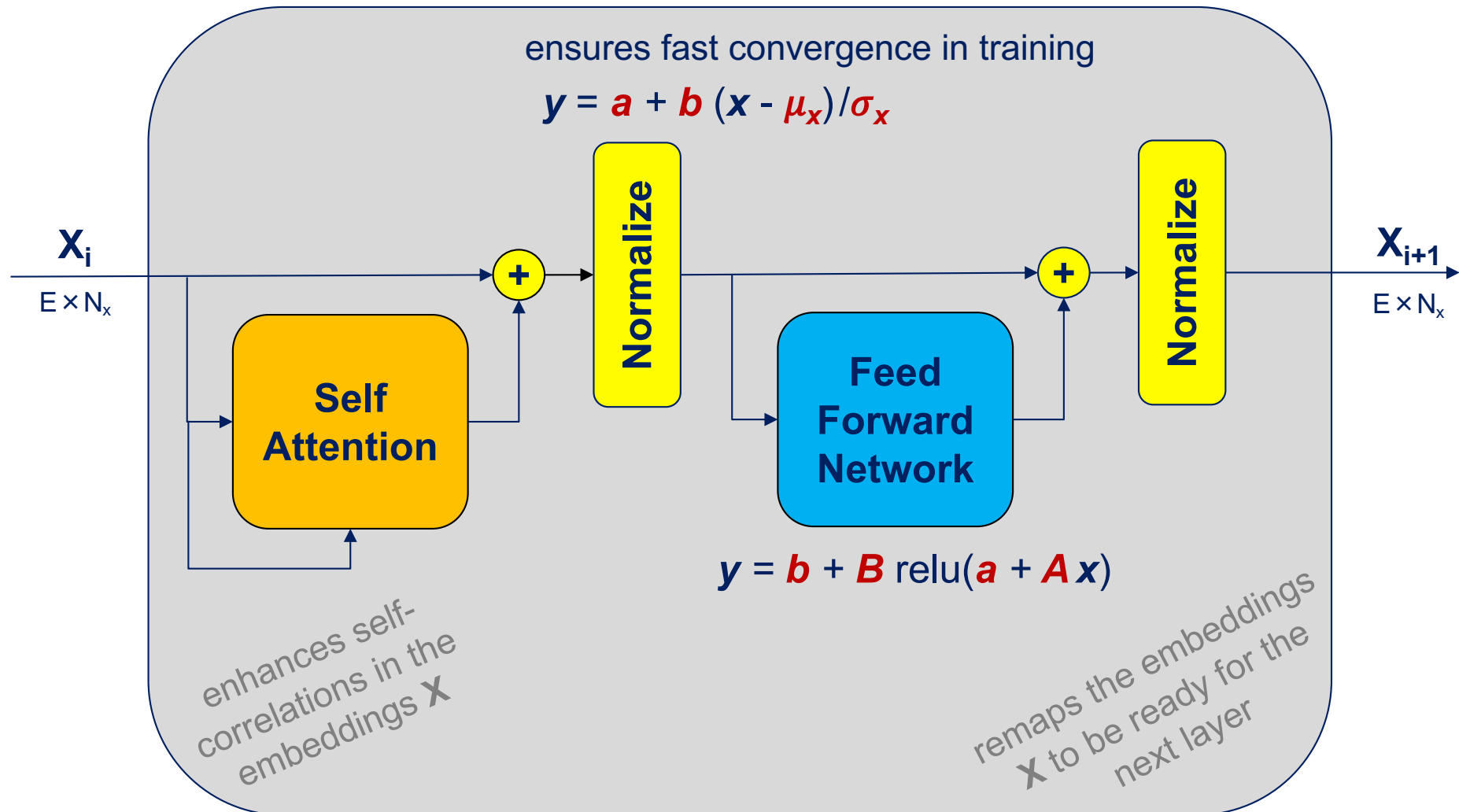
Visualizing Attention

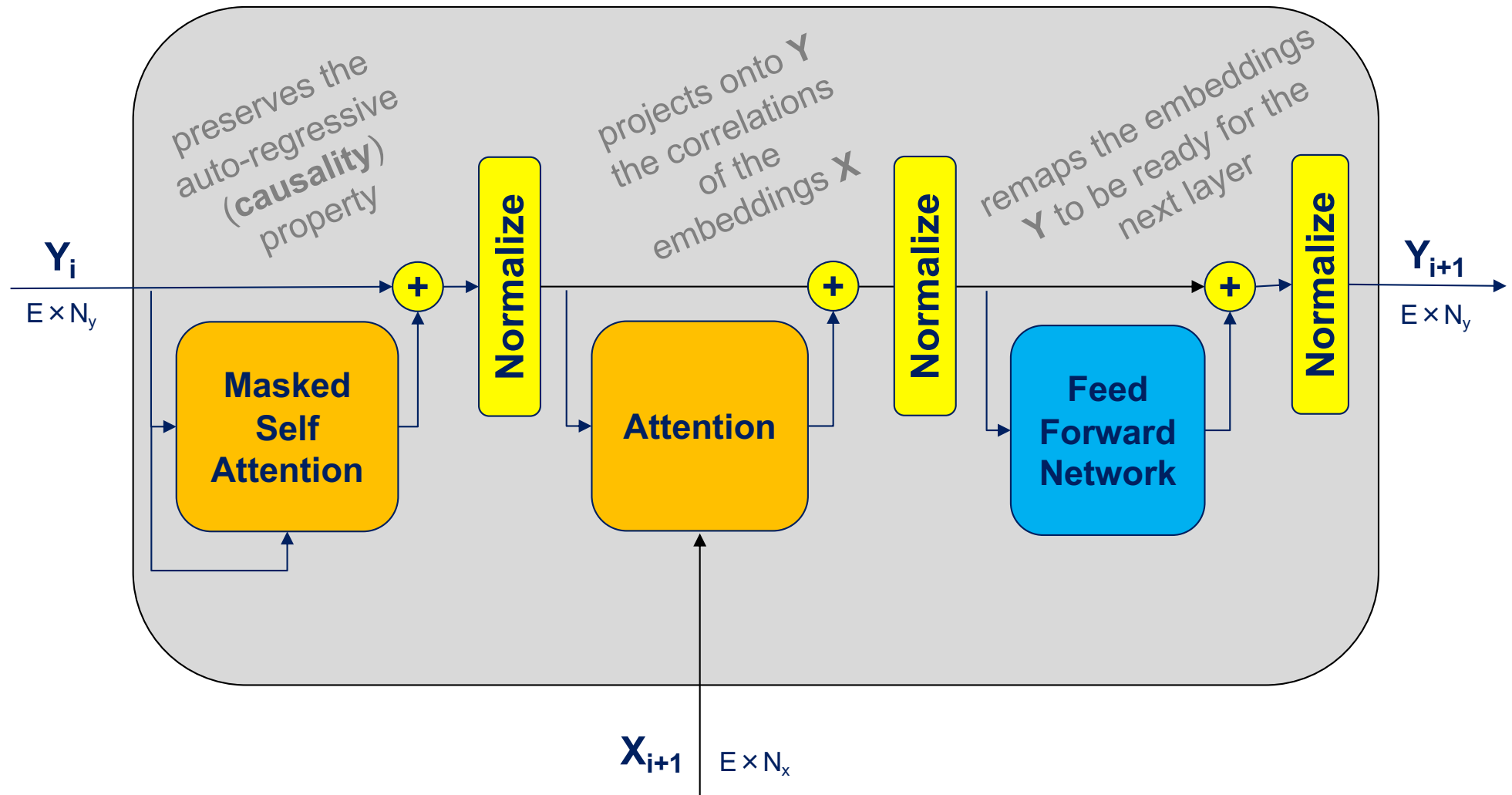
in a translation experiment (X English, Y French)



