



JOHNS HOPKINS

WHITING SCHOOL
of ENGINEERING

Topic Modeling (LDA)

1/29/24

Recap

- Last class:
 - Metrics to measure differences in word usage across subsets of corpora
 - Log Odds with Dirichlet Prior (Fightin' Words)
 - PMI Scores
- Today
 - Topic modeling (LDA)
 - Inference method 1: Gibbs sampling
 - Practical considerations



JOHNS HOPKINS

WHITING SCHOOL
of ENGINEERING

LDA Introduction

Odds ratio in Congressional data

Top Republican Words	Score		Top Democrat Words	Score
spending	-66.26		republican	56.63
obamacare	-59.90		wealthiest	40.78
government	-47.92		rhode	39.43
going	-45.33		women	38.16
that	-44.58		pollution	33.66
trillion	-43.43		republicans	32.86
taxes	-42.39		gun	32.45
you	-40.85		investments	32.22
administration	-39.07		families	31.93
debt	-38.92		violence	30.88

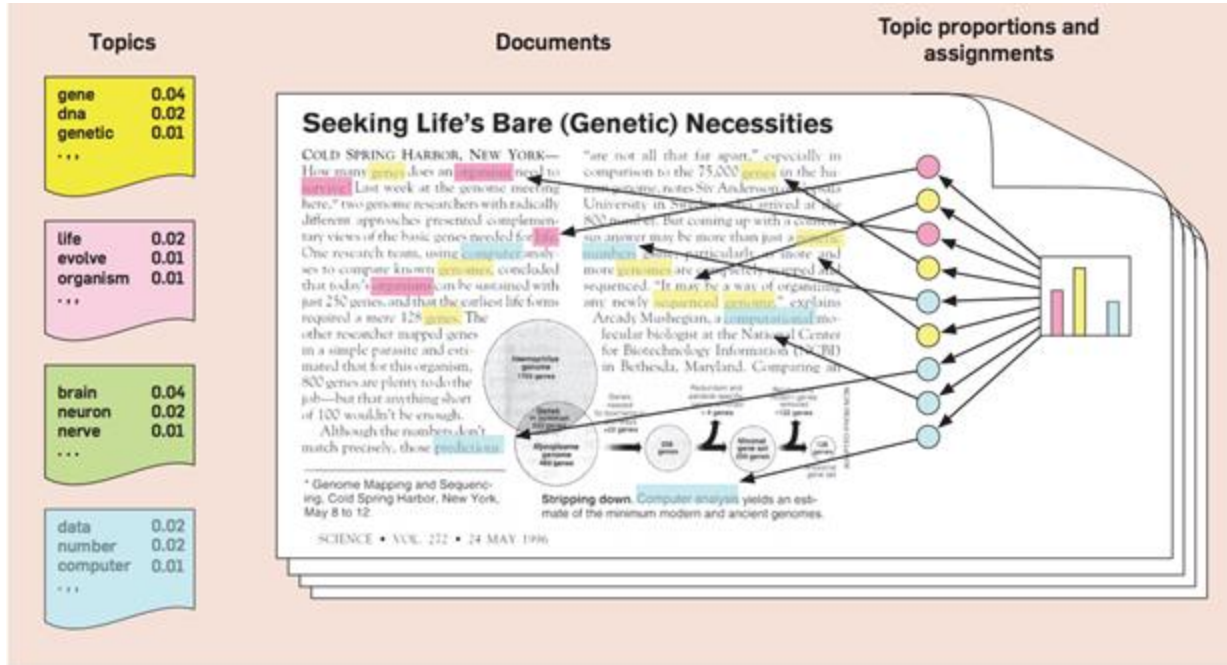
Probably all about budget and government spending

“Gun violence” is probably one topic

Topic Modeling: Motivation

- Sometimes we care about specific words (more on this later)
- Often we want to group words into broader *topics*
 - But we don't know these topics in advance, we need to discover them from the data

Latent Dirichlet Allocation



- Assume each document contains a mixture of "topics"
- Each topic uses mixtures of vocabulary words
- **Goal: recover topic and vocabulary distributions**

Definitions

	Topic 1	Topic 2	...	Topic 30
administration	0.01	0.12	...	0.02
advertising	0.02	0.001	...	0.25
debt	0.1	0.001	...	0.01
...
government	0.01	0.15	...	0.01
...
spending	0.12	0.01	...	0.03
taxes	0.15	0.02	...	0.35
trillion	0.19	0.003	...	0.02

Each “topic” is defined by ϕ , a multinomial distribution over the entire vocabulary

	Doc 1	Doc 2	...	Doc N
Topic 1	0.10	0.60	...	
Topic 3	0.02	0.05	...	
Topic 4	0.30	0.1	...	
...
Topic 15	0.20	0.01	...	0.40
...
Topic 28	0.01	0.03	...	0.20
Topic 29	0.25	0.15	...	
Topic 30	0.03	0.01	...	

