# Collapsed Gibbs sampling in LDA
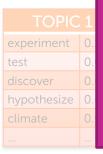
# "Collapsed" Gibbs sampling for LDA

Based on special structure of LDA model, can sample **just** indicator variables $z_{iw}$

- No need to sample other parameters
  - corpus-wide topic vocab distributions
  - per-doc topic proportions

Often leads to much better performance because examining uncertainty in smaller space

# Collapsed Gibbs sampling for LDA



**Never draw topic vocab distributions or doc topic proportions**

Randomly reassign $z_{iw}$ based on current assignments $z_{jv}$ of all other words **in document and corpus**

Machine Learning Specialization

# Select a document

| | | | | |
|---|---|---|---|---|
| epilepsy | dynamic | Bayesian | EEG | model |

5 word document

Machine Learning Specialization

# Randomly assign topics

| 3 | 2 | 1 | 3 | 1 |
|---|---|---|---|---|
| epilepsy | dynamic | Bayesian | EEG | model |

(one possible approach)

Machine Learning Specialization

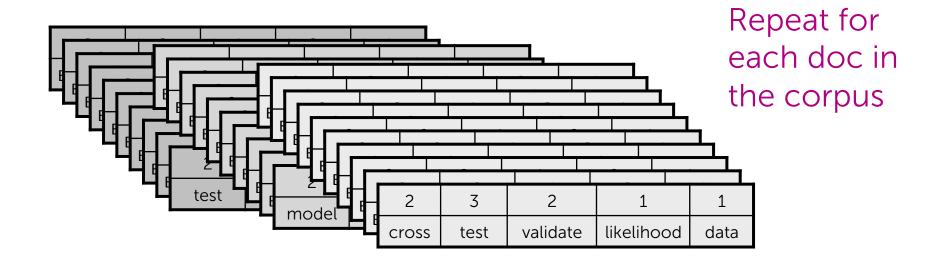# Randomly assign topics

| 3 | 2 | 1 | 3 | 1 |
|---|---|---|---|---|
| epilepsy | dynamic | Bayesian | EEG | model |

Repeat for each doc in the corpus

| 2 | 3 | 2 | 1 | 1 |
|---|---|---|---|---|
| cross | test | validate | likelihood | data |

# Maintain local statistics

| 3 | 2 | 1 | 3 | 1 |
|---|---|---|---|---|
| epilepsy | dynamic | Bayesian | EEG | model |

|  | Topic 1 | Topic 2 | Topic 3 |
|---|---|---|---|
| Doc i | 2 | 1 | 2 |

Machine Learning Specialization

# Maintain global statistics

| 3 | 2 | 1 | 3 | 1 |
|---|---|---|---|---|
| epilepsy | dynamic | Bayesian | EEG | model |

|  | Topic 1 | Topic 2 | Topic 3 |
|---|---|---|---|
| Doc i | 2 | 1 | 2 |

|  | 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 |
| ... |  |  |  |

Total counts from **all** docs

# Randomly reassign topics

| 3 | 2̶ | 1 | 3 | 1 |
|---|---|---|---|---|
| epilepsy | dynamic | Bayesian | EEG | model |

|  | Topic 1 | Topic 2 | Topic 3 |
|---|---|---|---|
| Doc i | 2 | 0̶ 1 | 2 |

|  | Topic 1 | Topic 2 | Topic 3 |
|---|---|---|---|
| epilepsy | 1 | 0 | 35 |
| Bayesian | 50 | 0 | 1 |
| model | 42 | 1 | 0 |
| EEG | 0 | 0 | 20 |
| dynamic | 10 | 7 8̶ | 1 |
| ... |  |  |  |

decrementing counts after removing current assignment

$z_{iw} = 2$

# Probability of new assignment

| 3 | ? | 1 | 3 | 1 |
|---|---|---|---|---|
| epilepsy | dynamic | Bayesian | EEG | model |

reassign with probability
$p(z_{iw} \mid \text{every other } z_{jv} \text{ in corpus, words in corpus})$

# Probability of new assignment

| 3 | ? | 1 | 3 | 1 |
|---|---|---|---|---|
| epilepsy | dynamic | Bayesian | EEG | model |