$$\text{softmax}(Y^T W_{1j} X)$$



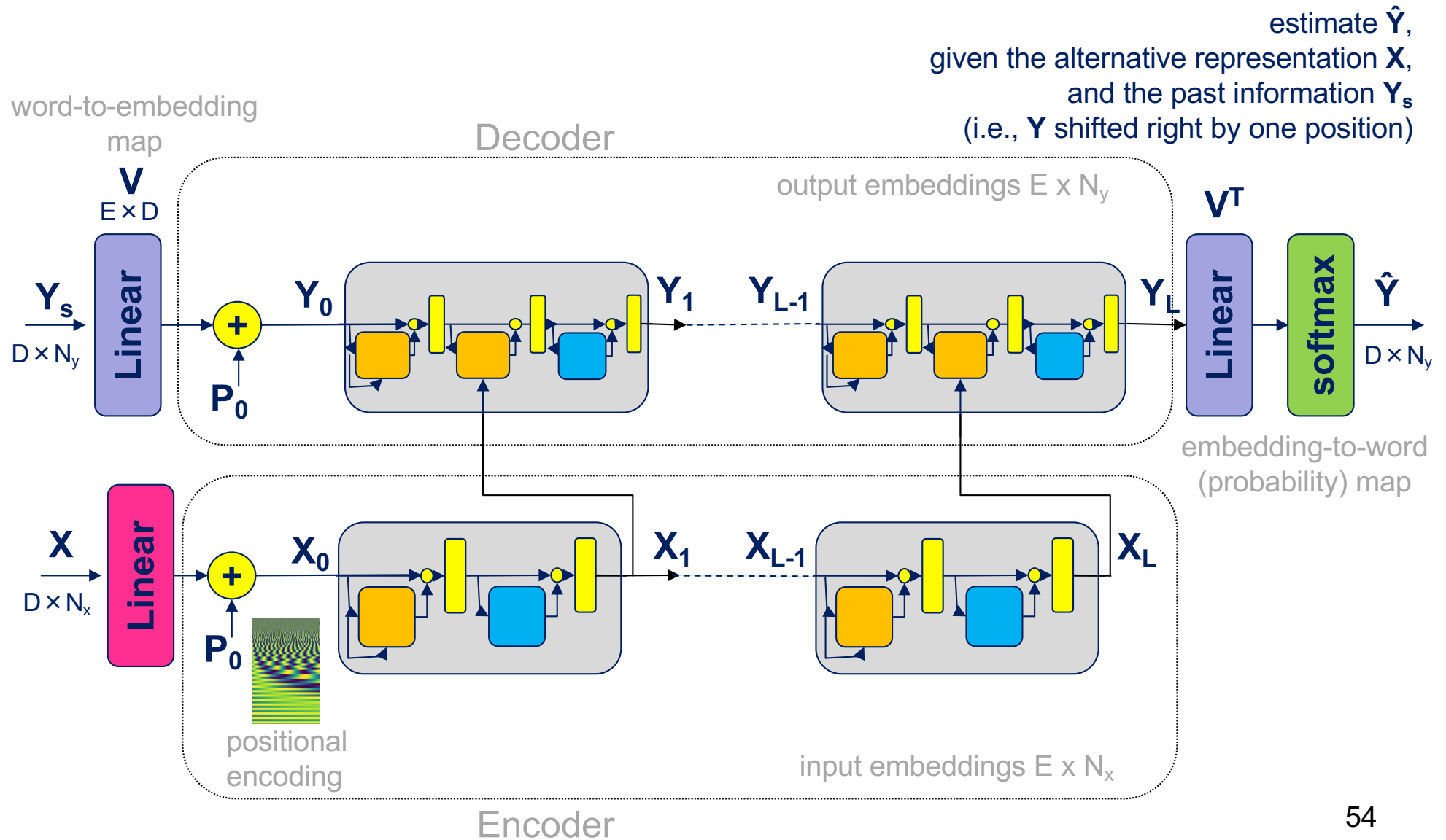






Transformer Architecture

Vaswani, Ashish, et al. "Attention is all you need" (2017)
Google's patent <https://patents.google.com/patent/US10452978B2/en>





UNIVERSITÀ
DEGLI STUDI
DI PADOVA

The Annotated Transformer

<http://nlp.seas.harvard.edu/2018/04/03/attention.html>
tensor2tensor library <https://github.com/tensorflow/tensor2tensor>