Each document has associated θ , a multinomial distribution over topics

LDA Generative Story

Basic idea:

- Assume a story for generating our data (sampling from distributions)
- Estimate the parameters of the distribution
- [There are other approaches to topic modeling, this is specifically LDA]

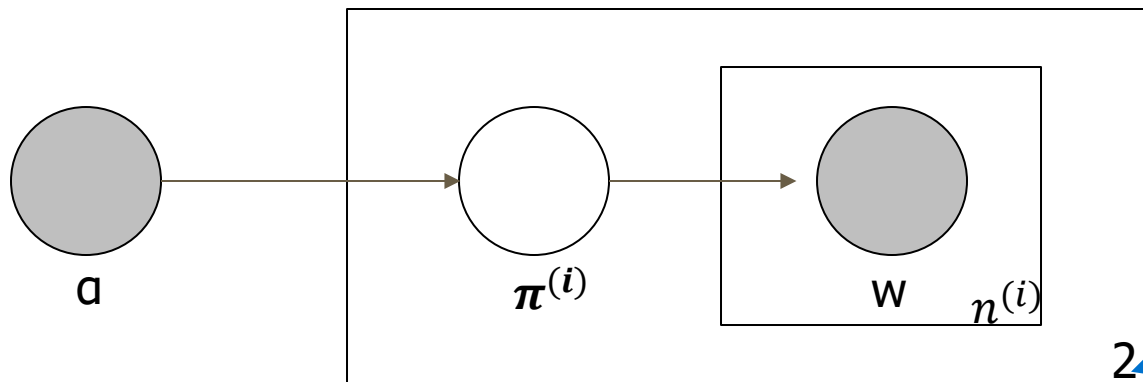
Fightin' Words Generative Story

Generative story for log-odds with a Dirichlet Prior:

1. Draw $\boldsymbol{\pi}^{(i)} \sim \text{Dirichlet}(\boldsymbol{\alpha})$
2. For $n^{(i)}$ steps:
 1. Draw $w \sim \text{Multinomial}(\boldsymbol{\pi}^{(i)})$

Plate Notation: Log-odds with Dirichlet prior

- Shaded circle: value we observe
- Rectangles: values that are repeated (with number in corner reflecting # of repetitions)



We drew π for Democrats and π for Republicans

LDA Generative Story

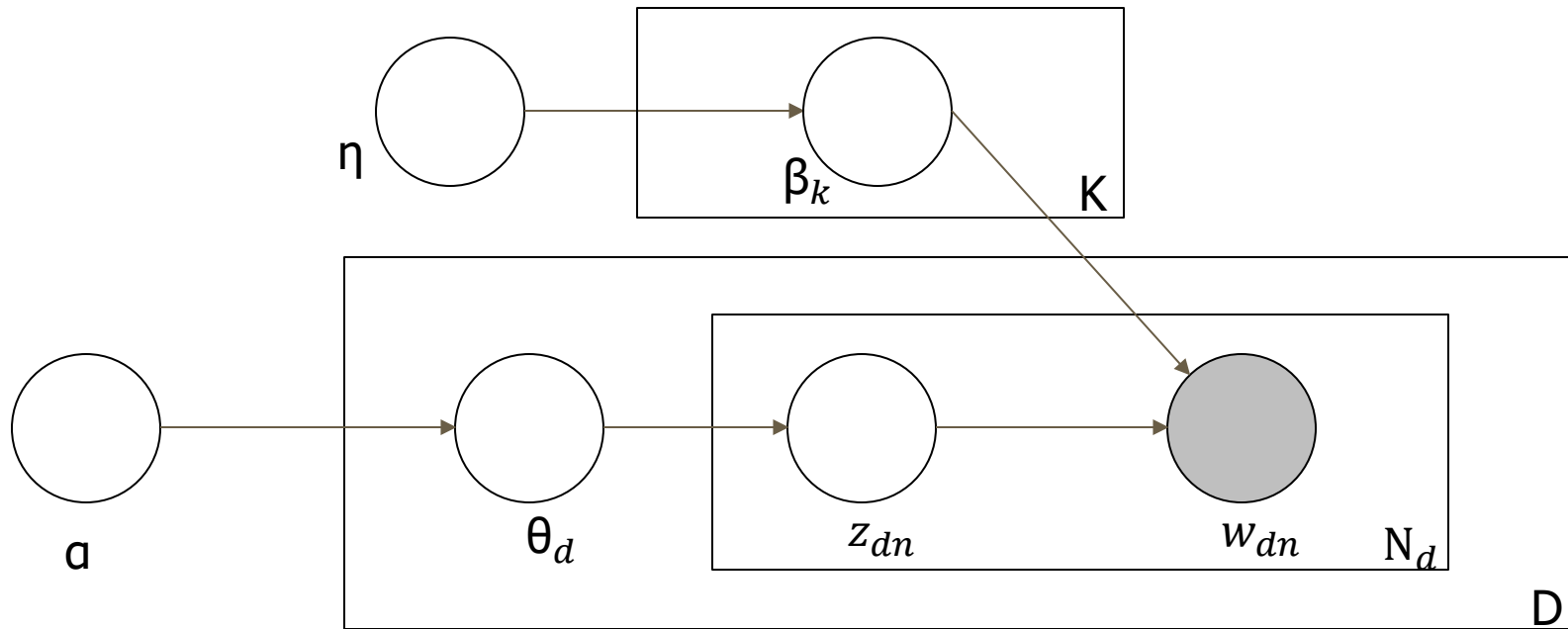
- For each topic k :
 - Draw $\beta_k \sim \text{Dir}(\eta)$
- For each document d :
 - Draw $\theta_d \sim \text{Dir}(\alpha)$
 - For each word in d :
 - Draw topic assignment $z \sim \text{Multinomial}(\theta_d)$
 - Draw $w \sim \text{Multinomial}(\beta_z)$

