

# Latent Dirichlet Allocation Overview

Kanru Wang

March 2020

# Usage and Main Idea

- Latent Dirichlet Allocation (LDA) is the one of the most popular algorithms for document clustering.
- E.g.
  - News article clustering, which facilitates recommendation
  - Large corpus topic discovery

# Usage and Main Idea

- Latent Dirichlet Allocation (LDA) is the one of the most popular algorithms for document clustering.
- E.g.
  - News article clustering, which facilitates recommendation
  - Large corpus topic discovery
- Main idea:
  - Assume that each document, if treated as bag-of-words, can be generated by a sampling process

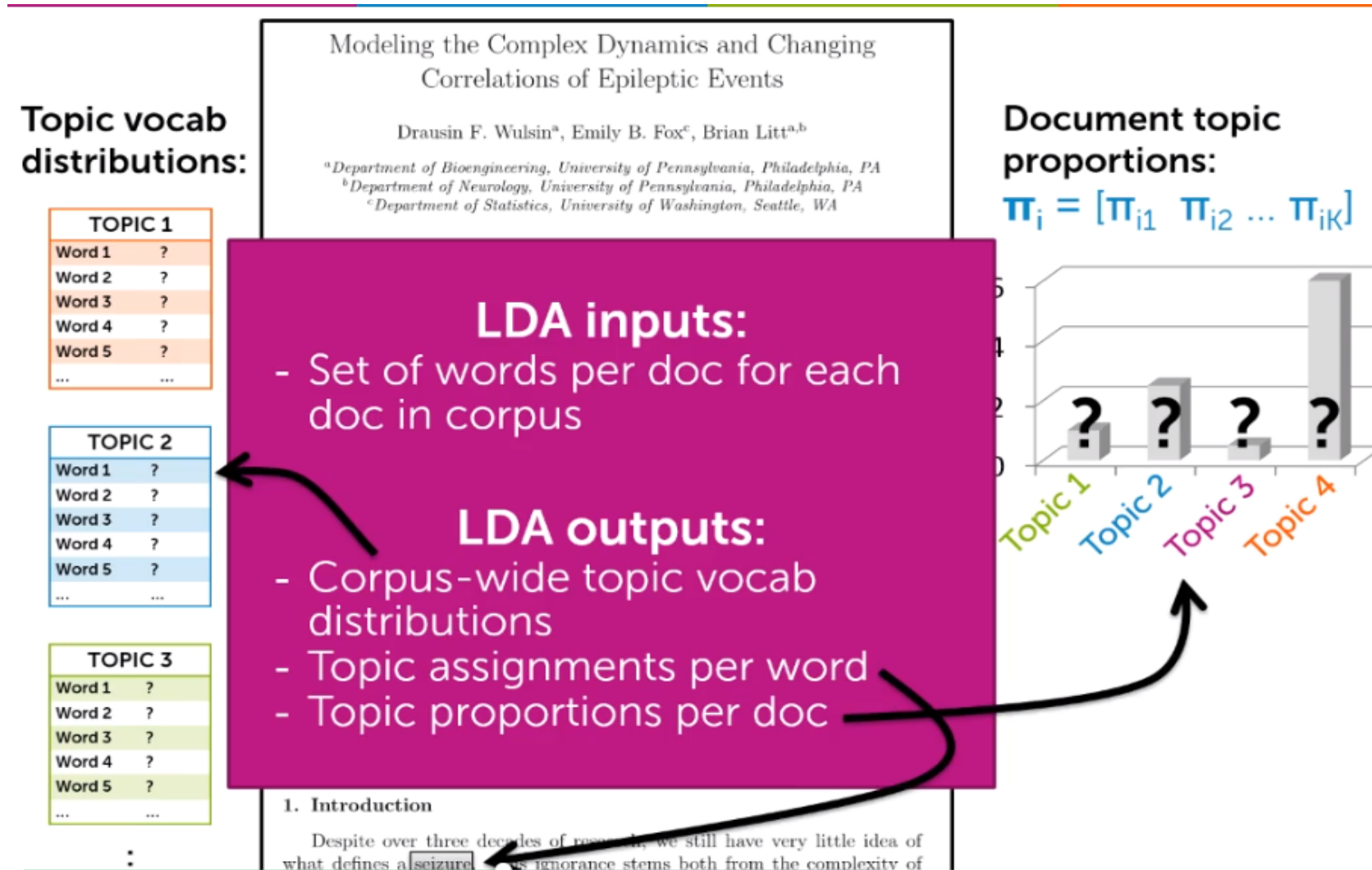
# Usage and Main Idea

- Latent Dirichlet Allocation (LDA) is the one of the most popular algorithms for document clustering.
- E.g.
  - News article clustering, which facilitates recommendation
  - Large corpus topic discovery
- Main idea:
  - Assume that each document, if treated as bag-of-words, can be generated by a sampling process
  - Each topic is just a specific distribution of words.  
(e.g. a Machine Learning topic may have: “model” 132, “rate” 118, ..., “cat” 3, “song” 1, “football” 0)
  - Each document has a specific distribution of topics.

# Usage and Main Idea

- Latent Dirichlet Allocation (LDA) is the one of the most popular algorithms for document clustering.
- E.g.
  - News article clustering, which facilitates recommendation
  - Large corpus topic discovery
- Main idea:
  - Assume that each document, if treated as bag-of-words, can be generated by a sampling process
  - Each topic is just a specific distribution of words.  
(e.g. a Machine Learning topic may have: “model” 132, “rate” 118, ..., “cat” 3, “song” 1, “football” 0)
  - Each document has a specific distribution of topics.
  - Each document is generated this way:
    - For each word-to-fill in the document (until the whole document is generated):
      1. Based on the document’s distribution of topics, select a topic at random
      2. Based on the chosen topic’s distribution of words, select a word at random