Members PI Code Publications

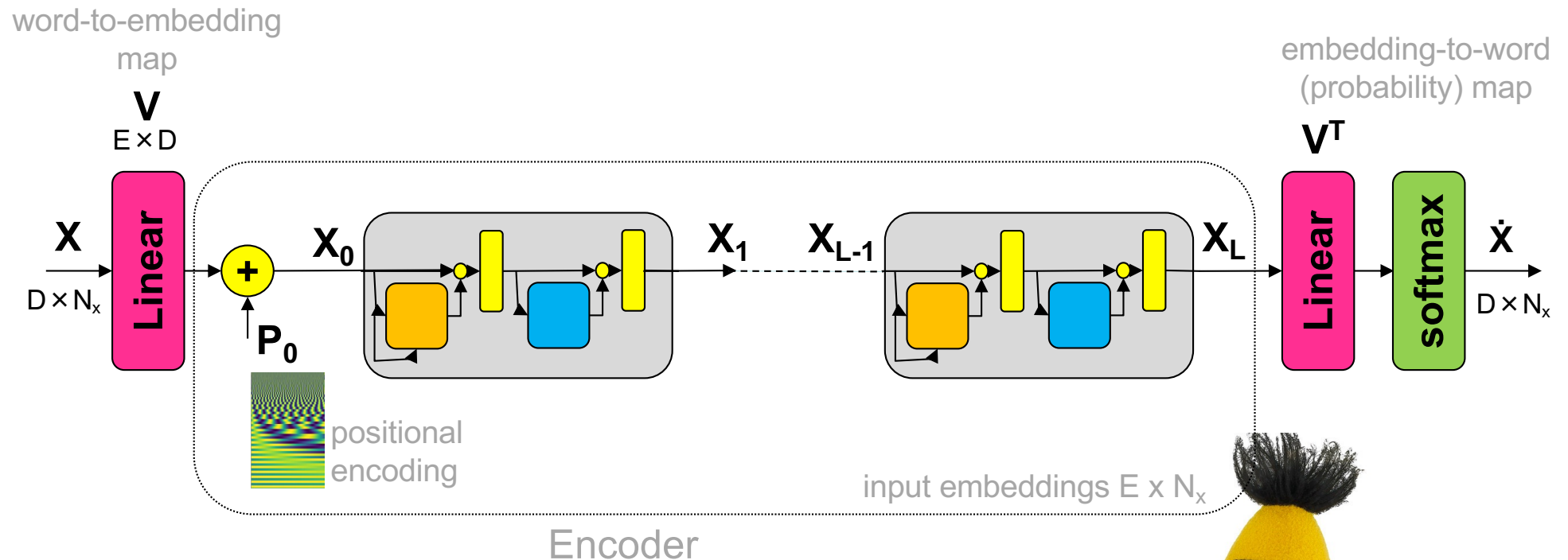
The Annotated Transformer

Apr 3, 2018

There is now a [new version](#) of this blog post updated for modern PyTorch.

```
from IPython.display import Image
Image(filename='images/aiayn.png')
```

Attention Is All You Need





BERT parameters

	Embeddings size E	Self-attention heads H	Head dimension K = E/H	FFN inner size I = 4E	Parameters per layer $12E^2+9E$	Layers L	Dictionary size D	Total parameters
BERT base	768	12	64	3072	7.1M	12	30.5K	110M
BERT large	1024	16	64	4096	12.6M	24	30.5K	340M

max tokens $N_x = 512$

Created by researchers at Google AI Language





BERT pre-training procedure

BooksCorpus (800M words) + English Wikipedia (2,500M words)

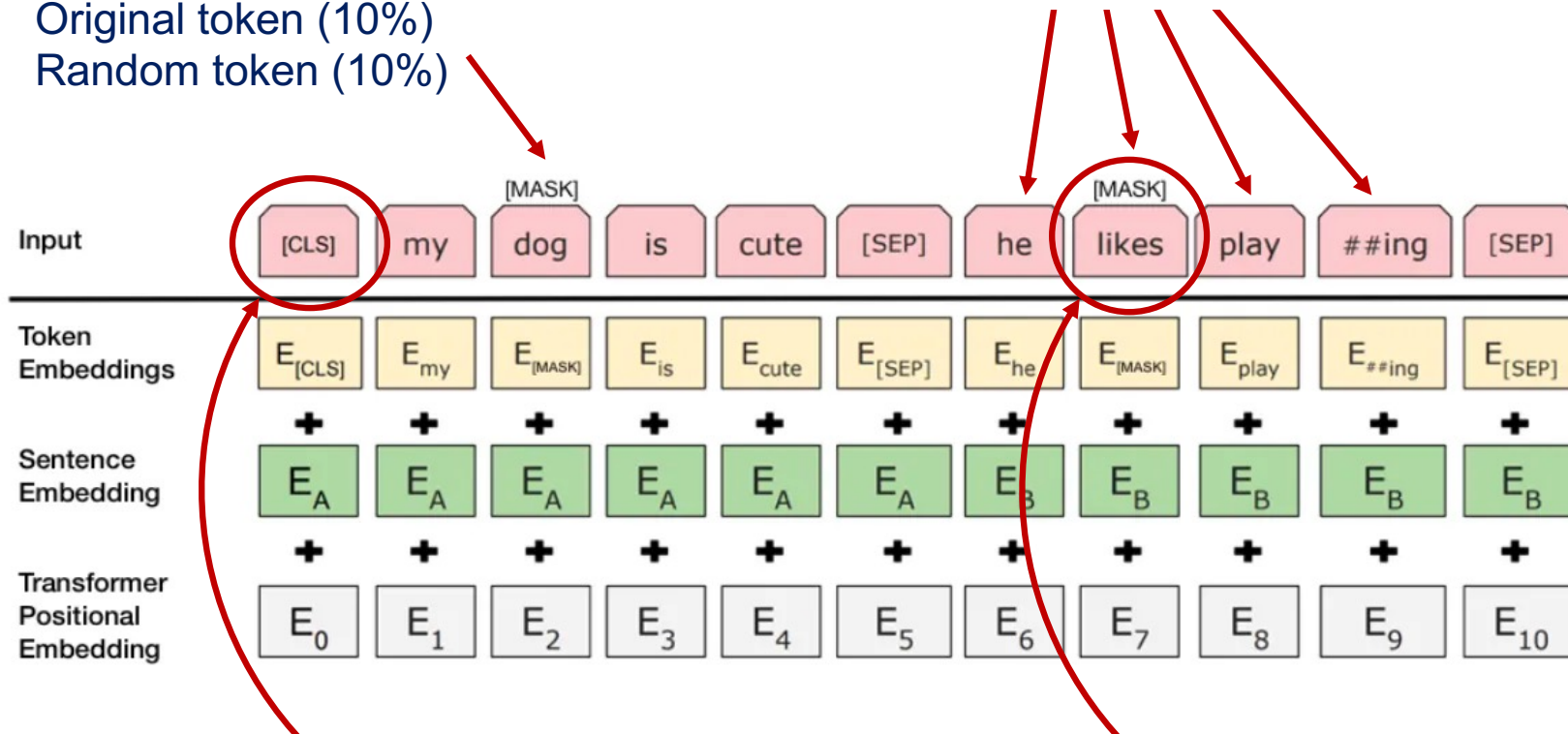
Masked Language Model

15% **masked tokens** replaced with:

- [MASK] token (80% of the times)
- Original token (10%)
- Random token (10%)

Next Sequence Prediction

- Next sequence (50% of the times)
- Random sequence (50%)



Output [CLS] fed into an additional output layer for softmax classification (of correct/wrong next sequence)

Output **masked tokens** fed into the output layer V^T and evaluated for probability of correct estimate