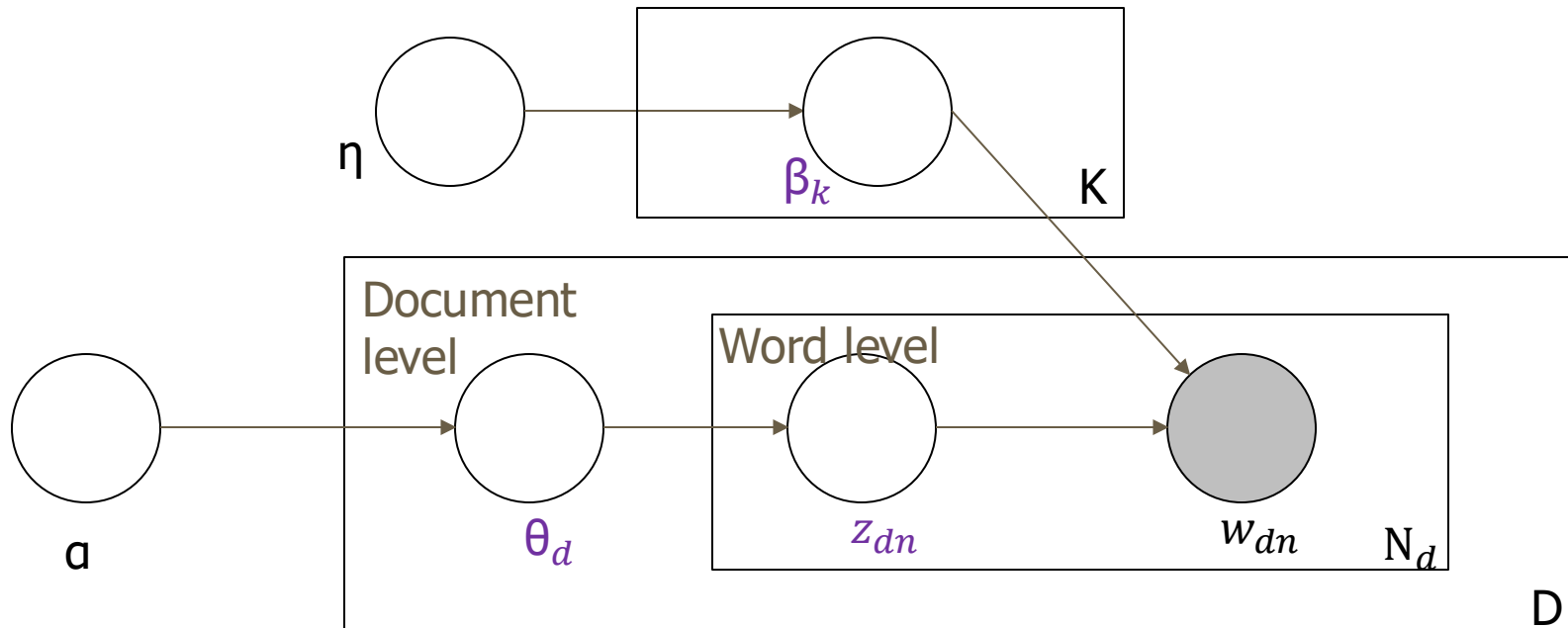
We use the data to estimate these two sets of parameters:

- β , a distribution over vocabulary (1 for each topic)
- θ , a distribution over topics (1 for each document)

LDA Generative Story

- For each topic k :
 - Draw $\beta_k \sim \text{Dir}(\eta)$
 - For each document d :
 - Draw $\theta_d \sim \text{Dir}(\alpha)$ → As long as θ_d is sparse, each document should be most affiliated with a few topics
 - For each word in d :
 - Draw topic assignment $z \sim \text{Multinomial}(\theta_d)$ → The document's topic influences what words are in it:
 - Draw $w \sim \text{Multinomial}(\beta_z)$
- β , a distribution over vocabulary (1 for each topic)
 - θ , a distribution over topics (1 for each document)
- words that co-occur in the same document should end up affiliated with the same topic
 - Documents with similar words will end up with similar topics





Variables we observe: D = number of documents; N = number of words per document, w words in document

Variables we want to estimate: θ , β , z are latent variables

Variables we choose: α , η are hyperparameters. K = number of topics

General Estimators [Heinrich, 2005]

Goal: estimate θ, β

$$p(\theta, \beta, z | w) = \frac{p(w | \theta, \beta, z)p(\theta, \beta, z)}{p(w)}$$

- MLE approach
 - Maximize likelihood: $p(w | \theta, \beta, z)$
- MAP approach
 - Maximize posterior: $p(\theta, \beta, z | w)$ OR $p(w | \theta, \beta, z)p(\theta, \beta, z)$
- Bayesian approach
 - Approximate posterior: $p(\theta, \beta, z | w)$
 - Take expectation of posterior to get point estimates

LDA: Bayesian Inference

- Goal: estimate θ, β
- Bayesian approach: we estimate full posterior distribution

$$p(\theta, \phi, z | w) = \frac{p(w | \theta, \beta, z)p(\theta, \beta, z)}{p(w)}$$

$p(w)$ is the probability of your data set occurring under *any* parameters -- this is intractable!

Solutions: Gibbs Sampling, Variational Inference

Quiz

1. How many elements does each θ_i have?
 - A. The number of words in document i (N_i)
 - B. The number of documents in the corpus (D)
 - C. The number of topics specified by the researcher (K)
 - D. The number of words in the vocabulary
2. How many elements does each β_j have?
 - A. The number of words in document i (N_i)
 - B. The number of documents in the corpus (D)
 - C. The number of topics specified by the researcher (K)
 - D. The number of words in the vocabulary
3. Which variables are observed?
 - A. D, N, w
 - B. θ, β, z
 - C. α, η, K





JOHNS HOPKINS
WHITING SCHOOL
of ENGINEERING

Gibbs Sampling

Gibbs Sampling

- A Markov chain Monte Carlo (MCMC) algorithm
 - Algorithms for drawing samples from a probability distribution
- We draw samples by constructing a Markov Chain: the probability of the next sample is calculated from the previous sample
- We construct the chain so that if we draw enough samples (the “burn-in” period), we eventually start drawing samples from our real target distribution
- Once we have samples from the target distribution, we can use them to estimate the parameters we care about

Gibbs Sampling

- Assume we know topic assignments z for all words in the corpus

Vastly available digitized text data has created new opportunities for understanding social phenomena. Relatedly, social issues like toxicity, discrimination, and propaganda frequently manifest in text, making text analyses critical for understanding and mitigating them. In this

- We know how many times each word has been assigned to each topic
- We know how many times each topic has been assigned to each document

	Topic 1	Topic 2	Topic 3	...
Social	5	0	2	
analyses	10	3	2	
Discrimination	1	10	2	
...	

	Topic 1	Topic 2	Topic 3	...
Doc 1	11	7	30	
Doc 2	2	22	1	
Doc 3	16	15	17	
...	



