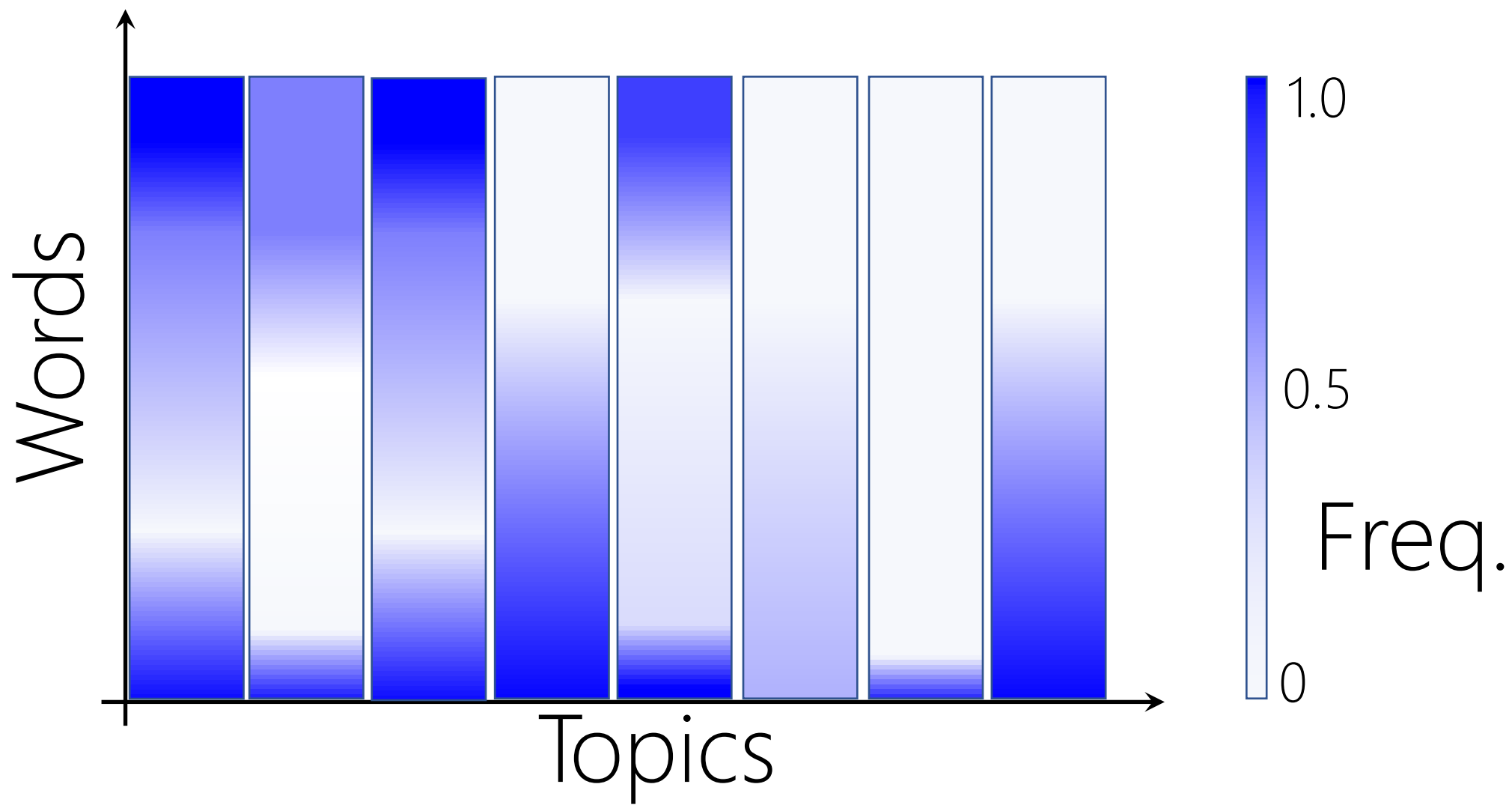


Probabilistic Topic Models

Probabilistic Topic Models

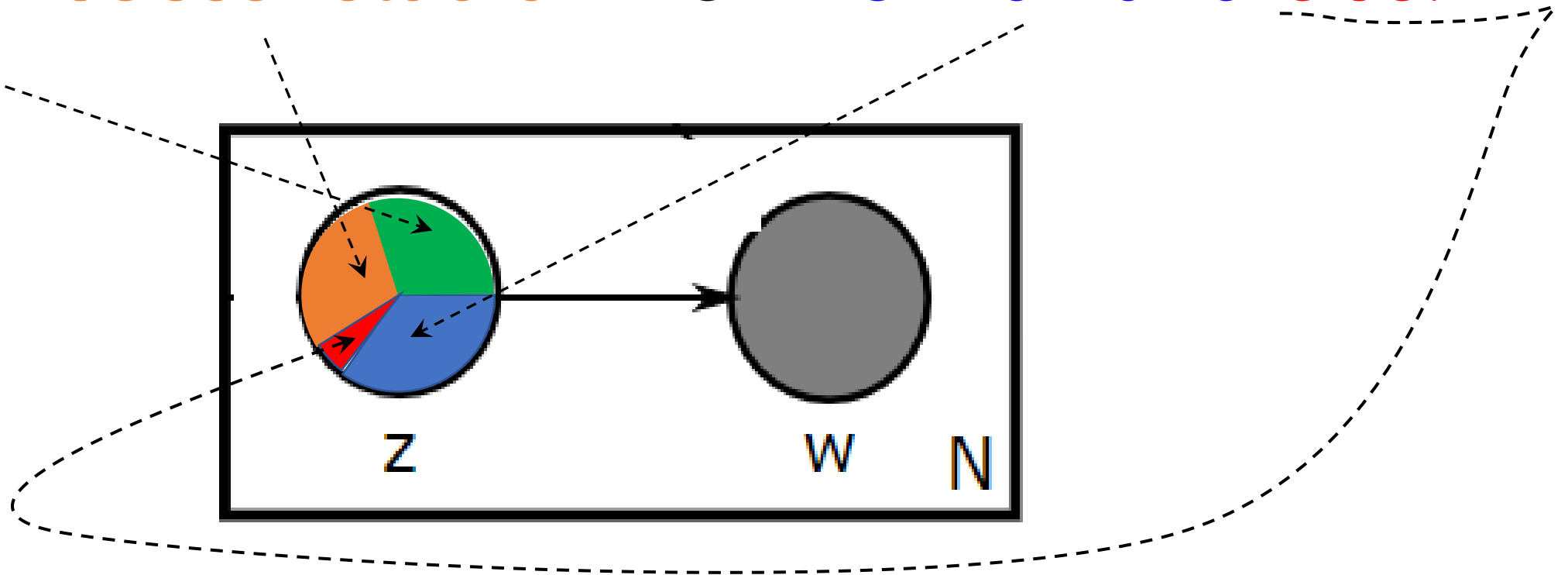


Probabilistic Topic Models

Let's write a document?

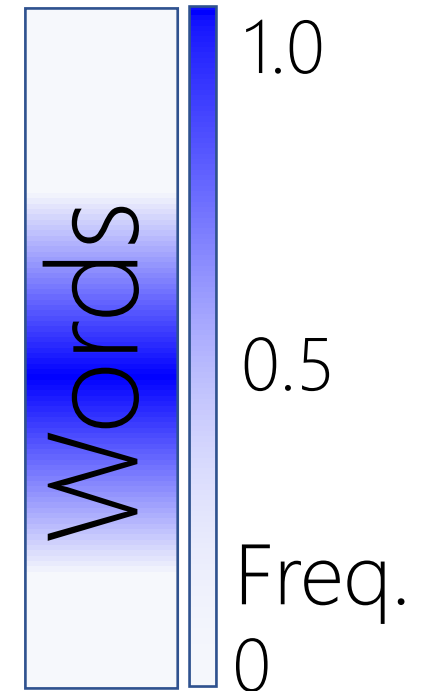
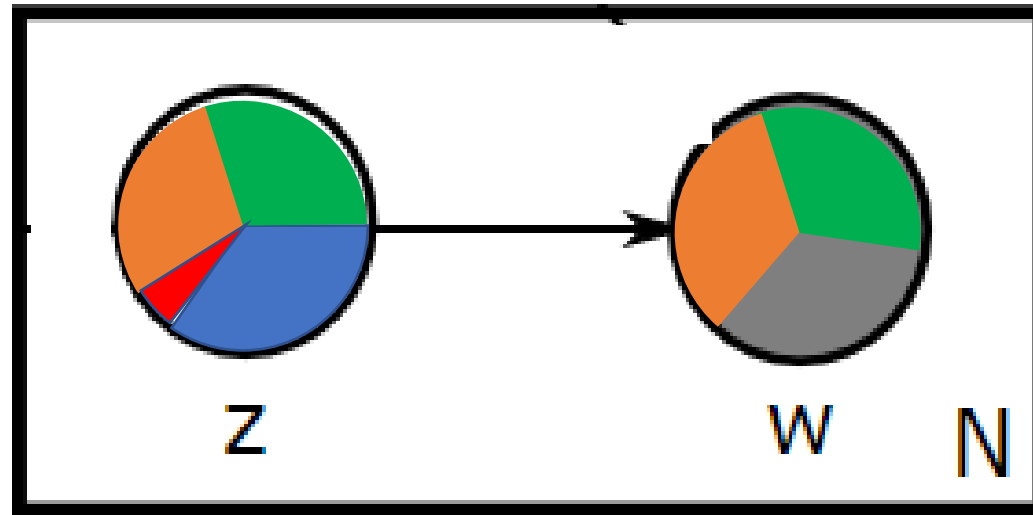
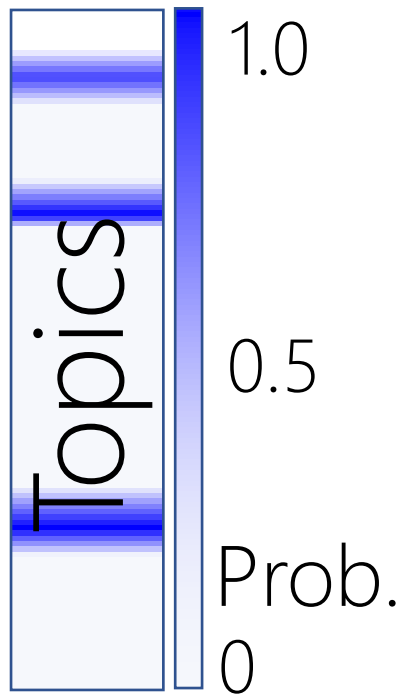
About what? *topic*?

Concert in soccer stadium for Nowruz and else.