Larger training corpora (10x larger)

training on BookCorpus + Wikipedia and also CC-News, OpenWebText, Stories

Dynamic masking

training data was duplicated 10 times so that each sequence is masked in 10 different ways over the 40 epochs of training

Full-sentences without NSP loss

full sentences sampled contiguously from one or more documents, such that the total length is at most 512 tokens

Large mini-batches

A larger byte-level BPE (byte pair encoding) of 50K subword units

a hybrid between character- and word-level representations that allows handling the large vocabularies common in natural language corpora

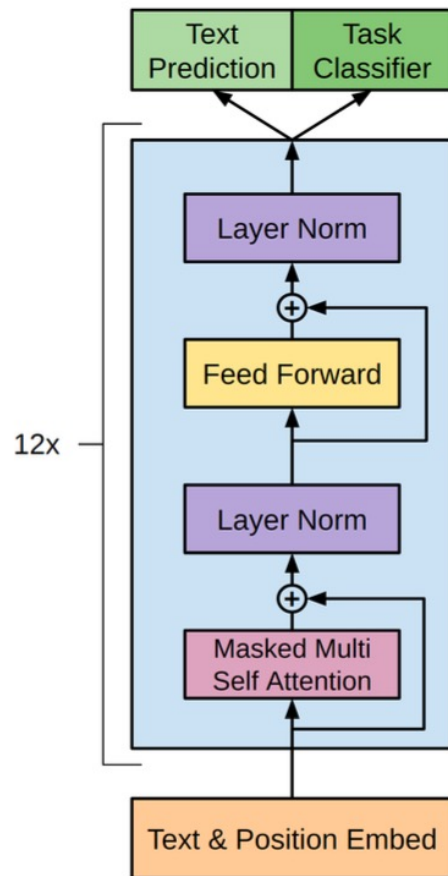


Generative Pre-Training (GPT)

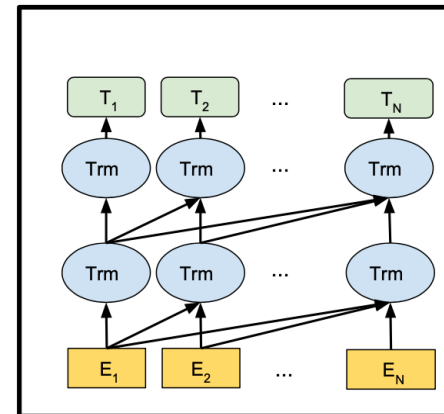
Radford, Alec, et al. "Improving language understanding by generative pre-training." (2018)

(unsupervised) pre-training on **Language Modelling (no mask)**

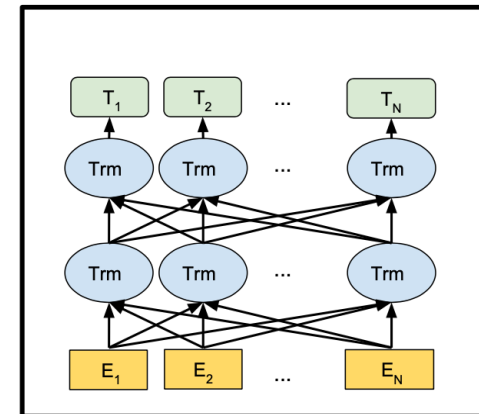
$$L_1(\mathcal{U}) = \sum_i \log P(u_i | u_{i-k}, \dots, u_{i-1}; \Theta)$$



GPT



BERT



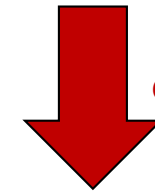
same parameters of BERT-base, but with **Masked Attention** trained on BookCorpus only



McCann et al. (2018)

language provides a flexible way to specify tasks, inputs, and outputs all as a sequence of symbols... it is therefore possible to train a single model with **sufficient capacity** to infer and perform many **different tasks**

model gets
complex!



data gets
larger!

Parameters	Layers	d_{model}	
117M	12	768	GPT, BERT-base
345M	24	1024	BERT-large
762M	36	1280	
1542M	48	1600	GPT-2

WebText

scraping all outbound links (45M links) from Reddit, a social media platform, which received at least 3 karma – exclude Wikipedia



increasingly larger data and model!

Model Name	n_{params}	n_{layers}	d_{model}	n_{heads}	d_{head}	Batch Size	Learning Rate
GPT-3 Small	125M	12	768	12	64	0.5M	6.0×10^{-4}
GPT-3 Medium	350M	24	1024	16	64	0.5M	3.0×10^{-4}
GPT-3 Large	760M	24	1536	16	96	0.5M	2.5×10^{-4}
GPT-3 XL	1.3B	24	2048	24	128	1M	2.0×10^{-4}
GPT-3 2.7B	2.7B	32	2560	32	80	1M	1.6×10^{-4}
GPT-3 6.7B	6.7B	32	4096	32	128	2M	1.2×10^{-4}
GPT-3 13B	13.0B	40	5140	40	128	2M	1.0×10^{-4}
GPT-3 175B or "GPT-3"	175.0B	96	12288	96	128	3.2M	0.6×10^{-4}

Layer normalization at the input (plus one at the output)

Sparse attention patterns

alternating dense and locally banded sparse attention patterns in the layers

Byte-level BPE (byte pair encoding) of 50K subword units

also prevent BPE from merging across character categories (to avoid dog, dog!, dog?)