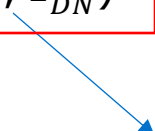
Gibbs Sampling

Vastly available digitized text data has created new opportunities for understanding social phenomena. Relatedly, social issues like toxicity, discrimination, and propaganda frequently manifest in text, making text analyses critical for understanding and mitigating them. In this

- One word at a time, remove the topic assignment and resample it

Remember the high-level: if we do this enough times, we start getting “good” topic assignments that we can use to estimate the parameters we care about

Gibbs Sampling

- Initialize z (e.g. randomly)
- for $t = 1$ to T do:  For each iteration
 - for $d = 1$ to D ; for $n = 1$ to N_d do:  For each word in the corpus
 - $z_{dn}^{(t+1)} \sim P(Z_{dn} \mid z_{11}^{(t+1)} \dots, z_{dn-1}^{(t+1)}, z_{dn+1}^{(t)}, \dots, z_{DN}^{(t)})$ 
 - end for
- end for

Gibbs Sampling

$$P(z_{dn} = k | z_{d,-n}, w, \alpha, \eta, \beta, \theta)$$

- We integrate out ϕ, β (we can do this because of conjugacy)

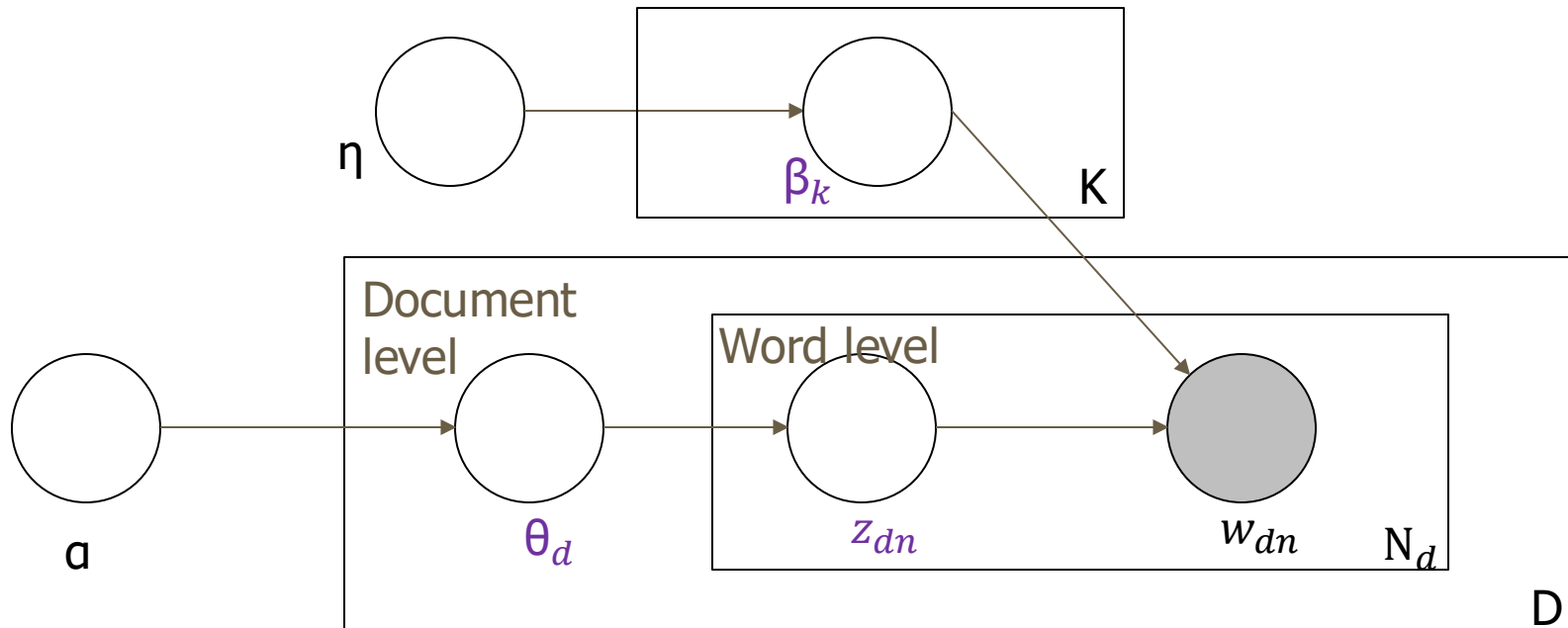
Number of times document d uses topic k

From prior

Number of times topic k uses word w_{dn}

From prior

$$P(z_{dn} = k | z_{d,-n}, w, \alpha, \eta) = \frac{c_{dk} + \alpha_k}{\sum_i^K c_{di} + \alpha_i} \frac{v_{kw_{dn}} + \eta_{w_{dn}}}{\sum_i v_{ki} + \eta_i}$$



Variables we observe: D = number of documents; N = number of words per document, w words in document

Variables we want to estimate: θ , β , z are latent variables

Variables we choose: α , η are hyperparameters. K = number of topics

Gibbs Sampling

$$P(z_{dn} = k | z_{d,-n}, w, \alpha, \eta) = \frac{\frac{c_{dk} + \alpha_k}{\sum_i^K c_{di} + \alpha_i}}{\frac{v_{kw_{dn}} + \eta_{w_{dn}}}{\sum_i v_{ki} + \eta_i}}$$

Prevalence of topic in document Prevalence of word in topic

What make a topic k more likely to be assigned to z_{dn} ? What properties does that mean we would expect to see in our final topic estimates?