  - Document -> (hidden) topics -> words

# Three useful outputs and how to get them



# Three useful outputs and how to get them

In order to generate a document that has the same bag-of-words representation as the actual document's, we need to optimize the weights in  
Topic Vocab Distributions  
and  
Document Topic Proportions

## Topic vocab distributions:

TOPIC 1	
Word 1	?
Word 2	?
Word 3	?
Word 4	?
Word 5	?
...	...

TOPIC 2	
Word 1	?
Word 2	?
Word 3	?
Word 4	?
Word 5	?
...	...

TOPIC 3	
Word 1	?
Word 2	?
Word 3	?
Word 4	?
Word 5	?
...	...

## Modeling the Complex Dynamics and Changing Correlations of Epileptic Events

Drausin F. Wulsin<sup>a</sup>, Emily B. Fox<sup>c</sup>, Brian Litt<sup>a,b</sup>

<sup>a</sup>Department of Bioengineering, University of Pennsylvania, Philadelphia, PA

<sup>b</sup>Department of Neurology, University of Pennsylvania, Philadelphia, PA

<sup>c</sup>Department of Statistics, University of Washington, Seattle, WA

## LDA inputs:

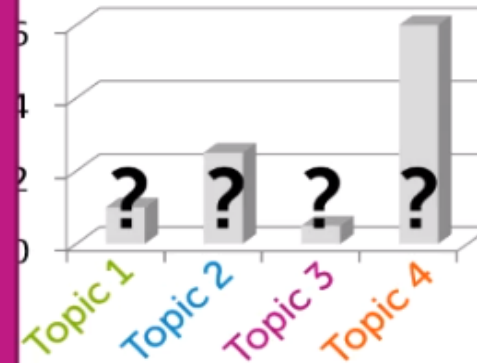
- Set of words per doc for each doc in corpus

## LDA outputs:

- Corpus-wide topic vocab distributions
- Topic assignments per word
- Topic proportions per doc

## Document topic proportions:

$$\pi_i = [\pi_{i1} \ \pi_{i2} \ \dots \ \pi_{iK}]$$



## 1. Introduction

Despite over three decades of research, we still have very little idea of what defines a seizure. This ignorance stems both from the complexity of

# Three useful outputs and how to get them

In order to generate a document that has the same bag-of-words representation as the actual document's, we need to optimize the weights in

Topic Vocab Distributions and Document Topic Proportions

There are 2 ways to optimize them:

1. Collapsed Gibbs Sampling
2. Variational Bayesian Inference (won't discuss)

**Topic vocab distributions:**

TOPIC 1	
Word 1	?
Word 2	?
Word 3	?
Word 4	?
Word 5	?
...	...

TOPIC 2	
Word 1	?
Word 2	?
Word 3	?
Word 4	?
Word 5	?
...	...

TOPIC 3	
Word 1	?
Word 2	?
Word 3	?
Word 4	?
Word 5	?
...	...