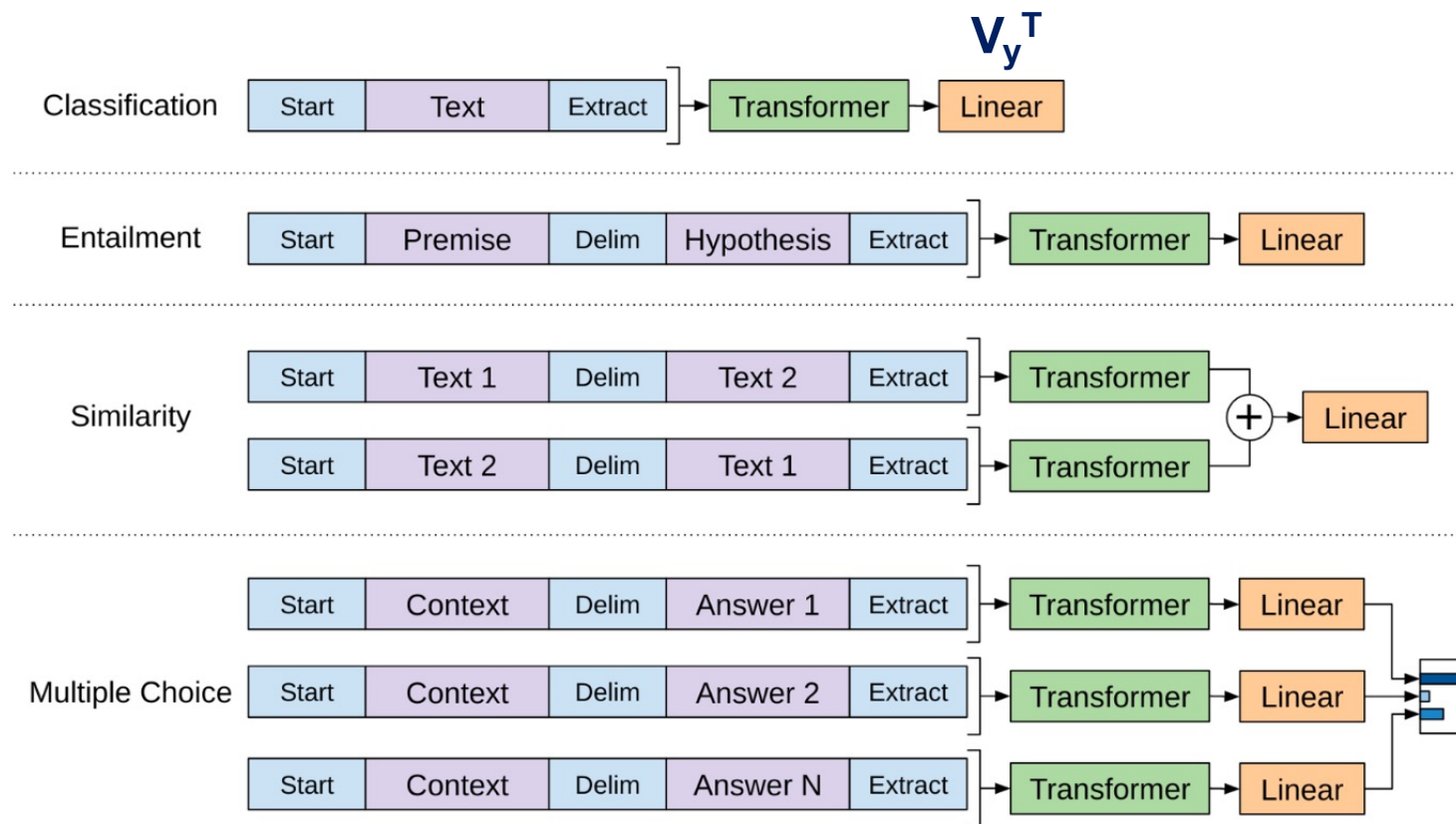
Modified initialization



Supervised fine-tuning

training on specific tasks

$$\log \text{softmax}(\mathbf{V}_y^T \mathbf{X}_L) \quad \xrightarrow{\text{Language Modelling loss}} \quad L_2(\mathcal{C}) + \lambda * L_1(\mathcal{C})$$





Task	Description	Possible approach
Masked language prediction	predict masked words in a text	This is what BERT model is pre-trained for
Text classification or Sentiment analysis	assign a label to a given sequence of text	Apply linear transform+softmax on K classes , and train the model for the specific classification task
Text translation	translate a text	Need to pre-train a full Transformer Architecture for this task
Summarization	generate a summary of a document	GPT example: context given by a document; then generate 100 tokens by top-2 random sampling (Fan et al., 2018), i.e., take at each step the most likely next word at random among the top-2 candidates; finally select first 3 sentences as abstract
Question answering	answer a question	GPT example: the context of the language model is seeded with example question answer pairs which helps the model infer the short answer style of the dataset
Document question answering	answer a question on a given text	GPT example: context seeded by a text; then as for question answering
Conversational	ChatBot	InstructGPT/ChatGPT: Fine-tuned models using reinforcement learning from human feedback