What happens if α is very high?

Gibbs Sampling

- Initialize z (e.g. randomly)
- for $t = 1$ to T do:
 - for $d = 1$ to D ; for $n = 1$ to N_d do:
 - $z_{dn}^{(t+1)} \sim P(Z_{dn} \mid z_{11}^{(t+1)} \dots, z_{dn-1}^{(t+1)}, z_{dn+1}^{(t)}, \dots, z_{DN}^{(t)})$
 - end for
- end for

- We can similarly use counts of topic assignments across multiple samples to estimate (β, θ)

Recap

- Goal: estimate θ, β
- Bayesian approach: we estimate full posterior distribution

$$p(\theta, \beta, z | w) = \frac{p(w | \theta, \beta, z)p(\theta, \beta, z)}{p(w)}$$

Gibbs Sampling:

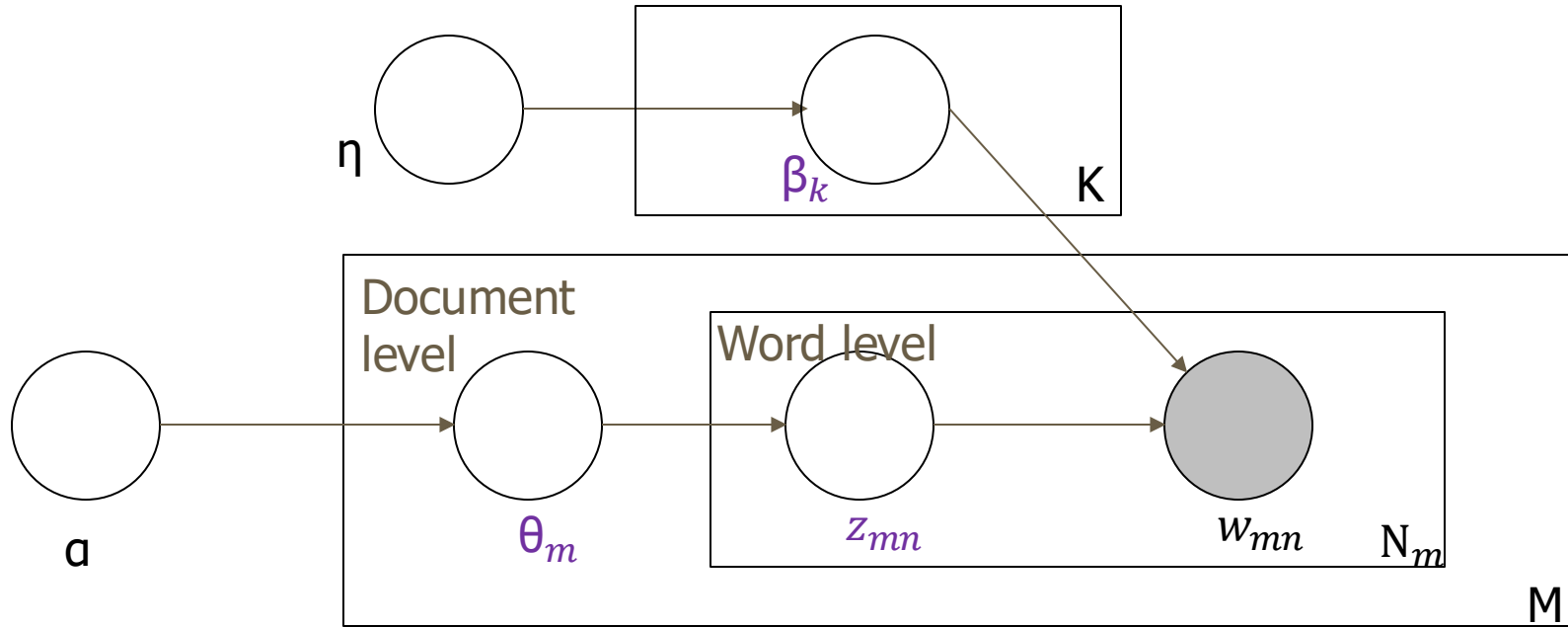
- We generate samples from the posterior distribution
- We estimate θ, β from those samples