Probabilistic Topic Models

- A document is about *all* topics but different distributions over topics
- A topic have *all* words but with different distributions over words



Latent Dirichlet Allocation

David M. Blei

*Computer Science Division
University of California
Berkeley, CA 94720, USA*

BLEI@CS.BERKELEY.EDU

Andrew Y. Ng

*Computer Science Department
Stanford University
Stanford, CA 94305, USA*

ANG@CS.STANFORD.EDU

Michael I. Jordan

*Computer Science Division and Department of Statistics
University of California
Berkeley, CA 94720, USA*

JORDAN@CS.BERKELEY.EDU

Editor: John Lafferty|

Probabilistic Topic Models

Assumptions:

1) #topics = $k \rightarrow Z = \{\text{topics } \mathbf{z}_i\}, |Z| = k$, for the whole corpus.

Probabilistic Topic Models

Assumptions:

- 1) #topics = $k \rightarrow Z = \{\text{topics } \mathbf{z}_i\}, |Z| = k$, for the whole corpus.
- 2) Each doc is about all k topics but with different distributions
 $C = \{\text{docs } \mathbf{d}_i \mid \mathbf{d}_i = [d_{i1}, d_{i2}, \dots, d_{ik}]\}_{i=1}^M = [\text{Matrix}]_{M \times k}$

Probabilistic Topic Models

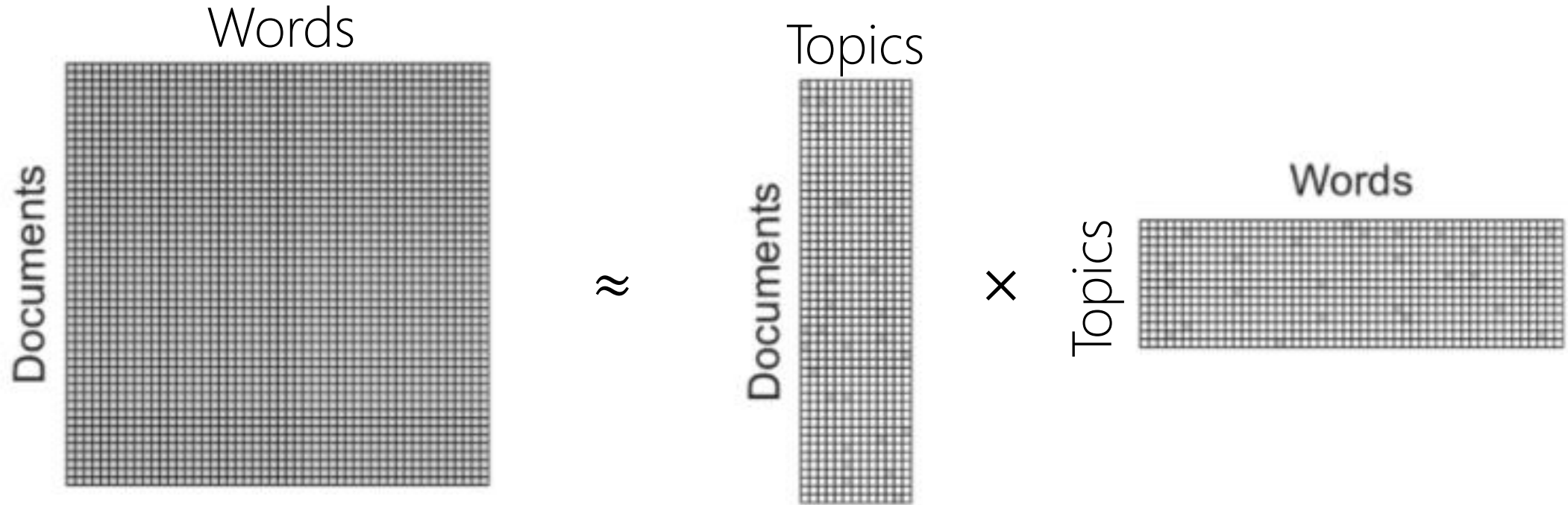
Assumptions:

1) #topics = $k \rightarrow Z = \{\text{topics } \mathbf{z}_i\}, |Z| = k$, for the whole corpus.

2) Each doc is about all k topics but with different distributions
 $C = \{\text{docs } \mathbf{d}_i \mid \mathbf{d}_i = [d_{i1}, d_{i2}, \dots, d_{ik}]\}_{i=1}^M = [\text{Matrix}]_{M \times k}$

3) Each topic \mathbf{z}_i $1 \leq i \leq k$, has all words but with different distributions
 $Z = \{\text{topics } \mathbf{z}_i \mid \mathbf{z}_i = [z_{i1}, z_{i2}, \dots, z_{iV}]\}_{i=1}^k = [\text{Matrix}]_{k \times V}$

Probabilistic Topic Models

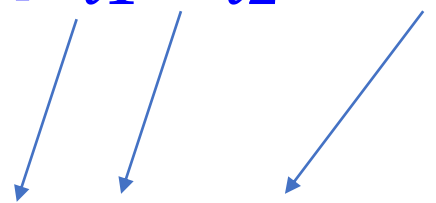


LDA can be seen as non-Negative Matrix Factorization
But learning algorithm is probabilistic!

Probabilistic Topic Models

Initializations:

- 1) For d_i in $C = \{\text{docs } d_i \mid d_i = [d_{i1}, d_{i2}, \dots, d_{ik}]\}_{i=1}^M = [\text{Matrix}]_{M \times k}$
- $d_i = [d_{i1}, d_{i2}, \dots, d_{ik}] \sim \text{Dirichlet}(\alpha = [\alpha_1, \alpha_2, \dots, \alpha_k])$


$$f(x_1, \dots, x_K; \alpha_1, \dots, \alpha_K) = \frac{1}{B(\alpha)} \prod_{i=1}^K x_i^{\alpha_i - 1}$$

where $\{x_k\}_{k=1}^{k=K}$ belong to the standard $K - 1$ **simplex**, or in other words:

$$\sum_{i=1}^K x_i = 1 \text{ and } x_i \geq 0 \text{ for all } i \in \{1, \dots, K\}$$

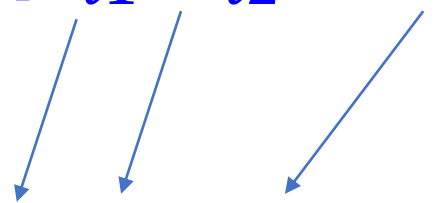


We've seen this before. where?

Probabilistic Topic Models

Initializations:

- 1) For d_i in $C = \{\text{docs } d_i \mid d_i = [d_{i1}, d_{i2}, \dots, d_{ik}]\}_{i=1}^M = [\text{Matrix}]_{M \times k}$
- $d_i = [d_{i1}, d_{i2}, \dots, d_{ik}] \sim \text{Dirichlet}(\alpha = [\alpha_1, \alpha_2, \dots, \alpha_k])$


$$f(x_1, \dots, x_K; \alpha_1, \dots, \alpha_K) = \frac{1}{B(\alpha)} \prod_{i=1}^K x_i^{\alpha_i - 1}$$

where $\{x_k\}_{k=1}^{k=K}$ belong to the standard $K - 1$ **simplex**, or in other words:

$$\sum_{i=1}^K x_i = 1 \text{ and } x_i \geq 0 \text{ for all } i \in \{1, \dots, K\}$$



Softmax! Transforming to vector of probs. But here we generate it!

Pólya's urn [\[edit source \]](#)

Consider an urn containing balls of K different colors. Initially, the urn contains α_1 balls of color 1, α_2 balls of color 2, and so on. Now perform N draws from the urn, where after each draw, the ball is placed back into the urn with an additional ball of the same color. In the limit as N approaches infinity, the proportions of different colored balls in the urn will be distributed as $\text{Dir}(\alpha_1, \dots, \alpha_K)$.^[20]

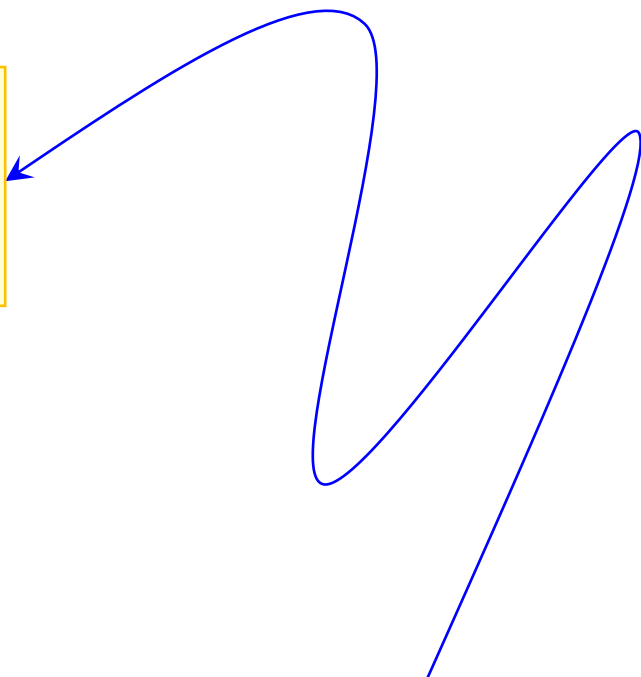
For a formal proof, note that the proportions of the different colored balls form a bounded $[0, 1]^K$ -valued [martingale](#), hence by the [martingale convergence theorem](#), these proportions converge [almost surely](#) and [in mean](#) to a limiting random vector. To see that this limiting vector has the above Dirichlet distribution, check that all mixed [moments](#) agree.

Each draw from the urn modifies the probability of drawing a ball of any one color from the urn in the future. This modification diminishes with the number of draws, since the relative effect of adding a new ball to the urn diminishes as the urn accumulates increasing numbers of balls.