Hugging Face

<https://huggingface.co/docs/transformers/v4.29.1/en/index>

State-of-the-art Machine Learning
for PyTorch, TensorFlow, and JAX



PyTorch



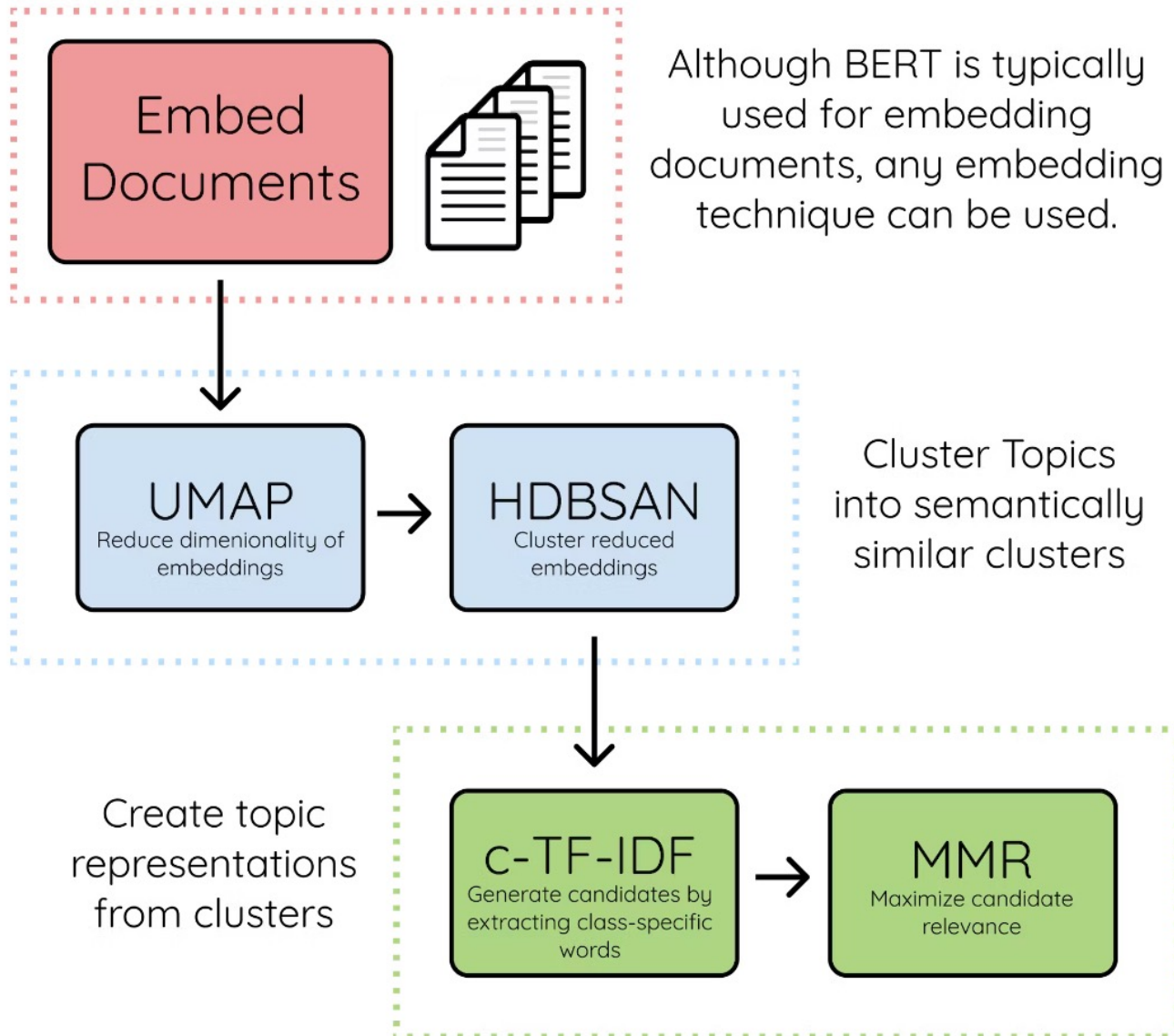
TensorFlow



ALBERT, BART, **BERT**, BigBird, BigBird-Pegasus, BioGpt, BLOOM, CamemBERT, CANINE, ConvBERT, CTRL, Data2VecText, DeBERTa, DeBERTa-v2, DistilBERT, ELECTRA, ERNIE, ErnieM, ESM, FlauBERT, FNet, Funnel Transformer, GPT-Sw3, **OpenAI GPT-2**, GPTBigCode, GPT Neo, GPT NeoX, GPT-J, I-BERT, LayoutLM, LayoutLMv2, LayoutLMv3, LED, LiLT, LLaMA, Longformer, LUKE, MarkupLM, mBART, MEGA, Megatron-BERT, MobileBERT, MPNet, MVP, Nezha, Nyströmformer, OpenLlama, **OpenAI GPT**, OPT, Perceiver, PLBart, QDQBert, Reformer, RemBERT, **RoBERTa**, RoBERTa-PreLayerNorm, RoCBert, RoFormer, SqueezeBERT, TAPAS, Transformer-XL, XLM, XLM-RoBERTa, XLM-RoBERTa-XL, XLNet, X-MOD, YOSO

BERT Topic

exploiting embeddings for topic detection

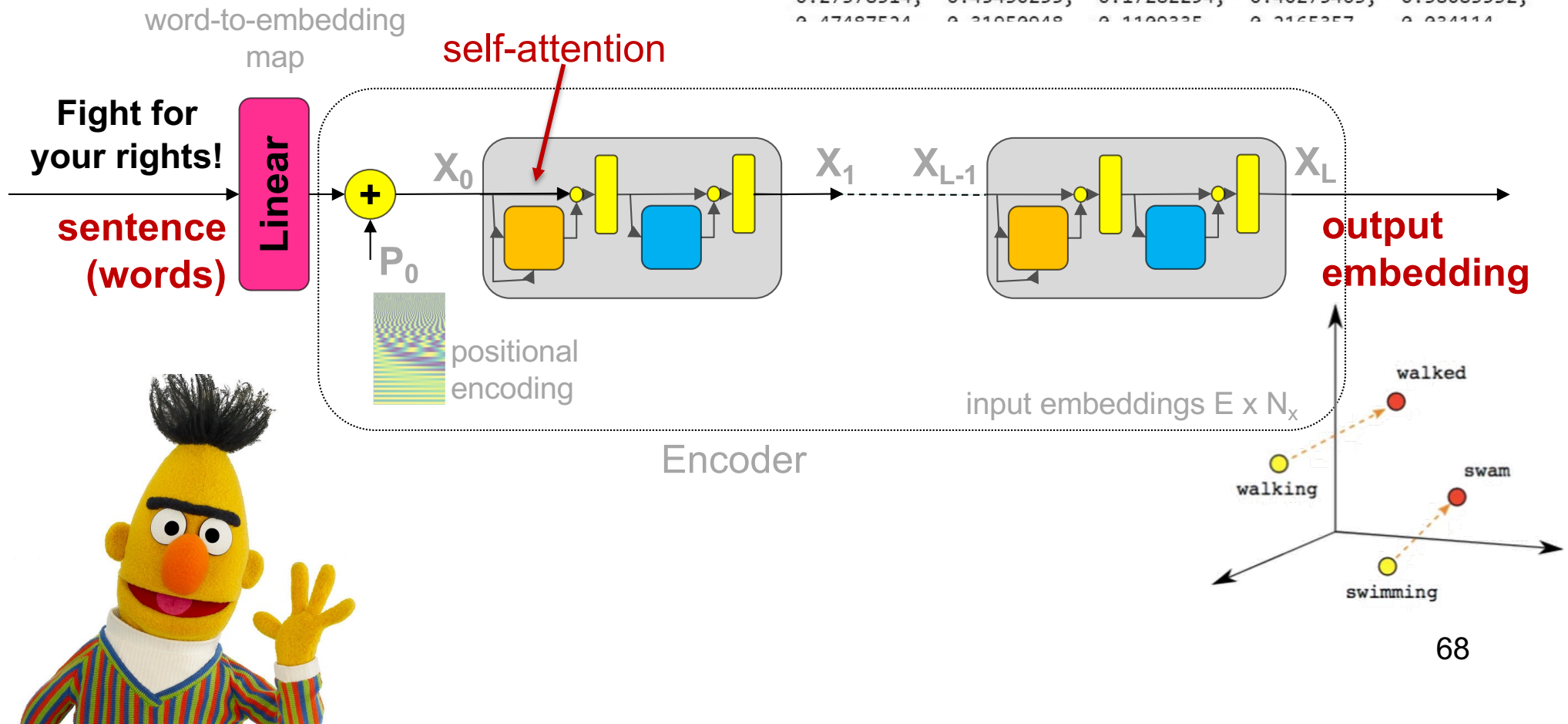




1. Embed documents

Using BERT

```
array([-0.5968882, -0.33086956, -0.32643065, -0.3670732,  0.628059 ,  
       -0.3692328, -0.37902787, -0.12308089, -0.38124698, -0.03940517,  
        0.2260839,  0.10852845, -0.2873811, -0.42781743,  0.06604357,  
       -0.07114276, -0.29775023, -0.99628943, -0.54497653, -0.11718027,  
       -0.15935768,  0.09587188, -0.2503798,  0.06768776,  0.3311586 ,  
        0.43098116,  0.06936899,  0.24311952,  0.14515282,  0.19245838,  
        0.10462623, -0.45676082,  0.5662387,  0.69908774,  0.48064467,  
        0.27378514, -0.45430255,  0.17282294, -0.40275463, -0.38083532,  
        0.17187534,  0.31050048,  0.1100335,  0.3165357,  0.0341114
```