JOHNS HOPKINS

WHITING SCHOOL
of ENGINEERING

LDA In Practice



Variables we observe: M = number of documents; N = number of words per document, w words in document

Variables we want to estimate: θ , β , z are latent variables

Variables we choose: α , η are hyperparameters. K = number of topics

Choosing α , η and K

- In practice, typically choose *symmetric* Dirichlet priors, e.g. $\alpha, \eta = [1, 1, 1, 1, \dots], [0.1, 0.1, 0.1, 0.1, \dots]$ but some research has explored alternatives
- In practice, try a few K values and judge if topics look reasonable, but there are approaches that estimate the best value

Sample Topics from NYT Corpus

#5	#6	#7	#8	#9	#10
10	0	he	court	had	sunday
30	tax	his	law	quarter	saturday
11	year	mr	case	points	friday
12	reports	said	federal	first	van
15	million	him	judge	second	weekend
13	credit	who	mr	year	gallery
14	taxes	had	lawyer	were	iowa
20	income	has	commission	last	duke
sept	included	when	legal	third	fair
16	500	not	lawyers	won	show

How do we describe a topic?

- Most probable words for each topic
- Words common in this topic *relative* to other topics
 - We could use PMI scores!
- Examining documents that contain high proportion of topic

LDA: Evaluation

- Held-out likelihood
 - Hold out some subset of your corpus
 - Compute the likelihood of the held-out data under the parameters you estimated
 - Says NOTHING about coherence of topics

Intruder Detection Tasks

Word Intrusion

1 / 10
floppy alphabet computer processor memory disk

2 / 10
molecule education study university school student

3 / 10
linguistics actor film comedy director movie

4 / 10
islands island bird coast portuguese mainland

Topic Intrusion

6 / 10

DOUGLAS HOFSTADTER

Douglas Richard Hofstadter (born February 15, 1945 in New York, New York) is an American academic whose research focuses on consciousness, thinking and creativity. He is best known for "[Show entire excerpt](#)", first published in

student	school	study	education	research	university	science	learn
human	life	scientific	science	scientist	experiment	work	idea
play	role	good	actor	star	career	show	performance
write	work	book	publish	life	friend	influence	father

Key idea: If topics are coherent, annotators should easily be able to identify the intruder

LDA: Evaluation

- Can we automate these judgements?
- Word intrusion detection:
 - Compute PMI scores between each pair of words in the set of real and intruder words for each topic; train an SVM model to learn intruder words
- [Similar heuristics for topic intrusion detection]
- But follow-up work has suggested these kinds metrics don't always correlate with human judgement

Practical advice for getting coherent topics

- Evaluate the topics by hand
- Hyperparameter selection (α , η , K):
 - Test different numbers of topics
 - Tune the parameter controlling the topic distributions
- Pay attention to your data:
 - Better for long documents
 - Keep stopwords, don't stem
 - Remove (high numbers of) duplicates
- Just use LDA!

<https://maria-antoniak.github.io/2022/07/27/topic-modeling-for-the-people.html>

[Pulling out the stops: Rethinking stopwords removal for topic models](#) by Alexandra Schofield, Måns Magnusson, and

David Mimno (EACL, 2017)

LDA: Advantages and Drawbacks

- When to use it
 - Initial investigation into unknown corpus
 - Concise description of corpus (dimensionality reduction)
 - [Features in downstream task]
- Limitations
 - Can't apply to specific questions (completely unsupervised)
 - Simplified word representations
 - BOW model
 - Can't take advantage of similar word
 - Strict assumptions
 - E.g. Independence assumptions

Today's takeaways

- Motivation behind topic modeling
- **LDA formulation**
- LDA inference
 - Gibbs sampling (overview of process, e.g. slide 24)
- Evaluation and practical considerations
- Next class:
 - Variational inference
 - LDA extensions
 - More in-depth use cases

Logistics

- HW 1 to be released by the end of this week
- We haven't yet covered all of it, but we have covered part 1

References

1. Blei, David M., Andrew Y. Ng, and Michael I. Jordan. "Latent dirichlet allocation." *Journal of machine Learning research* 3.Jan (2003): 993-1022.