Probabilistic Topic Models

Initializations:

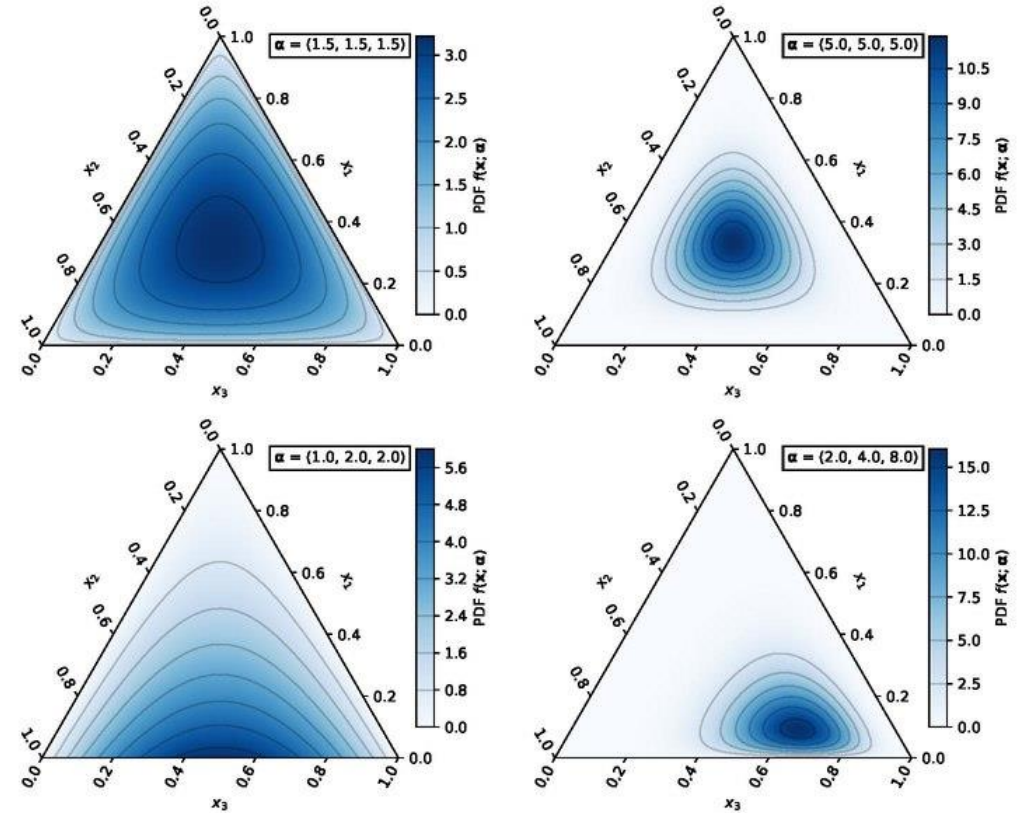
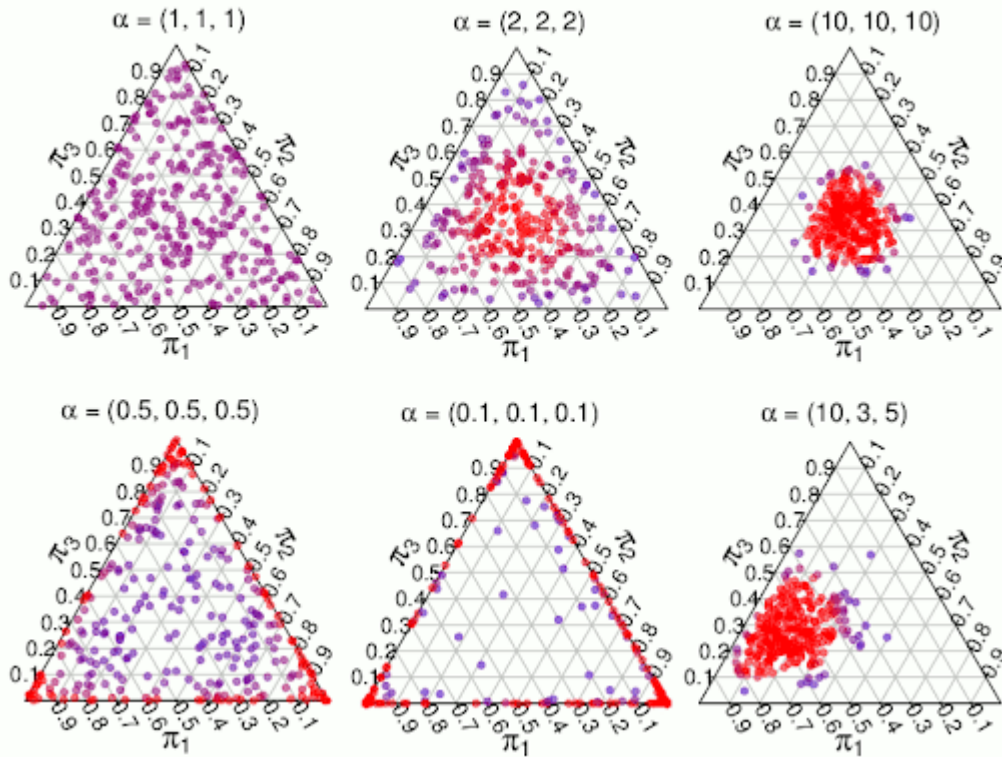
- 1) For \mathbf{d}_i in $C = \{\text{docs } \mathbf{d}_i \mid \mathbf{d}_i = [d_{i1}, d_{i2}, \dots, d_{ik}]\}_{i=1}^M = [\text{Matrix}]_{M \times k}$
- $\mathbf{d}_i = [d_{i1}, d_{i2}, \dots, d_{ik}] \sim \text{Dirichlet}(\alpha = [\alpha_1, \alpha_2, \dots, \alpha_k])$

$$f(x_1, \dots, x_K; \alpha_1, \dots, \alpha_K) = \frac{1}{B(\boldsymbol{\alpha})} \prod_{i=1}^K x_i^{\alpha_i - 1}$$


We've seen this before. where?

Probabilistic Topic Models

Draws from a 3-dimensional Dirichlet with different α

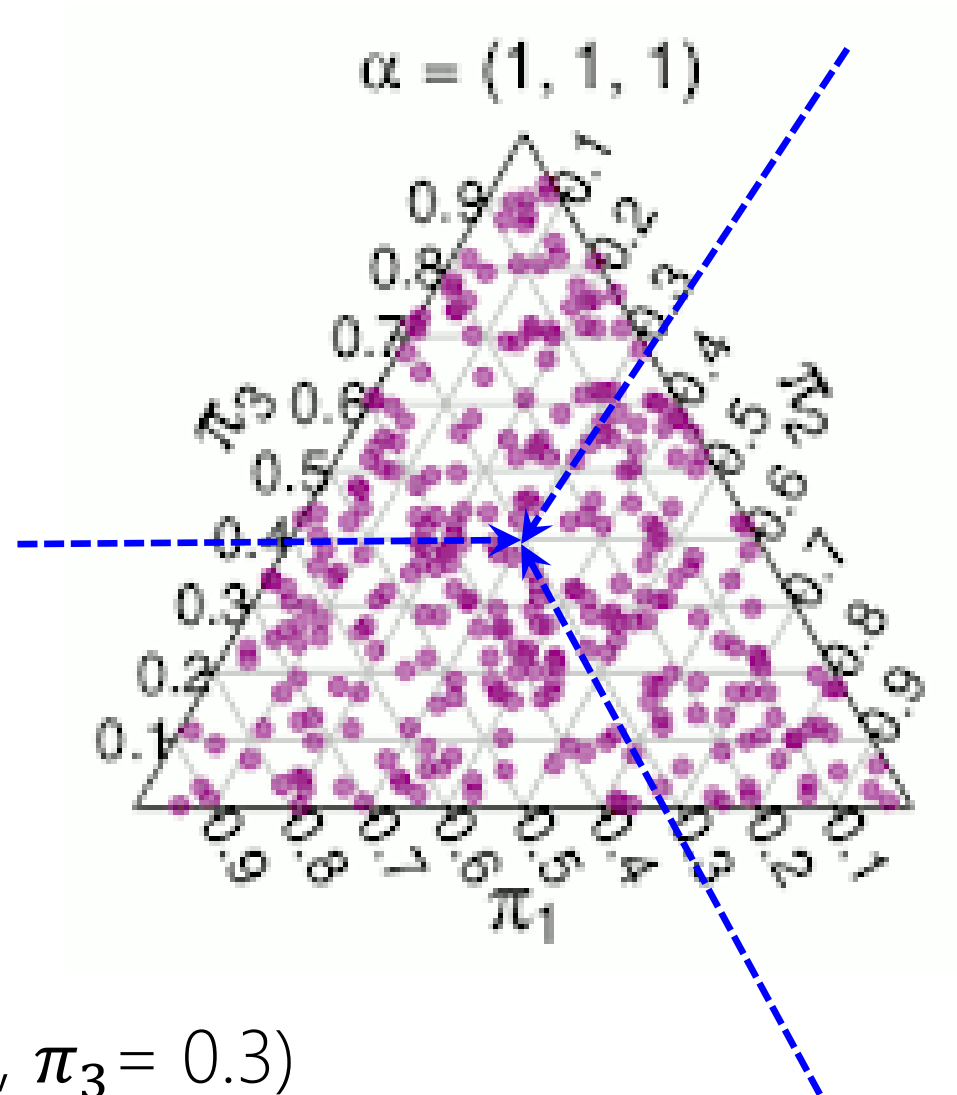


<http://www.sumsar.net/blog/2015/04/the-non-parametric-bootstrap-as-a-bayesian-model/>

https://en.wikipedia.org/wiki/Dirichlet_distribution

Probabilistic Topic Models

Random chance to
 π_1, π_2, π_3



$(\pi_1 = 0.3, \pi_2 = 0.4, \pi_3 = 0.3)$

Simplex

From Wikipedia, the free encyclopedia

For other uses, see [Simplex \(disambiguation\)](#).

In [geometry](#), a **simplex** (plural: **simplexes** or **simplices**) is a generalization of the notion of a [triangle](#) or [tetrahedron](#) to arbitrary [dimensions](#). The simplex is so-named because it represents the simplest possible [polytope](#) in any given space.

For example,

- a 0-simplex is a [point](#),
- a 1-simplex is a [line segment](#),
- a 2-simplex is a [triangle](#),
- a 3-simplex is a [tetrahedron](#),
- a 4-simplex is a [5-cell](#).

Specifically, a ***k*-simplex** is a *k*-dimensional [polytope](#) which is the [convex hull](#) of its *k* + 1 [vertices](#). More formally, suppose the *k* + 1 points $u_0, \dots, u_k \in \mathbb{R}^k$ are [affinely independent](#), which means $u_1 - u_0, \dots, u_k - u_0$ are [linearly independent](#). Then, the simplex determined by them is the set of points

$$C = \left\{ \theta_0 u_0 + \dots + \theta_k u_k \mid \sum_{i=0}^k \theta_i = 1 \text{ and } \theta_i \geq 0 \text{ for } i = 0, \dots, k \right\}.$$

A **regular simplex**^[1] is a simplex that is also a [regular polytope](#). A regular *k*-simplex may be constructed from a regular (*k* − 1)-simplex by connecting a new vertex to all original vertices by the common edge length.