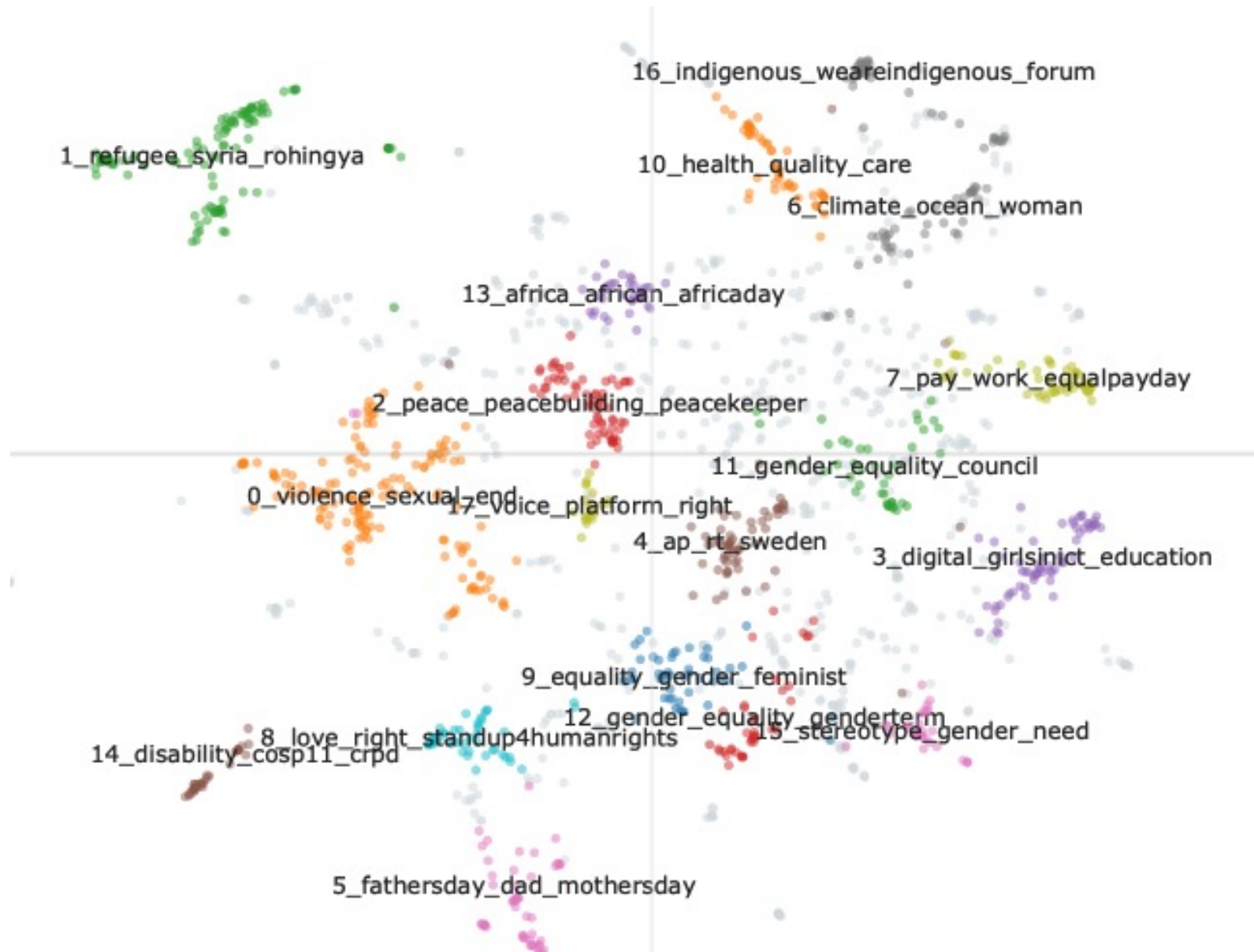




2. Reinterpret embeddings

Using UMAP

```
bert_model.visualize_documents(docs)
```





```
!pip install bertopic
from bertopic import BERTopic
from sentence_transformers import SentenceTransformer
```

```
sentence_model = SentenceTransformer("all-MiniLM-L6-v2")
bert_model = BERTopic(embedding_model=sentence_model,
                      min_topic_size=20, nr_topics='auto')
```

initialise model

```
docs = list(df2["text_sup_clean"])
topics, probabilities = bert_model.fit_transform(docs)
```

fit model

```
topics = bert_model.reduce_outliers(docs, topics)
```

reduce outliers

```
# extract community assignments
C = sps.csr_matrix((len(topics), max(topics)+2))
for i in range(C.shape[1]):
    C[np.array(topics)==(i-1), i] = 1

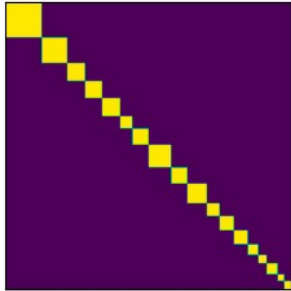
# remove zero assignments
C = C[:, np.unique(scipy.sparse.find(C)[1])]

```

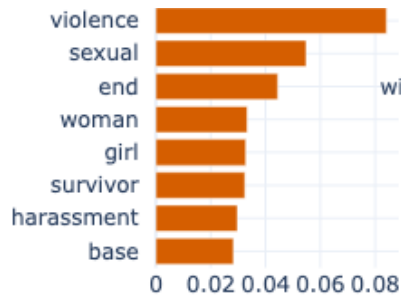
extract C from topic
assignment



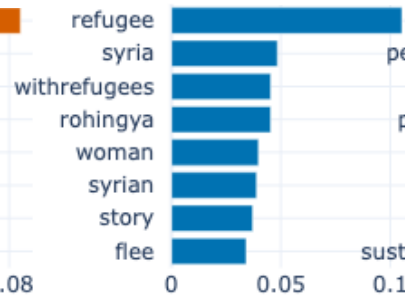
bert_model.visualize_barchart() #metoo2018