Optional sources for more depth:

- Gibbs Sampling:

- Jordan Boyd-Graber's Introduction: <https://www.youtube.com/watch?v=u7l5hhmdc0M>
- <https://api.drum.lib.umd.edu/server/api/core/bitstreams/a36ce44d-0732-427d-8a81-a18c9b0b4dfa/content>
- UMD Technical Report: <https://drum.lib.umd.edu/items/d5aa258e-d2ac-4529-8831-ec0e08a5f2cc>

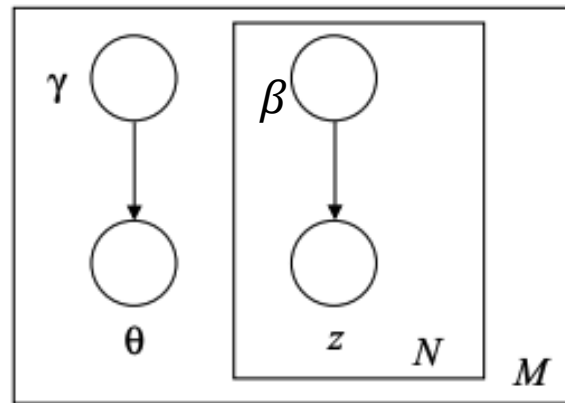
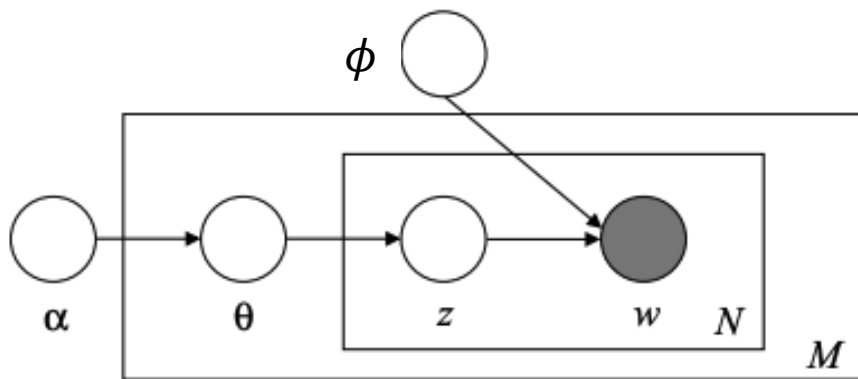
- Variational Inference

- Blei, David M., Alp Kucukelbir, and Jon D. McAuliffe. "Variational inference: A review for statisticians." *Journal of the American statistical Association* 112.518 (2017): 859-877.
- Xanda Schofield and Jordan Boyd-Graber
 - <https://www.youtube.com/watch?v=-tKmyHoVZ-g>
 - <https://www.youtube.com/watch?v=smfWKhDcaoA>
- David Blei Lecture Notes <https://www.cs.princeton.edu/courses/archive/fall11/cos597C/lectures/variational-inference-i.pdf>

Choose q

$$p(\theta, \phi, z | w) = \frac{p(w | \theta, \phi, z)p(\theta, \phi, z)}{p(w)}$$

- Choose $q(\theta, \phi, z) = \prod_{i=1}^K q(\phi_i | \lambda_i) \prod_{m=1}^M q_m(\theta_m, z_m | \gamma_m, \beta_m)$
 - where $q_m(\theta, z) = q(\theta | \gamma) \prod_i q(z_n | \beta_n)$
 - Assume $q(\theta | \gamma)$ is a Dirichlet distribution with variational parameters γ
 - Assume $q(\phi | \lambda)$ is a Dirichlet distribution with variational parameters λ
 - Assume $q(z_n | \beta_n)$ is a multinomial (categorical) distribution with variational parameters β_n



Full procedure

- Choose q
 - For each document: $q(z, \theta) = q(\theta|\gamma) \prod_i^N q(z_n|\beta_n)$
- For each iteration
 - For each word in each document n
 - For each possible topic assignment i
 - $\beta_{ni} \propto \phi_{i w_n} \exp\{E_q[\log(\theta_i)|\gamma]\}$
 - $\gamma = \alpha + \sum_{n=1}^N \beta_n$
 - [Check that ELBO increased; e.g. q is moving closer to p]
- End at convergence

[Use q to approximate posterior: we can take expectations of q to estimate parameters]

Choose q

$$p(\theta, \phi, z | w) = \frac{p(w | \theta, \phi, z)p(\theta, \phi, z)}{p(w)}$$

[let's ignore ϕ]

- For each document $q(\theta, z) = q(\theta | \gamma) \prod_i q(z_n | \beta_n)$
 - Assume $q(\theta | \gamma)$ is a Dirichlet distribution with variational parameters γ
 - Assume $q(z_n | \beta_n)$ is a multinomial (categorical) distribution with variational parameters β_n