The **standard simplex** or **probability simplex** ^[2] is the simplex whose vertices are the *k* standard unit vectors and the origin, or

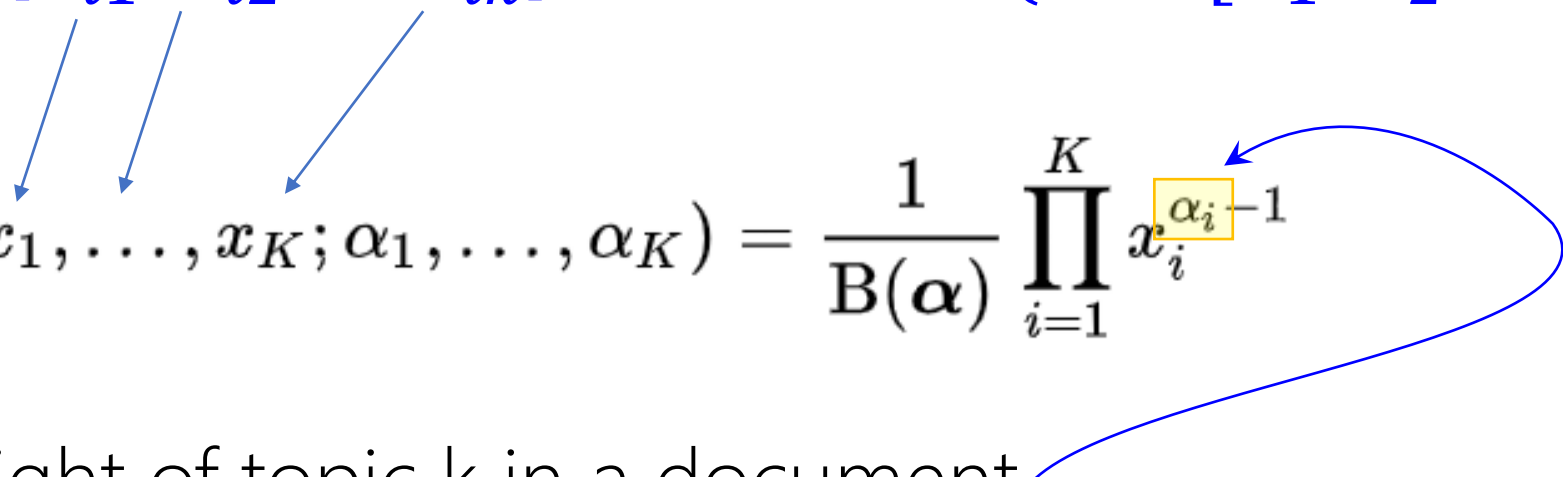
$$\{x \in \mathbb{R}^k : x_0 + \dots + x_{k-1} = 1, x_i \geq 0 \text{ for } i = 0, \dots, k - 1\}.$$



Probabilistic Topic Models

Initializations:

- 1) For d_i in $C = \{\text{docs } d_i \mid d_i = [d_{i1}, d_{i2}, \dots, d_{ik}]\}_{i=1}^M = [\text{Matrix}]_{M \times k}$
- $d_i = [d_{i1}, d_{i2}, \dots, d_{ik}] \sim \text{Dirichlet}(\alpha = [\alpha_1, \alpha_2, \dots, \alpha_k])$


$$f(x_1, \dots, x_K; \alpha_1, \dots, \alpha_K) = \frac{1}{B(\boldsymbol{\alpha})} \prod_{i=1}^K x_i^{\alpha_i - 1}$$

Prior weight of topic k in a document

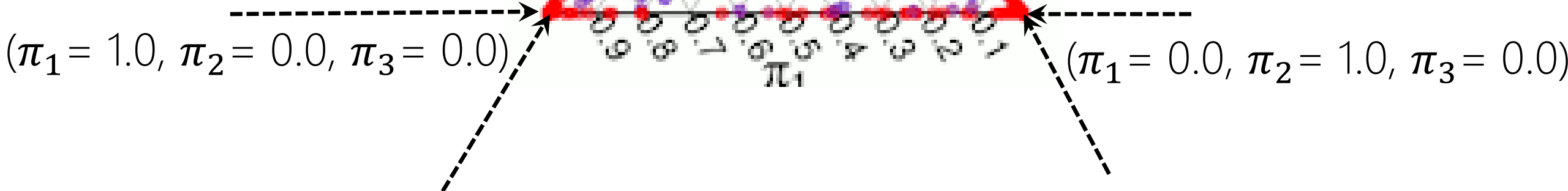
- Initially the same for all topics
- Normally $\alpha_{1 \leq j \leq k} < 1$, e.g. 0.1,
- Prefer sparse topic distributions, i.e., few topics per document

Probabilistic Topic Models

Sparse chance to π_1, π_2, π_3

$$\alpha = (0.1, 0.1, 0.1)$$

$$(\pi_1 = 0.0, \pi_2 = 0.0, \pi_3 = 1.0)$$

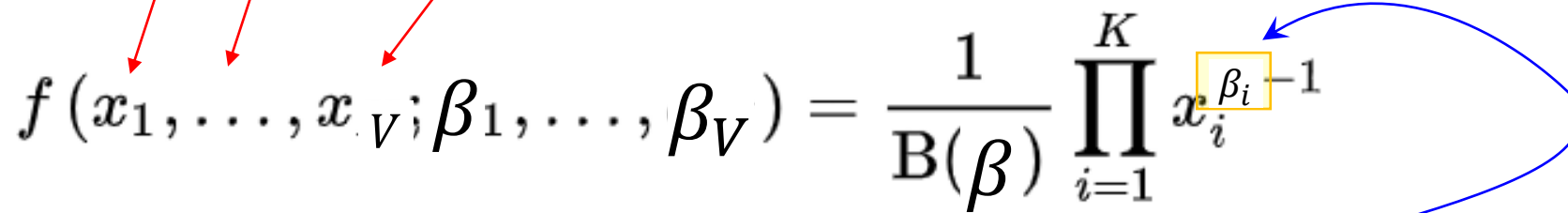


Probabilistic Topic Models

Initializations:

2) For \mathbf{z}_i in $Z = \{\text{topics } \mathbf{z}_i \mid \mathbf{z}_i = [z_{i1}, z_{i2}, \dots, z_{iV}]\}_{i=1}^K = [\text{Matrix}]_{K \times V}$

○ $\mathbf{z}_i = [z_{i1}, z_{i2}, \dots, z_{iV}] \sim \text{Dirichlet}(\beta = [\beta_1, \beta_2, \dots, \beta_V])$


$$f(x_1, \dots, x_V; \beta_1, \dots, \beta_V) = \frac{1}{B(\beta)} \prod_{i=1}^K x_i^{\beta_i - 1}$$

Prior weight of topic k in a document

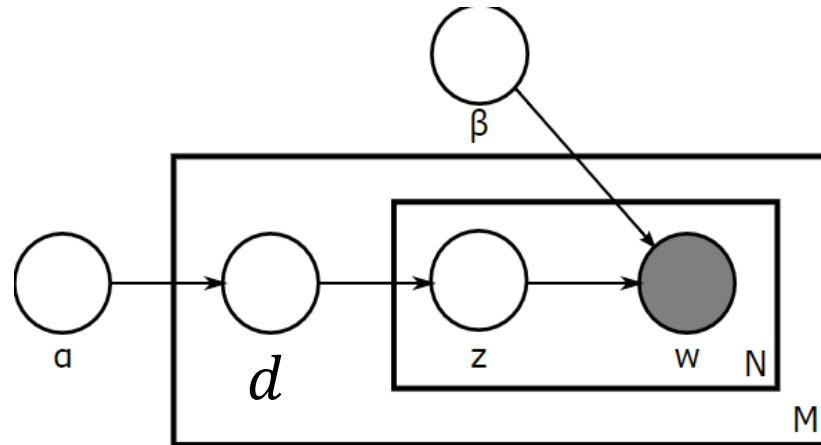
- Initially the same for all topics
- Normally $\beta_{1 \leq j \leq V} \ll 1$, e.g. 0.001,
- Very sparse word distributions, i.e., few words per topics

Probabilistic Topic Models

Generative Process:

$$p(\boldsymbol{\theta}, \mathbf{z}, \mathbf{w} \mid \boldsymbol{\alpha}, \boldsymbol{\beta}) = p(\boldsymbol{\theta} \mid \boldsymbol{\alpha}) \prod_{n=1}^N p(z_n \mid \boldsymbol{\theta}) p(w_n \mid z_n, \boldsymbol{\beta}),$$

Plate Notation
Graphical Model

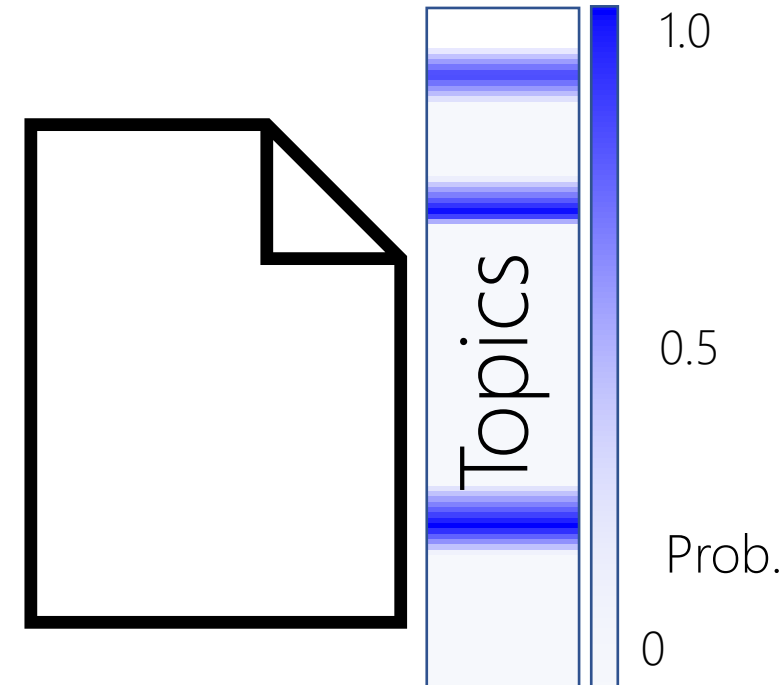
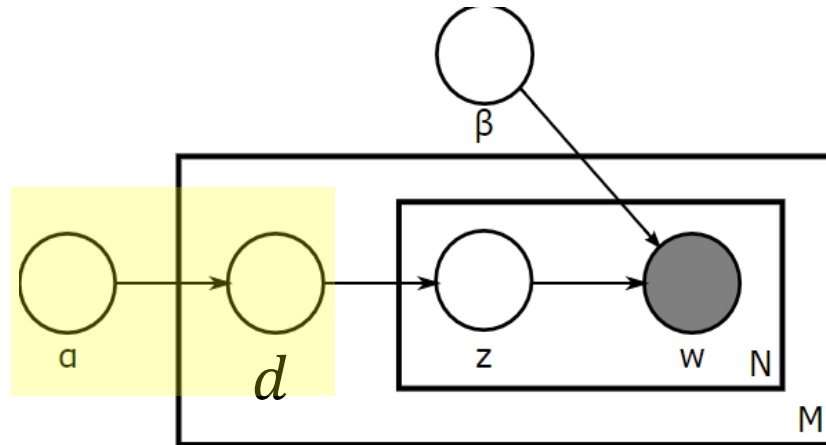


Probabilistic Topic Models

Generative Process:

$$p(\boldsymbol{\theta}, \mathbf{z}, \mathbf{w} \mid \boldsymbol{\alpha}, \boldsymbol{\beta}) = p(\boldsymbol{\theta} \mid \boldsymbol{\alpha}) \prod_{n=1}^N p(z_n \mid \boldsymbol{\theta}) p(w_n \mid z_n, \boldsymbol{\beta}),$$