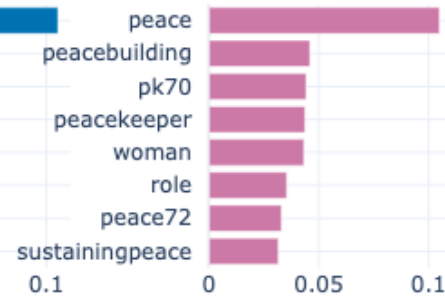
sexual violence



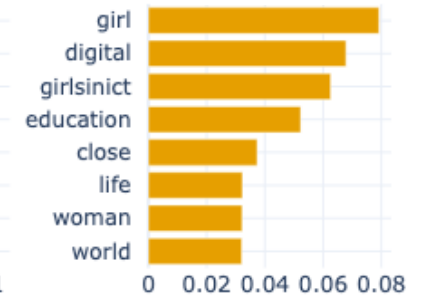
refugees



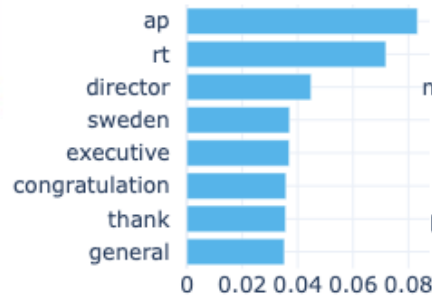
peace



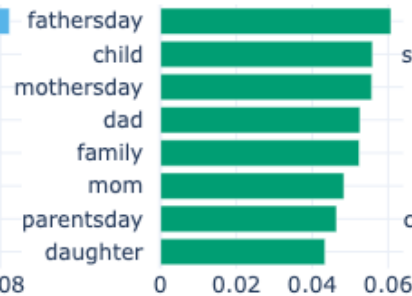
girlsinit



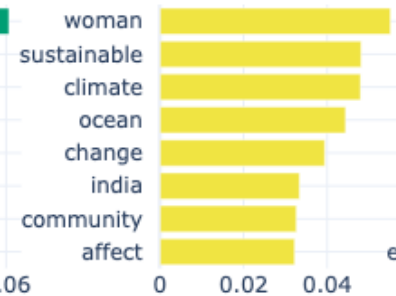
executive director



mothersday



sustainability



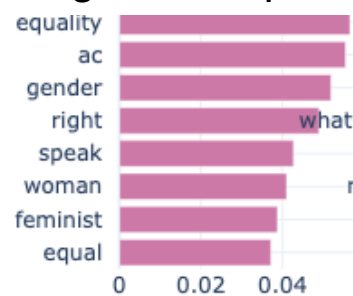
equal pay



discrimination



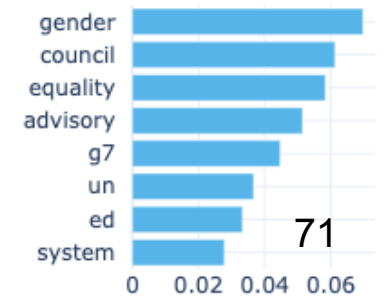
gender equality



whatwomenwant



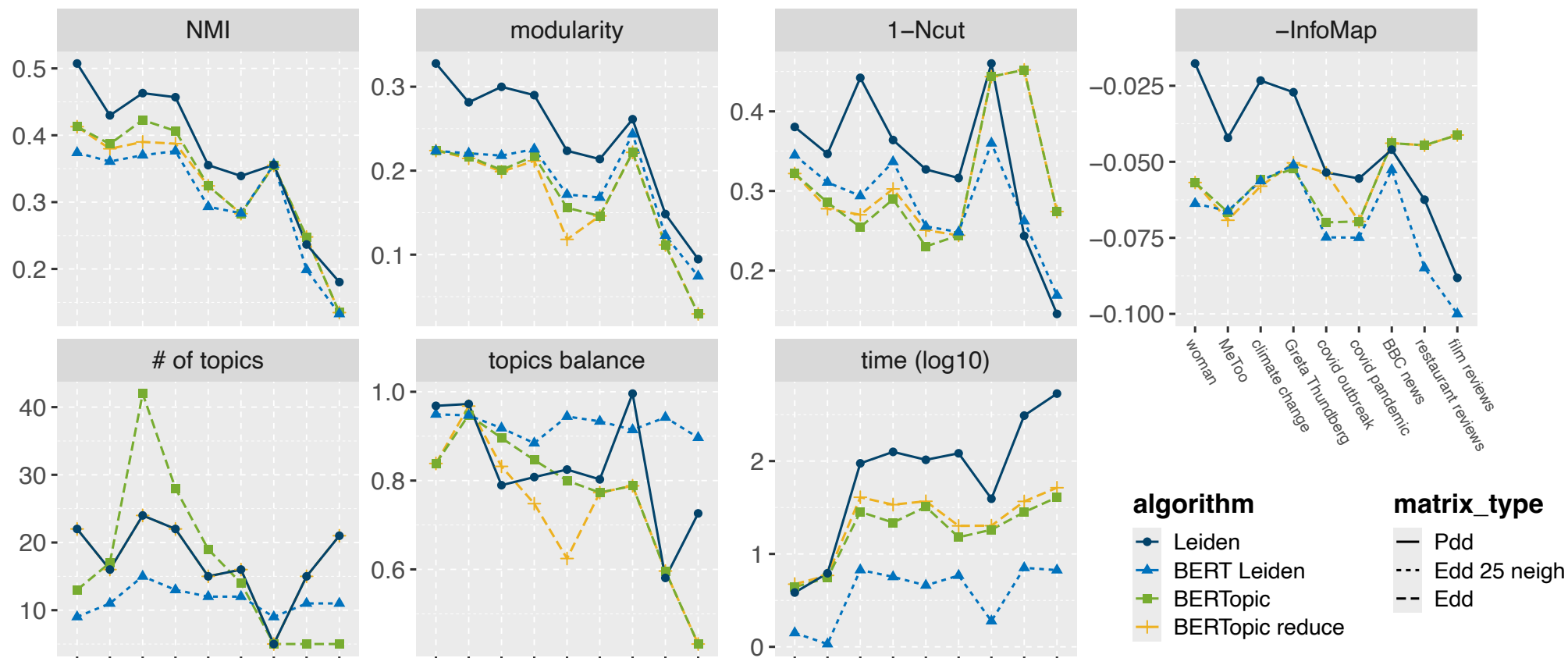
politics





BERTopic at work

On a semantic network





- ❑ Naturally provides a hard topic assignment
- ❑ Useful tool
- ❑ More readable output with deep cleaned text but same performance
- ❑ Comparison – with Louvain
 - weaker in general, especially in modularity
 - equivalent NMI = relevant topics
 - lower modularity = the documents that identify the topics are less distinguishable
 - higher complexity involved
 - less balanced topics, but generally meaningful
 - topics correlated with Louvain

Wrap-up

on topic detection



- ❑ What available tools should be used
 - Louvain & BERTopic
 - compare their performance through NMI, modularity, etc.
- ❑ What available tools should **NOT** be used
 - InfoMap, NMF & LDA
 - they show poor performance
- ❑ What would be nice to see implemented
 - soft Louvain made fast
 - advanced SMBs
 - deep learning models