Modeling the Complex Dynamics and Changing Correlations of Epileptic Events

Drausin F. Wulsin<sup>a</sup>, Emily B. Fox<sup>c</sup>, Brian Litt<sup>a,b</sup>

<sup>a</sup>Department of Bioengineering, University of Pennsylvania, Philadelphia, PA

<sup>b</sup>Department of Neurology, University of Pennsylvania, Philadelphia, PA

<sup>c</sup>Department of Statistics, University of Washington, Seattle, WA

**LDA inputs:**

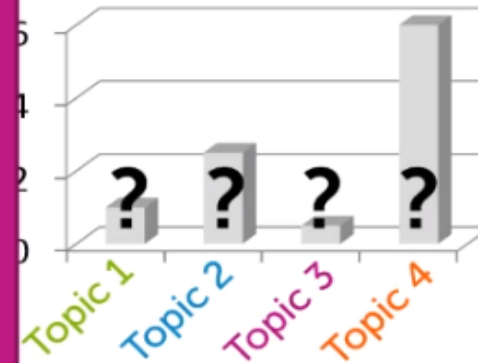
- Set of words per doc for each doc in corpus

**LDA outputs:**

- Corpus-wide topic vocab distributions
- Topic assignments per word
- Topic proportions per doc

**Document topic proportions:**

$$\pi_i = [\pi_{i1} \ \pi_{i2} \ \dots \ \pi_{iK}]$$



## 1. Introduction

Despite over three decades of research, we still have very little idea of what defines a seizure. This ignorance stems both from the complexity of



# Collapsed Gibbs Sampling

## Maintain global statistics

3	2	1	3	1
epilepsy	dynamic	Bayesian	EEG	model

	Topic 1	Topic 2	Topic 3
epilepsy	1	0	35
Bayesian	50	0	1
model	42	1	0
EEG	0	0	20
dynamic	10	8	1
...			

	Topic 1	Topic 2	Topic 3
Doc i	2	1	2

Total  
counts  
from **all**  
docs

- For example we have a document that has just five words, “Epilepsy dynamic Bayesian EEG model”.
- The five numbers in the top table are topic indicators.
- At the start of the entire algorithm, the topic assignment of each word is random.

# Collapsed Gibbs Sampling

## Randomly reassign topics

- Need to reassign every word, for each document, for many iterations
- Let's start from reassigning the word "dynamic".

3	<del>2</del>	1	3	1
epilepsy	dynamic	Bayesian	EEG	model

	Topic 1	Topic 2	Topic 3
epilepsy	1	0	35
Bayesian	50	0	1
model	42	1	0
EEG	0	0	20
dynamic	10	<del>7</del>	1
...			

	Topic 1	Topic 2	Topic 3
Doc i	2	<del>0</del> 1	2

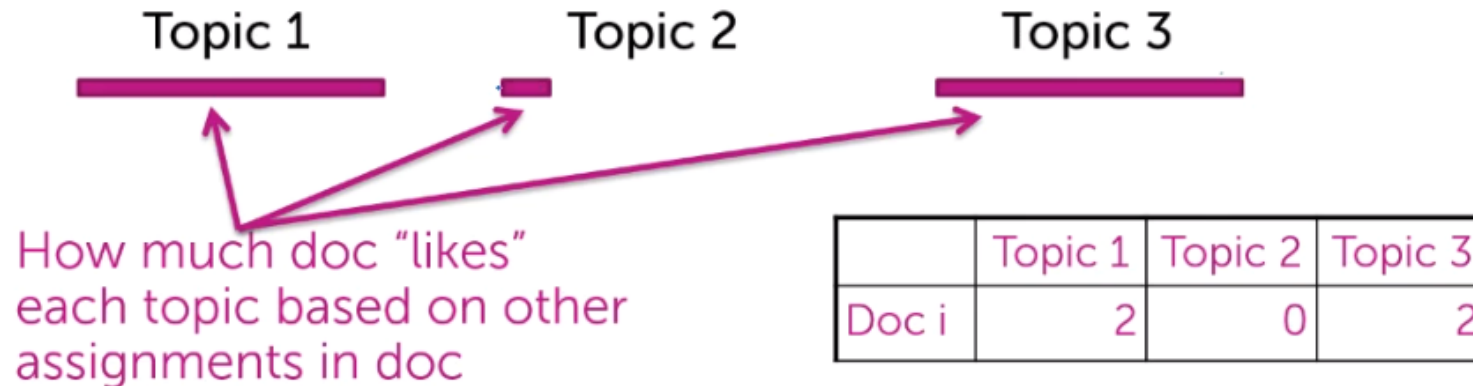
decrementing  
counts  
after removing  
current assignment  
 $Z_{iw} = 2$

# Collapsed Gibbs Sampling

## Probability of new assignment

3	?	1	3	1
epilepsy	dynamic	Bayesian	EEG	model

- K is number of topics.
- The calculation result becomes the bar length.