Topic 1        Topic 2        Topic 3

How much doc "likes" each topic based on other assignments in doc

| | Topic 1 | Topic 2 | Topic 3 |
|---|---|---|---|
| Doc i | 2 | 0 | 2 |

# current assignments to topic k in doc i

# words in doc i

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

smoothing param

from Bayes prior

ignore current word

Machine Learning Specialization

# Probability of new assignment

| 3 | ? | 1 | 3 | 1 |
|---|---|---|---|---|
| epilepsy | dynamic | Bayesian | EEG | model |

**Topic 1**          **Topic 2**          **Topic 3**

How much each topic
likes the word "dynamic"
based on assignments in
other docs in corpus

|  | Topic 1 | Topic 2 | Topic 3 |
|---|---|---|---|
| dynamic | 10 | 7 | 1 |

\# 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*

# Probability of new assignment

| 3 | ? | 1 | 3 | 1 |
|---|---|---|---|---|
| epilepsy | dynamic | Bayesian | EEG | model |

**Topic 1**    **Topic 2**    **Topic 3**

|  | Topic 1 | Topic 2 | Topic 3 |
|---|---|---|---|
| dynamic | 10 | 7 | 1 |

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

# Probability of new assignment

| 3 | ? | 1 | 3 | 1 |
|---|---|---|---|---|
| epilepsy | dynamic | Bayesian | EEG | model |

**Topic 1**

**Topic 2**

**Topic 3**

Topic fits word **and** document

Topic fits word, but not doc

Topic fits doc, but not word

# Probability of new assignment

| 3 | ? | 1 | 3 | 1 |
|---|---|---|---|---|
| epilepsy | dynamic | Bayesian | EEG | model |

Topic 1

Topic 2

Topic 3

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

# Randomly draw a new topic indicator

| 3 | ? | 1 | 3 | 1 |
|---|---|---|---|---|
| epilepsy | dynamic | Bayesian | EEG | model |

**Topic 1**        **Topic 2**        **Topic 3**

**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

$$\frac{n_{ik} + \alpha}{N_i - 1 + K\alpha} \quad \frac{m_{\text{dynamic},k} + \gamma}{\sum_{w \in V} m_{w,k} + V\gamma}$$

How much topic likes word

# Update counts

| 3 | 1 | 1 | 3 | 1 |
|---|---|---|---|---|
| epilepsy | dynamic | Bayesian | EEG | model |

|  | Topic 1 | Topic 2 | Topic 3 |
|---|---|---|---|
| Doc i | 3 ~~2~~ | 0 | 2 |

|  | Topic 1 | Topic 2 | Topic 3 |
|---|---|---|---|
| epilepsy | 1 | 0 | 35 |
| Bayesian | 50 | 0 | 1 |
| model | 42 | 1 | 0 |
| EEG | 0 | 0 | 20 |
| dynamic | 11 ~~10~~ | 7 | 1 |
| … |  |  |  |

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

# Geometrically...

| 3 | 1 | 1 | 3 | 1 |
|---|---|---|---|---|
| epilepsy | dynamic | Bayesian | EEG | model |

Topic 1          Topic 2                    Topic 3

Increase popularity of "dynamic" in topic 1 (**corpus-wide**)

Increase popularity of topic 1 **in doc i**

# Iterate through all words/docs

| 3 | 1 | 1 | 3 | 1 |
|---|---|---|---|---|
| epilepsy | dynamic | Bayesian | EEG | model |

|  | Topic 1 | Topic 2 | Topic 3 |
|---|---|---|---|
| Doc i | 2 | 0 | 2 |

|  | Topic 1 | Topic 2 | Topic 3 |
|---|---|---|---|
| epilepsy | 1 | 0 | 35 |
| Bayesian | 50 | 0 | 1 |
| model | 42 | 1 | 0 |
| EEG | 0 | 0 | 20 |
| dynamic | 10 | 7 | 1 |
| ... |  |  |  |

Machine Learning Specialization

# Using samples from collapsed Gibbs

# What to do with the collapsed samples?



From "best" sample of $\{z_{iw}\}$, can infer:

Machine Learning Specialization

# What to do with the collapsed samples?



From "best" sample of $\{z_{iw}\}$, can infer:
1. Topics from conditional distribution...
   need corpus-wide info

# What to do with the collapsed samples?



From "best" sample of $\{z_{iw}\}$, can infer:
1. Topics from conditional distribution...
   need corpus-wide info
2. Document "embedding"...
   need doc info only

Machine Learning Specialization

# Embedding new documents



**Simple approach:**
1. Fix topics based on training set collapsed sampling
2. Run uncollapsed sampler on new doc(s) only

Machine Learning Specialization