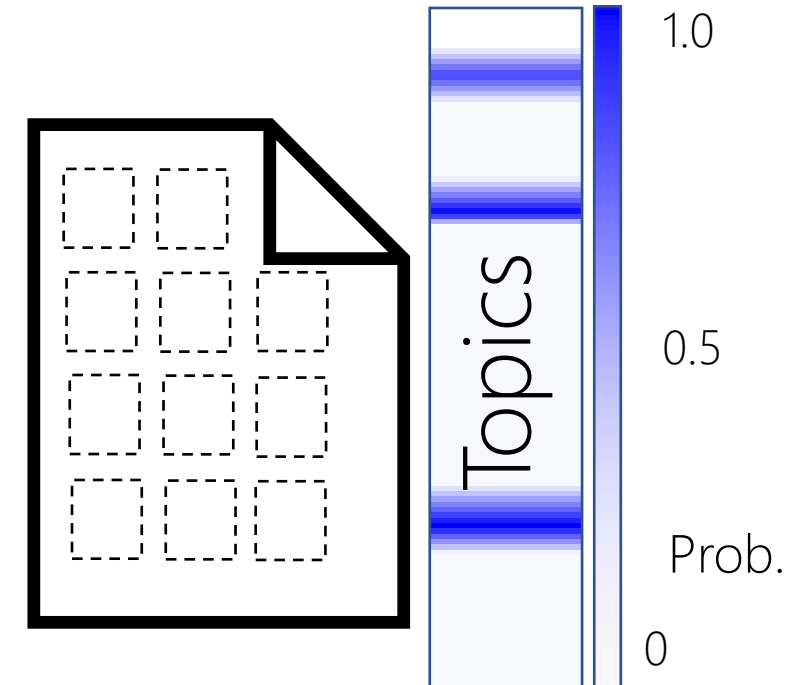
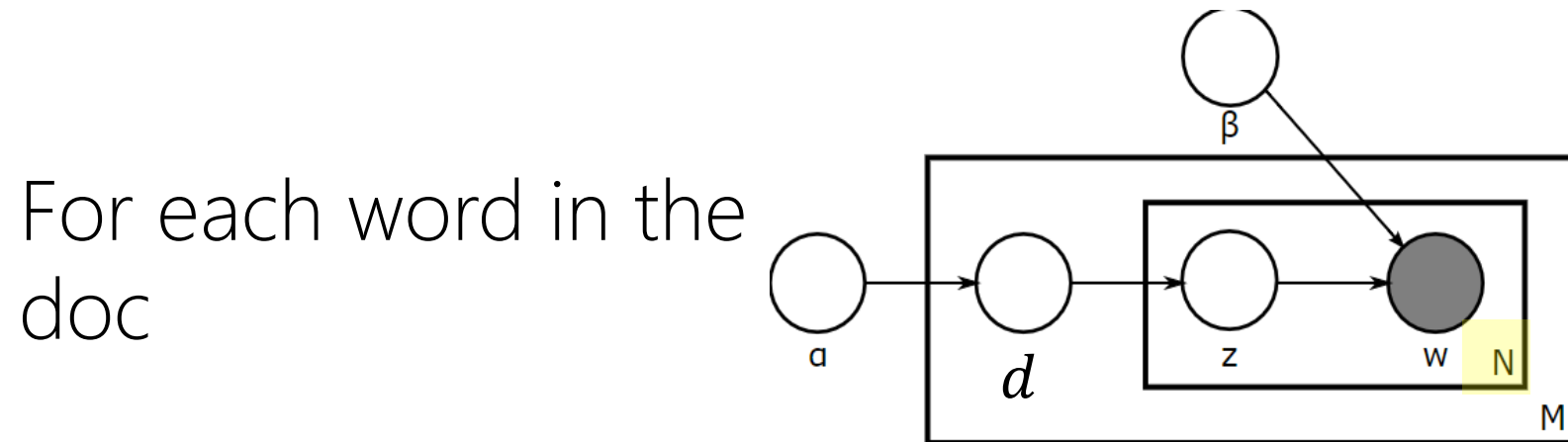
Select the topic
distribution for the
doc



Probabilistic Topic Models

Generative Process:

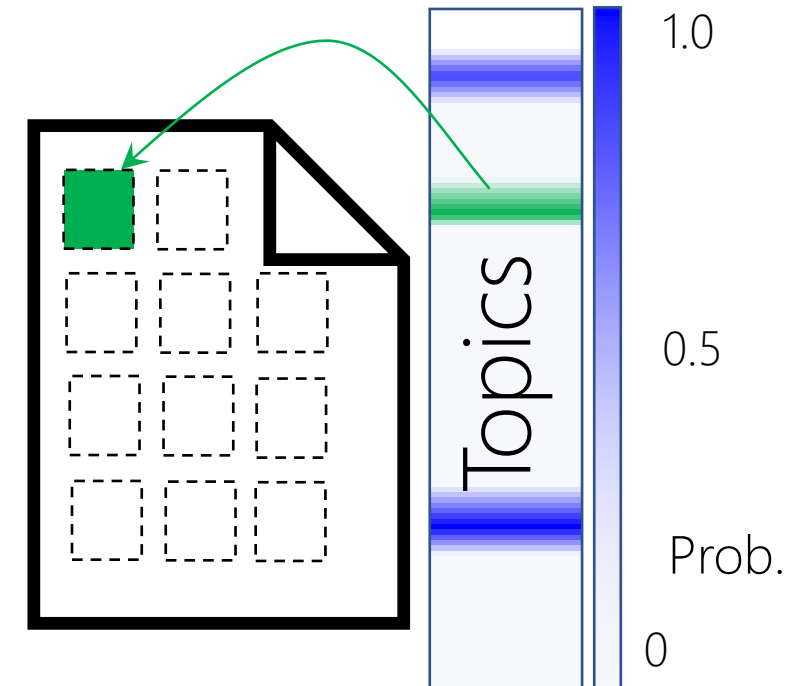
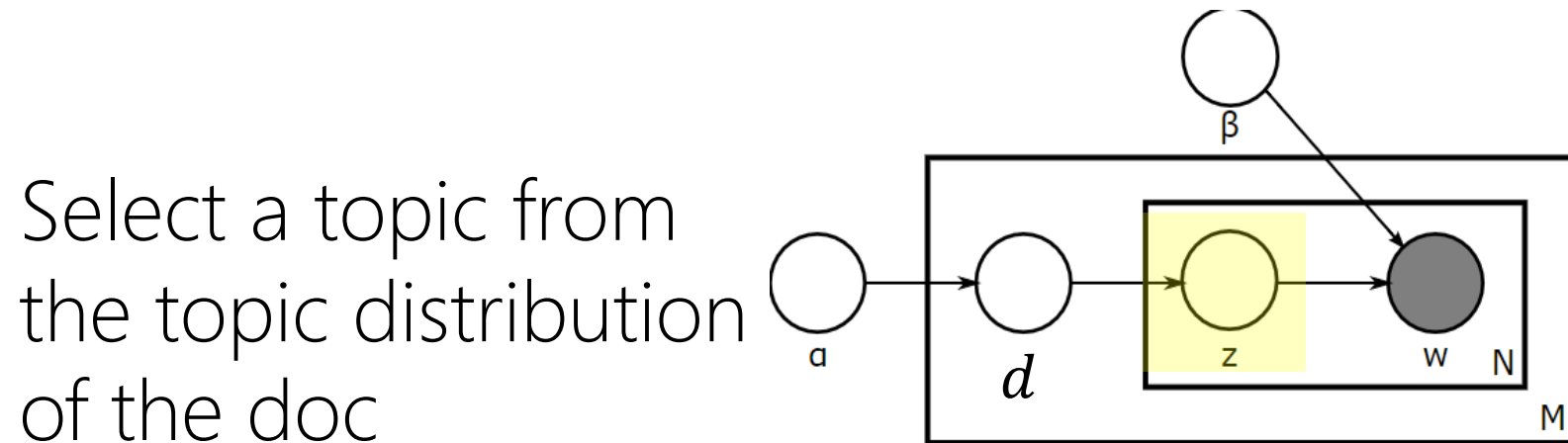
$$p(\boldsymbol{\theta}, \mathbf{z}, \mathbf{w} \mid \boldsymbol{\alpha}, \boldsymbol{\beta}) = p(\boldsymbol{\theta} \mid \boldsymbol{\alpha}) \prod_{n=1}^N p(z_n \mid \boldsymbol{\theta}) p(w_n \mid z_n, \boldsymbol{\beta}),$$



Probabilistic Topic Models

Generative Process:

$$p(\boldsymbol{\theta}, \mathbf{z}, \mathbf{w} \mid \boldsymbol{\alpha}, \boldsymbol{\beta}) = p(\boldsymbol{\theta} \mid \boldsymbol{\alpha}) \prod_{n=1}^N p(z_n \mid \boldsymbol{\theta}) p(w_n \mid z_n, \boldsymbol{\beta}),$$