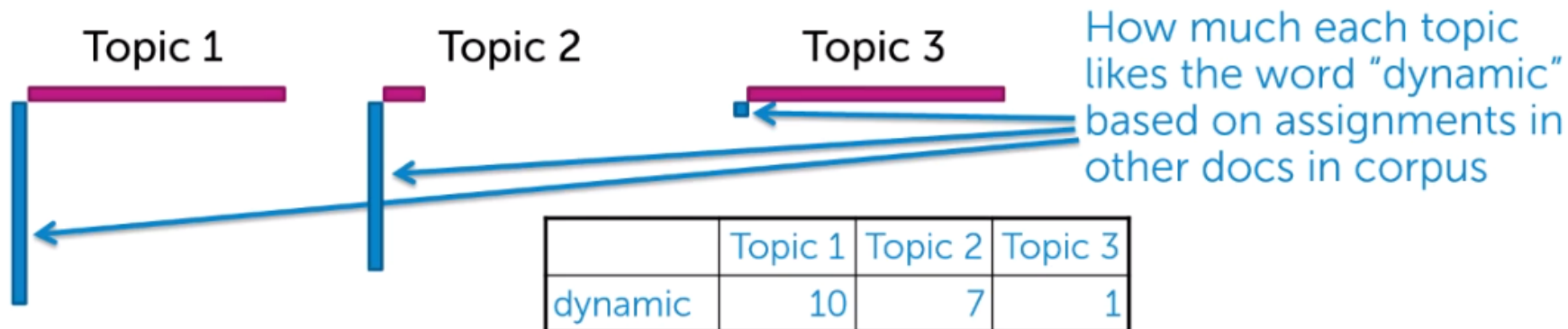
# current assignments to topic k in doc i  $\rightarrow n_{ik} + \alpha$  smoothing param from Bayes prior

# words in doc i  $\rightarrow N_i - 1 + K\alpha$  ignore current word

# Collapsed Gibbs Sampling

## Probability of new assignment

3	?	1	3	1
epilepsy	dynamic	Bayesian	EEG	model



# assignments  
**corpus-wide** of  
word "dynamic"  
to topic k

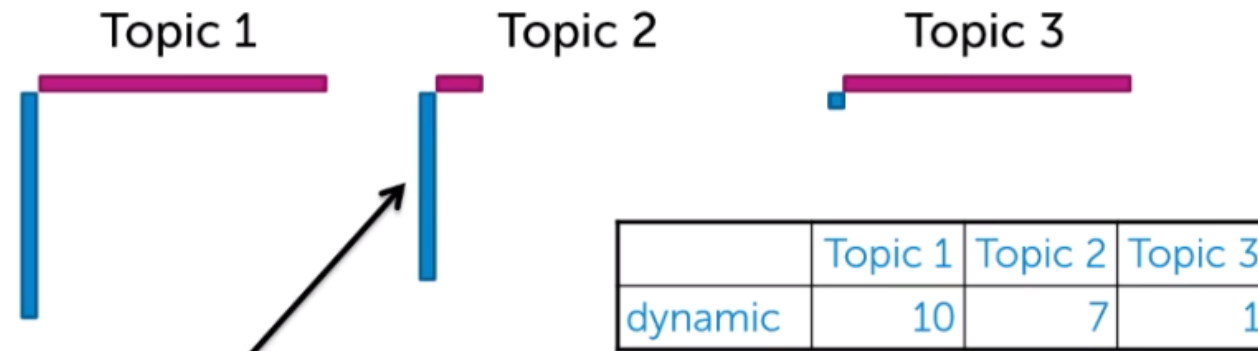
$$\frac{m_{\text{dynamic},k} + \gamma}{\sum_{w \in V} m_{w,k} + V\gamma}$$

smoothing param *from Bayes prior*

*size of vocab*

- $\text{sum}(m_{w,k})$  equals to the sum of all vocabulary count of a topic.
- $V$  is the size of vocabulary.
- The calculation result becomes the bar length.