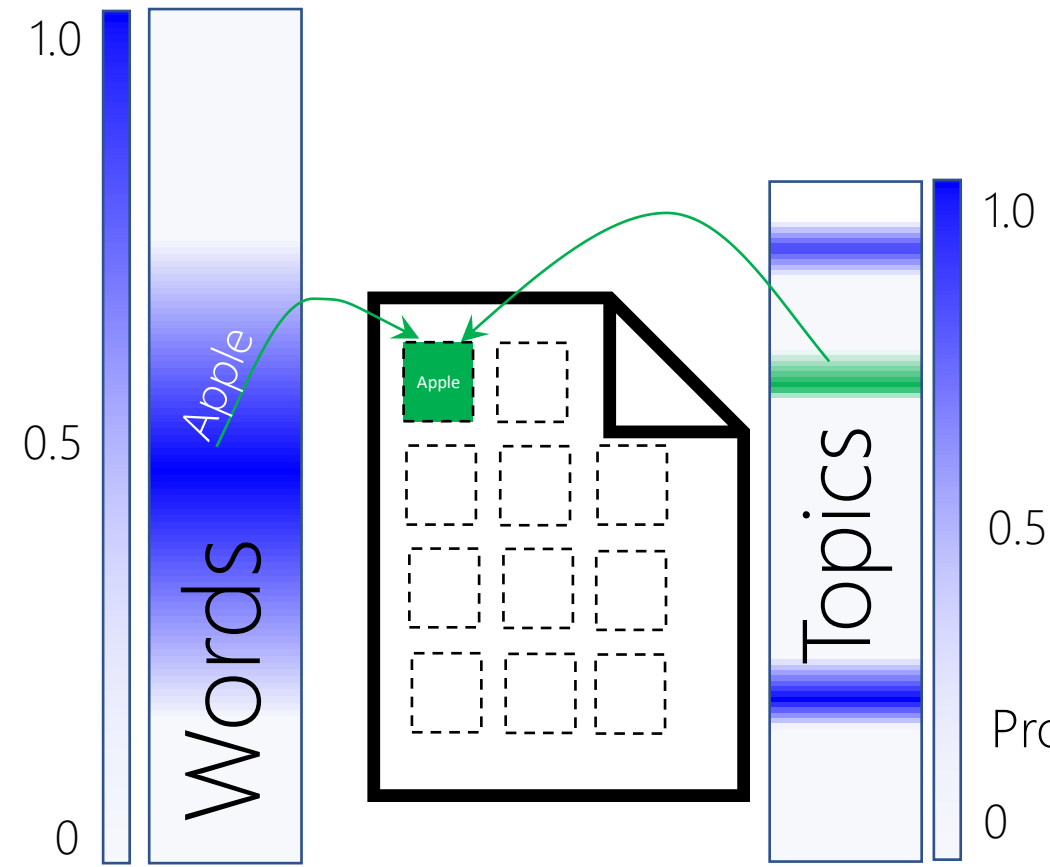
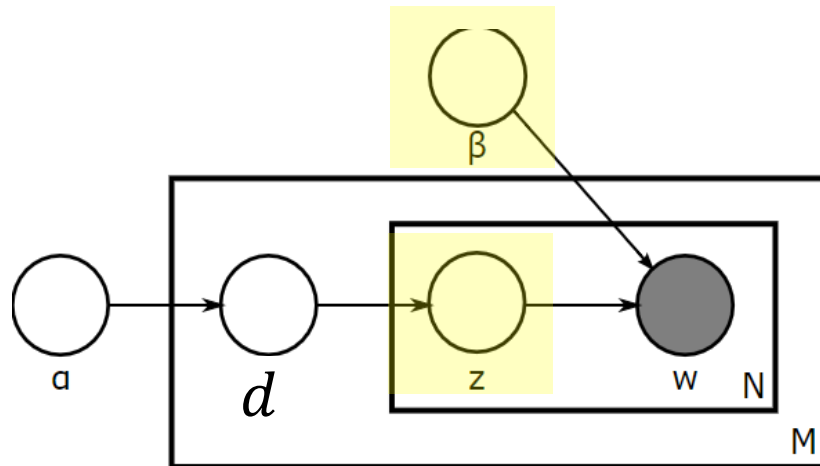


Probabilistic Topic Models

Generative Process:

$$p(\boldsymbol{\theta}, \mathbf{z}, \mathbf{w} \mid \boldsymbol{\alpha}, \boldsymbol{\beta}) = p(\boldsymbol{\theta} \mid \boldsymbol{\alpha}) \prod_{n=1}^N p(z_n \mid \boldsymbol{\theta}) p(w_n \mid z_n, \boldsymbol{\beta}),$$

Select the word from
the topic distribution



Probabilistic Topic Models

Generative Process:

For d_i in C

For $j = 0: N$

$z =$ Choose a topic z_i in Z based on $[d_{i1}, d_{i2}, \dots, d_{ik}]$

$w =$ Choose a word w_j in V based on $z = [z_1, z_2, \dots, z_V]$

Probabilistic Topic Models

Generative Process:

For d_i in C

For $j = 0: N$

$z =$ Choose a topic z_i in Z based on $[d_{i1}, d_{i2}, \dots, d_{ik}]$

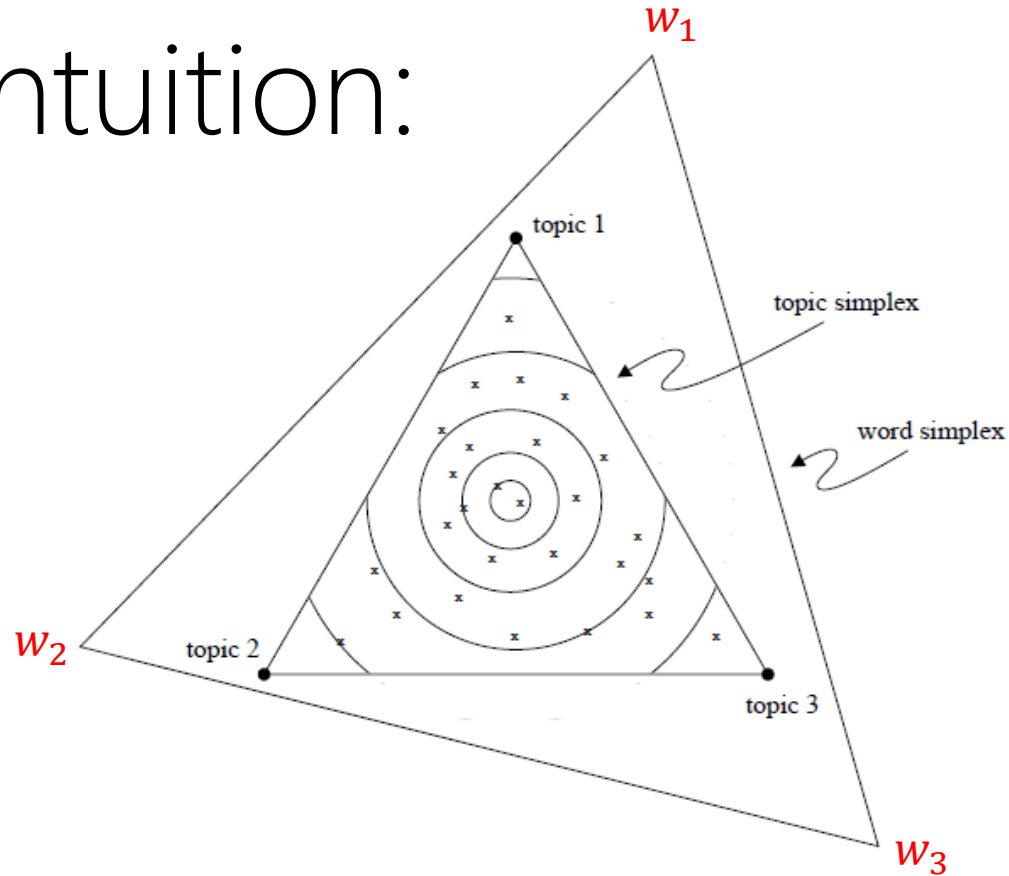
$w =$ Choose a word w_j in V based on $z = [z_1, z_2, \dots, z_V]$

$$p(\theta, \mathbf{z}, \mathbf{w} \mid \alpha, \beta) = p(\theta \mid \alpha) \prod_{n=1}^N p(z_n \mid \theta) p(w_n \mid z_n, \beta),$$

- Bernoulli distribution:
1 out of two mutually exclusive options = success (p) or failure ($q = 1 - p$)
- Generalized Bernoulli (Categorical) distribution = 1 out of N mutually exclusive options p_i

Probabilistic Topic Models

Geometric Intuition:



For d_i in C

For $j = 0: N$

$z =$ Choose a topic z_i in Z based on $[d_{i1}, d_{i2}, \dots, d_{ik}]$

$w =$ Choose a word w_j in V based on $z = [z_1, z_2, \dots, z_V]$

Probabilistic Topic Models

Optimization

Generated Docs \leftrightarrow Observed Docs

$$\ell(\alpha, \beta) = \sum_{d=1}^M \log p(\mathbf{w}_d | \alpha, \beta).$$

For d_i in C

For $j = 0: N$

\mathbf{z} = Choose a topic \mathbf{z}_i in Z based on $d_i = [d_{i1}, d_{i2}, \dots, d_{ik}]$

\mathbf{w} = Choose a word \mathbf{w}_j in V based on $\mathbf{z} = [z_1, z_2, \dots, z_V]$

Probabilistic Topic Models

Question: what's the difference?

For d_i in C

$z =$ Choose a topic z_i in Z based on $[d_{i1}, d_{i2}, \dots, d_{ik}]$

For $j = 0: N$

$w =$ Choose a word w_j in V based on $z = [z_1, z_2, \dots, z_V]$

Probabilistic Topic Models

Question: what's the difference?

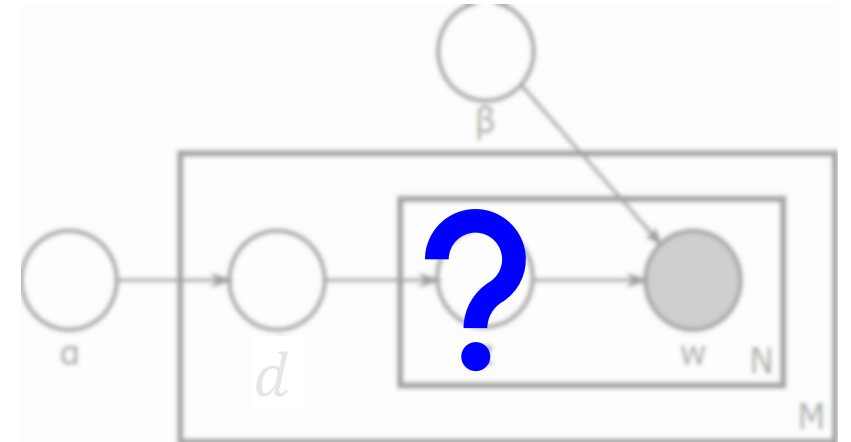
For d_i in C

$z =$ Choose a topic z_i in Z based on $[d_{i1}, d_{i2}, \dots, d_{ik}]$

For $j = 0: N$

$w =$ Choose a word w_j in V based on $z = [z_1, z_2, \dots, z_V]$

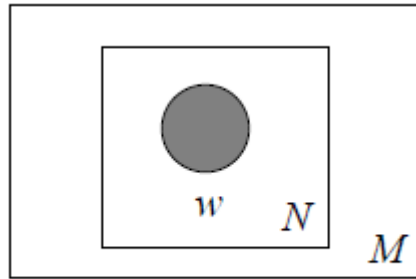
$$p(\theta, \mathbf{z}, \mathbf{w} \mid \alpha, \beta) = p(\theta \mid \alpha) p(z_n \mid \theta) \prod_{n=1}^N p(w_n \mid z_n, \beta)$$



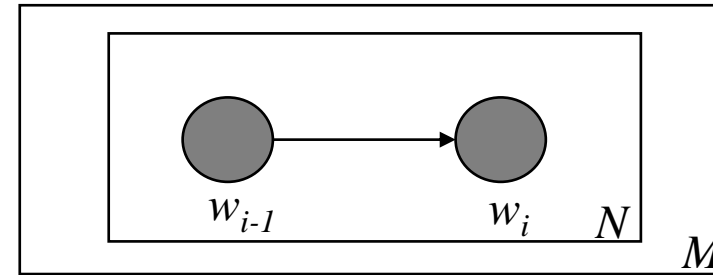
Probabilistic Topic Models

LDA vs. or as a n-gram LM

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



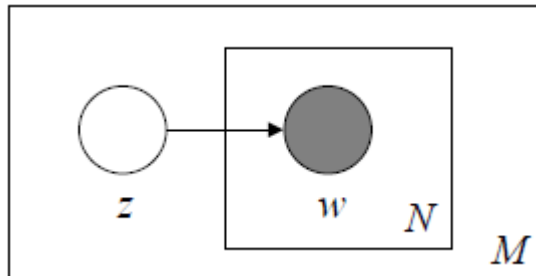
(a) unigram



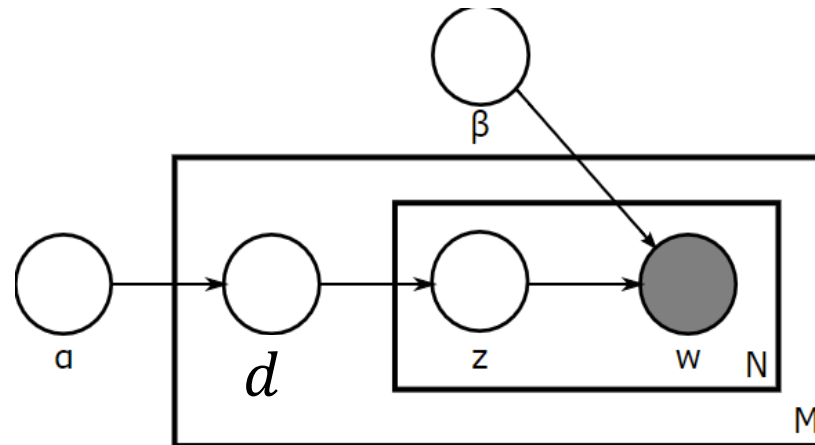
(b) bigram

$$p(\mathbf{w}) = ?$$

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



(c) mixture of unigrams



Probabilistic Topic Models: API

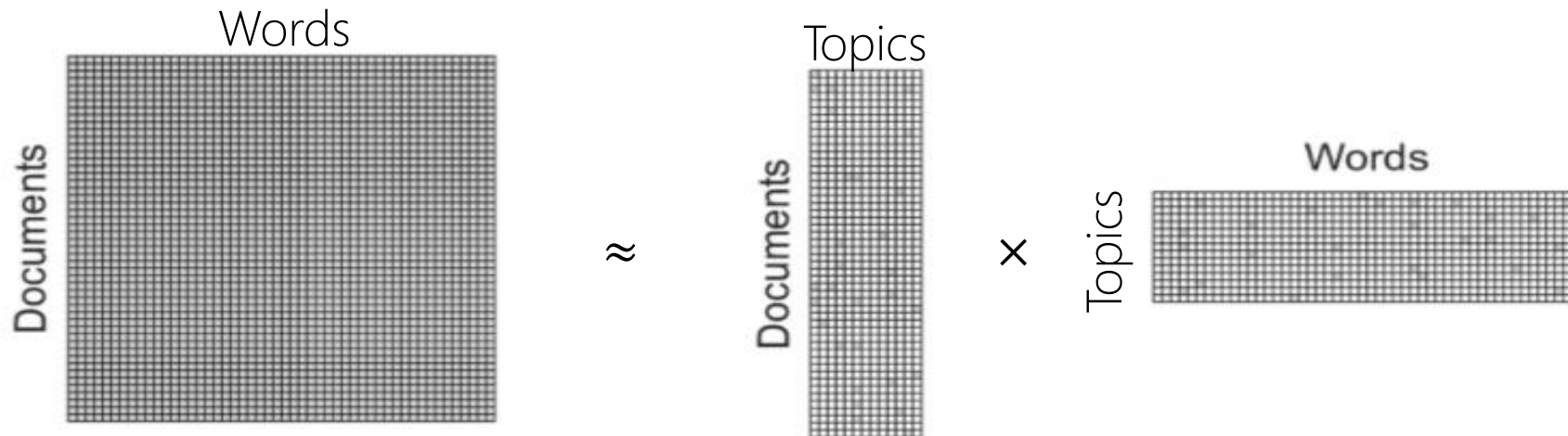
- MALLET (<http://mallet.cs.umass.edu/topics.php>)
- Gensim (<https://radimrehurek.com/gensim/>)
- Gensim wrapper for MALLET

Probabilistic Topic Models: Applications

- Corpus Summarization
- Document Clustering / Classification
- Document Summarization
- User Modeling

Probabilistic Topic Models: Applications

- Corpus Summarization
 - Corpus \rightarrow Topics \rightarrow Top-k Words



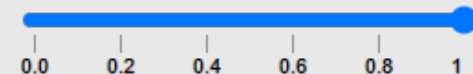
Probabilistic Topic Models: Applications

1	Week	NASA	Space	Science	Scientist	.com	Phys.org	Planet	Community	Time					
	0.085	0.054	0.051	0.048	0.032	0.031	0.026	0.02	0.019	0.018					
2	Message	Email	Day	England	Time	France	Manchester	Assassination	Venezuela	Game				Configuration:	
	0.041	0.038	0.037	0.035	0.026	0.024	0.024	0.024	0.02	0.02				# of Topics	40
3	Quebec	Pará	Latin	Delaware	System	Mexico	Sweden	Travel	Como	Louisiana				Dataset	1/14 zabel
	0.079	0.043	0.036	0.027	0.026	0.026	0.021	0.021	0.021	0.017				Representation	USER
4	DONT	Ur	Fuck	International	Music	RT!	God	Human	Airport	Feces				Topic Detector	MALLET
	0.039	0.033	0.026	0.019	0.018	0.016	0.015	0.015	0.014	0.014				Preprocessing	TagME
5	Hootsuite	Party	LI	Sea	Light-year	.co	Life	People	Shit	Human				# of Users	1619
	0.077	0.067	0.043	0.043	0.036	0.024	0.022	0.017	0.015	0.015				# of Tweets	167572
6	Street	China	Business	Reuters	Stock	Ireland	Bank	Reut	Economics	0				Filter	Yes
	0.127	0.042	0.038	0.031	0.027	0.027	0.016	0.016	0.015	0.014				TagME threshold	0.01
7	Mail	Sunday	Time	People	Love	Person	Hope	Thought	Nice	Christmas					
	0.057	0.05	0.048	0.042	0.029	0.026	0.023	0.023	0.018	0.014					
8	MSNBC	News	White	Committee	World	Fox	Hootsuite	President	Party	PBS					
	0.079	0.055	0.036	0.033	0.029	0.029	0.026	0.024	0.023	0.019					
9	BBC	News	World	Sky	Ireland	Death	TGR	Murder	.uk	PH					
	0.317	0.261	0.172	0.026	0.006	0.006	0.006	0.006	0.005	0.005					
10	Foursquare	Wales	Facebook	Engagement	Time	Wedding	International	Family	Shanghai	CNNMoney					
	0.077	0.043	0.036	0.035	0.033	0.031	0.021	0.019	0.017	0.016					

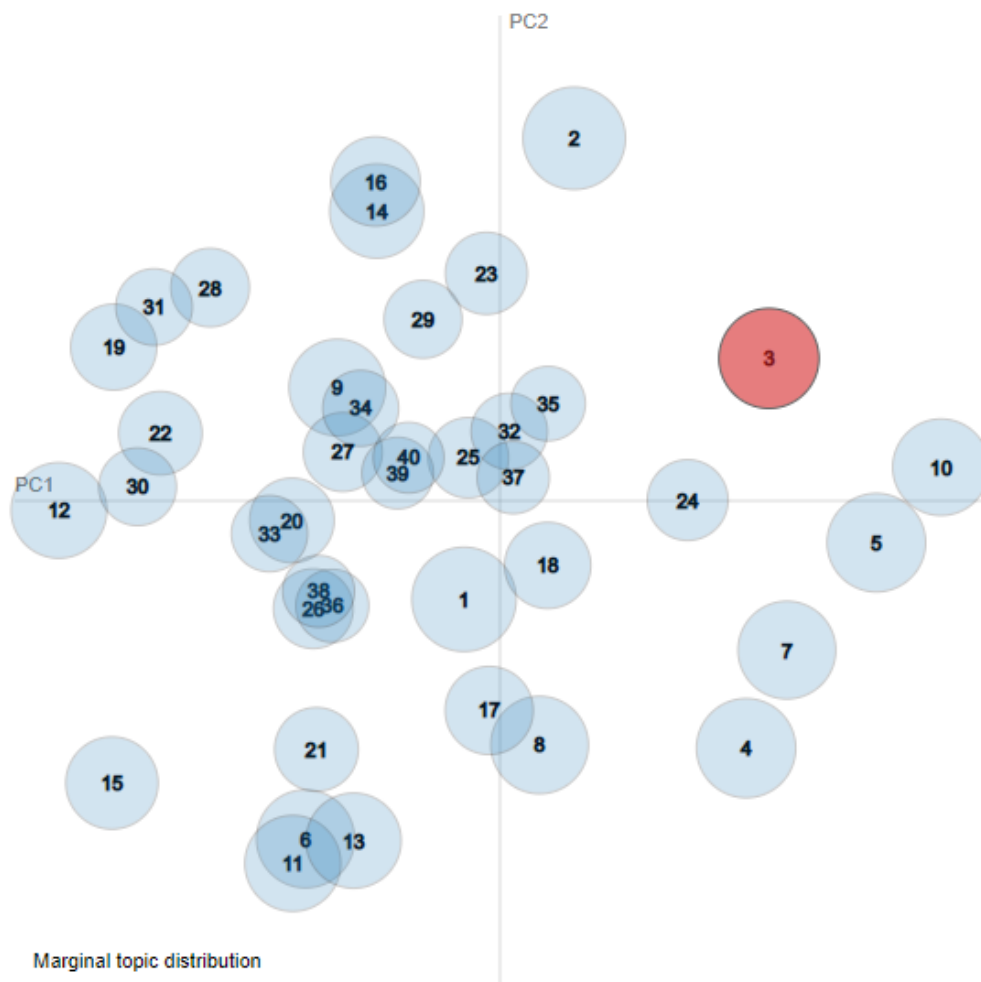
In [207]: code_model['vis']