# Collapsed Gibbs Sampling



Topic 2 also really likes "dynamic",  
but in a different context...  
e.g., a topic on fluid dynamics

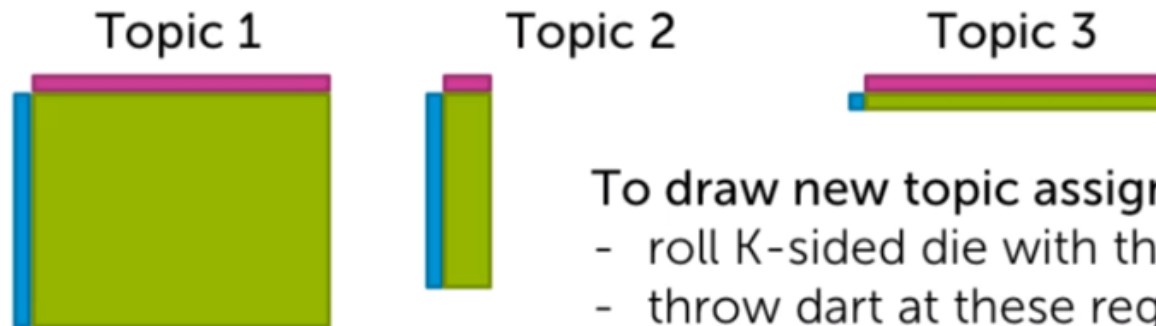


# Collapsed Gibbs Sampling

Randomly draw a new topic indicator

- Two terms are multiplied together to get the area (probability).

3	?	1	3	1
epilepsy	dynamic	Bayesian	EEG	model



To draw new topic assignment (equivalently):

- roll K-sided die with these probabilities
- throw dart at these regions

Normalize this product of terms over K possible topics!

How much doc likes topic

$$\frac{n_{ik} + \alpha}{N_i - 1 + K\alpha}$$

$$\frac{m_{\text{dynamic},k} + \gamma}{\sum_{w \in V} m_{w,k} + V\gamma}$$

How much topic likes word

# Collapsed Gibbs Sampling

## Update counts

3	1	1	3	1
epilepsy	dynamic	Bayesian	EEG	model

	Topic 1	Topic 2	Topic 3
epilepsy	1	0	35
Bayesian	50	0	1
model	42	1	0
EEG	0	0	20
dynamic	11	7	1
...			

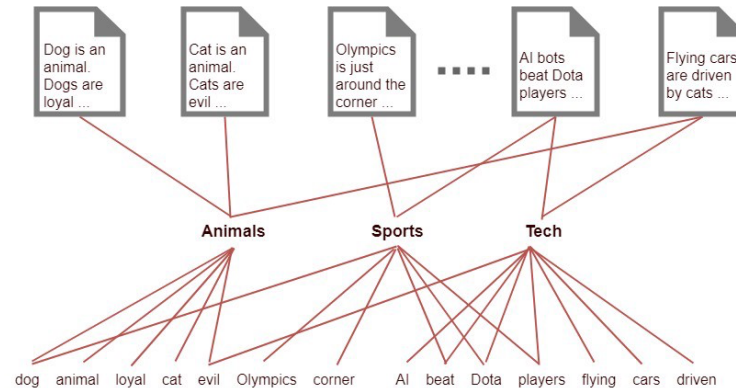
	Topic 1	Topic 2	Topic 3
Doc i	3	0	2

increment counts  
based on new  
assignment of  
 $z_{iw}=1$

- After updating the topic assignment of the word “dynamic”, go to the next word “Bayesian” in this five-word document “Epilepsy dynamic Bayesian EEG model”.
- After we finish this doc, go to the next doc.
- After we finish the corpus, go through the corpus again.

# So why called Latent Dirichlet Allocation?

- Latent  
The topics are hidden.



- Dirichlet  
E.g. A machine can produce different dice with different biased weights, and each dice itself is a distribution as we get multiple values when we roll a dice. This is what it means to be a distribution of distributions and this is what Dirichlet is.  
Here, in the context of topic modelling, the Dirichlet is the  
[first step] distribution of topics in documents, and  
[second step] distribution of words in the topic.  
(Notice that this is the document generating process, which is different from the Gibbs Sampling.)
- Allocation  
We allocate topics to the documents, and words (of the document) to topics.