Out[207]: Selected Topic: Previous Topic Next Topic Clear Topic

Slide to adjust relevance metric:⁽²⁾
 $\lambda = 1$



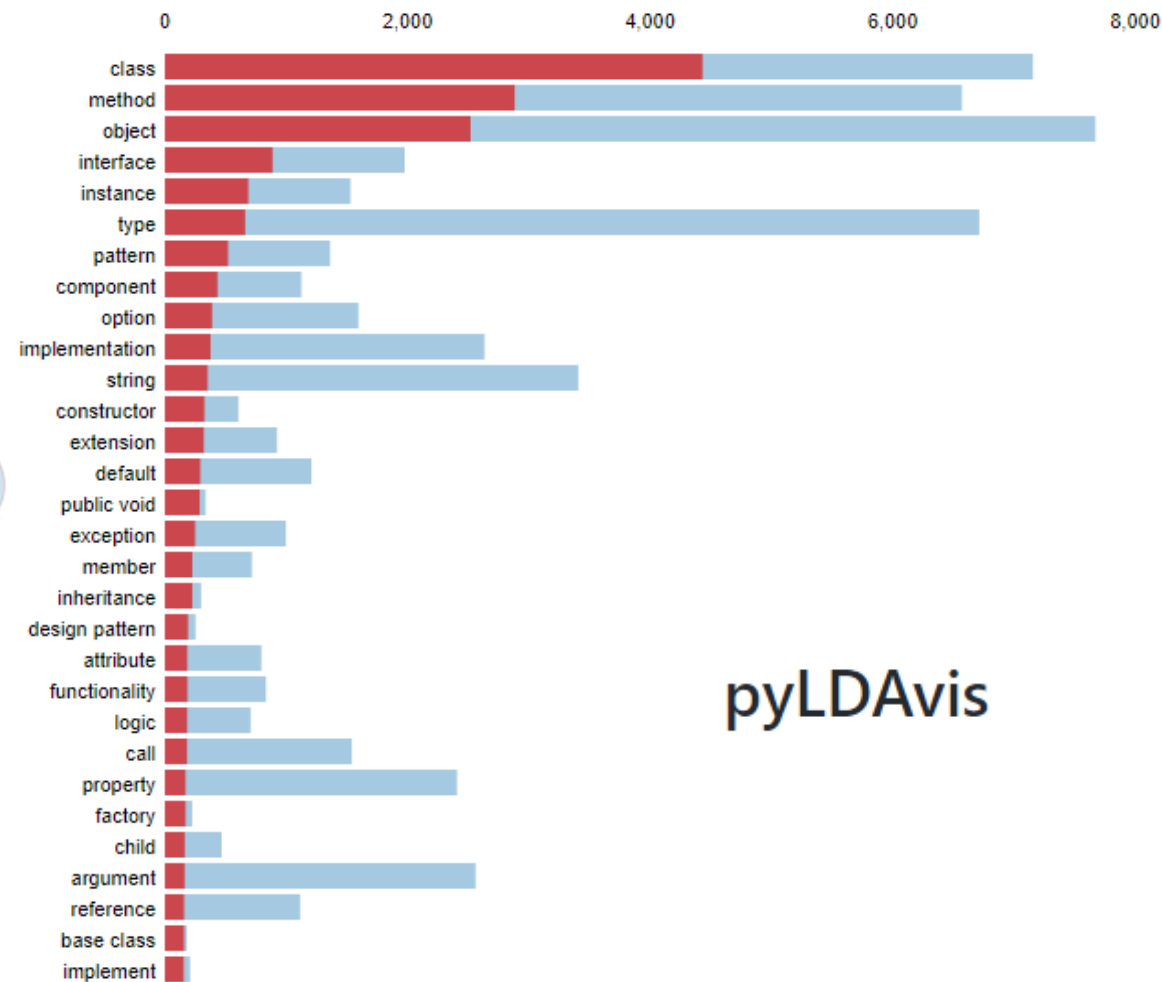
Intertopic Distance Map (via multidimensional scaling)



Marginal topic distribution



Top-30 Most Relevant Terms for Topic 3 (3.3% of tokens)



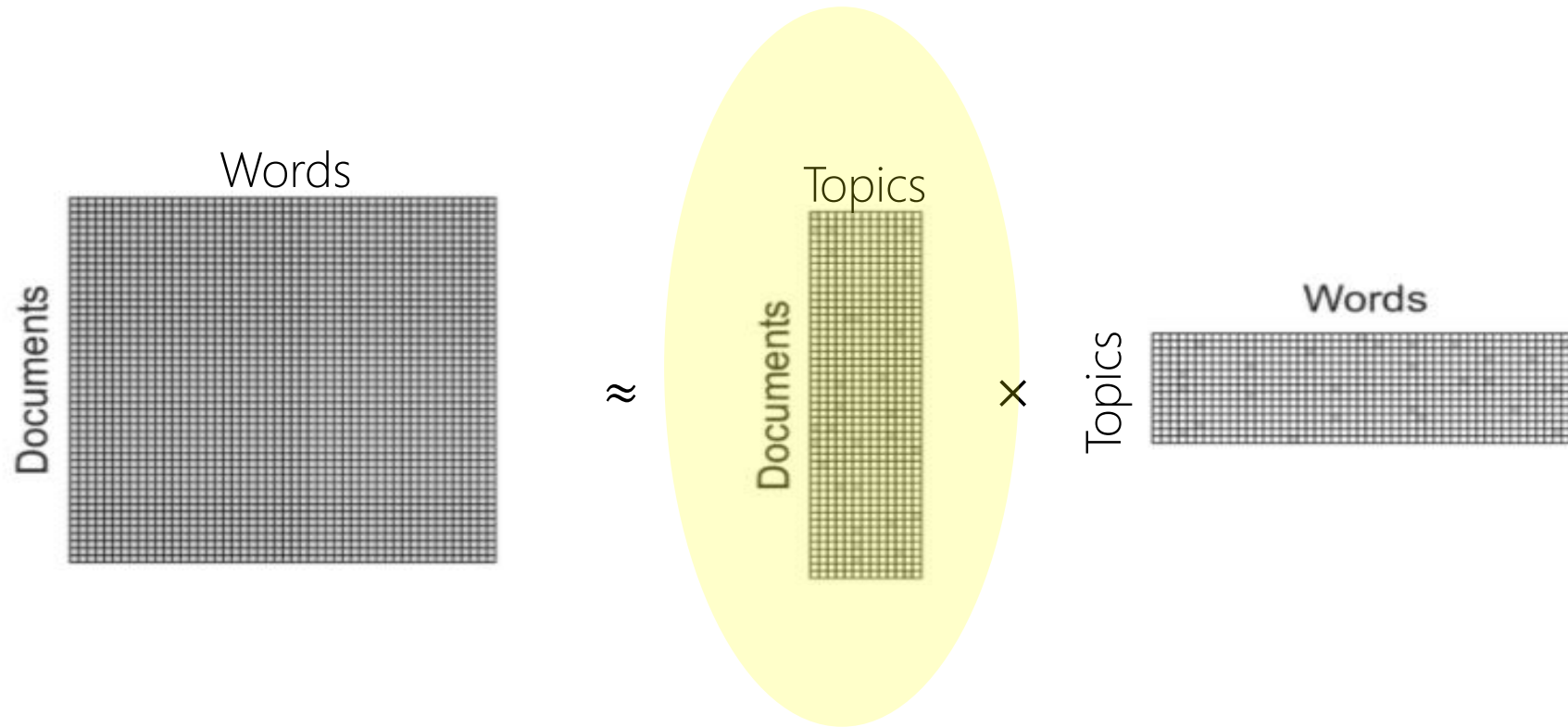
pyLDAvis

Overall term frequency
Estimated term frequency within the selected topic

1 $\text{saliency}(\text{term } w) = \text{frequency}(w) * (\sum_t \text{f}(t|w) * \log(\text{f}(t|w)/\text{f}(t)))$ for topics t ; see Chuang et al (2012)

Probabilistic Topic Models: Applications

- Document Clustering / Classification



Probabilistic Topic Models: Applications

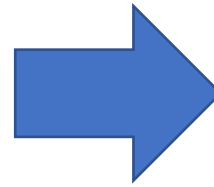
- Document Clustering / Classification
 - Each document is a feature vector of topics
 - Similar documents have similar vectors/topics

[illegible]

Probabilistic Topic Models: Applications

- Document Summary

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



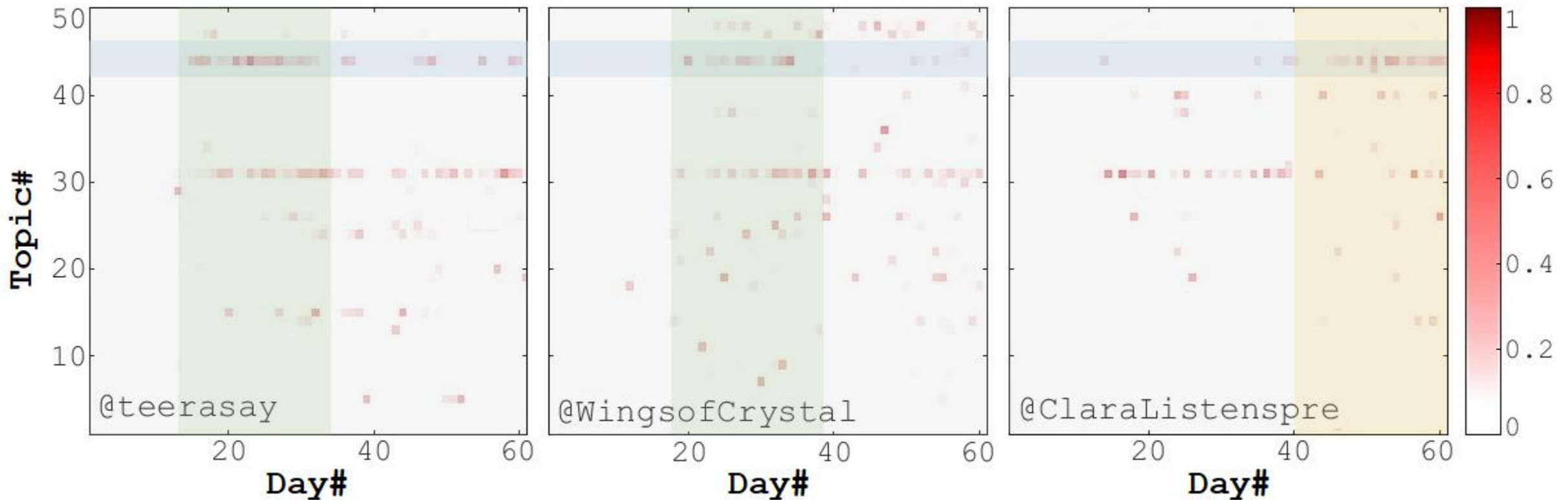
NEW	MILLION	CHILDREN	SCHOOL
FILM	TAX	WOMEN	STUDENTS
SHOW	PROGRAM	PEOPLE	SCHOOLS
MUSIC	BUDGET	CHILD	EDUCATION
MOVIE	BILLION	YEARS	TEACHERS
PLAY	FEDERAL	FAMILIES	HIGH
ACTOR	NEW	WORK	PUBLIC
FIRST	STATE	PARENTS	TEACHER
YORK	PLAN	SAYS	BENNETT
OPERA	MONEY	FAMILY	MANIGAT
THEATER	PROGRAMS	WELFARE	NAMPHY
ACTRESS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

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.

Probabilistic Topic Models: Applications

- User Modeling (will be discussed more ...)



Distributed Memory Model of Paragraph Vectors (PV-DM) aka Doc2Vec
