

# *AN INTRO TO BAYESIAN ML (PART II)*

**Prof. Chris Jermaine**  
**cmj4@cs.rice.edu**

# OK, So What Does This Have To Do W Text?

- Can apply same methodology to learning topics in text
- A “topic” is a set of words that appear to together with high prob
- Idea:
  - If you can analyze a corpus...
  - And figure out a set of  $k$  topics...
  - As well as how prevalent each topic is in each document
  - You then know a lot about the corpus
  - In our case, can use this prevalence info to search the corpus
  - Two docs have similar topic compositions? Then they are similar!

# OK, So What Does This Have To Do W Text?

- This is exactly what “LDA” does
  - Stands for “Latent Dirichlet Allocation”... fancy!
- As in all BML, we first need a generative process
  - After that, can talk about “learning”
- Basic LDA setup
  - LDA will generate  $n$  random documents given a dictionary
  - Dictionary is of size `num_words`
  - Best shown thru an example
  - In our example: dictionary will have: (0, “bad”) (1, “I”) (2, “can’t”) (3, “stand”) (4, “comp 215”), (5, “to”) (6, “leave”) (7, “love”) (8, “beer”) (9, “humanities”) (10, “classes”)

# LDA Step One

- Generate each of the  $k$  “topics”
  - Each topic is represented by a vector of probabilities
  - The  $w$ th entry in the vector is associated with the  $w$ th word in the dictionary
  - $\text{word\_probs}_t[w]$  is the probability that topic  $t$  would produce word  $w$
  - Vector is sampled from a Dirichlet (alpha) distribution
  - So, for each  $t$  in  $\{0 \dots k - 1\}$ ,  $\text{word\_probs}_t \sim \text{Dirichlet}(\alpha)$

# LDA Step One

- Generate each of the  $k$  “topics”
  - Each topic is represented by a vector of probabilities
  - The  $w$ th entry in the vector is associated with the  $w$ th word in the dictionary
  - $\text{word\_probs}_t[w]$  is the probability that topic  $t$  would produce word  $w$
  - Vector is sampled from a Dirichlet (alpha) distribution
  - So, for each  $t$  in  $\{0 \dots k - 1\}$ ,  $\text{word\_probs}_t \sim \text{Dirichlet}(\alpha)$
- Ex:  $k = 3$ 
  - $\text{word\_probs}_0 = (.2, .2, .2, .2, 0, 0, 0, 0, .2, 0, 0)$
  - $\text{word\_probs}_1 = (0, .2, .2, .2, 0, 0, 0, 0, 0, .2, .2)$
  - $\text{word\_probs}_2 = (0, .2, .2, 0, .2, 0, .2, .2, 0, 0, 0)$

# LDA Step Two

- Generate the topic proportions for each document
  - Each topic “controls” a subset of the words in a document
  - $\text{topic\_probs}_d[t]$  is the probability that an arbitrary word in document  $d$  will be controlled by topic  $t$
  - Vector is sampled from a Dirichlet (beta) distribution
  - So, for each  $d$  in  $\{0 \dots n - 1\}$ ,  $\text{topic\_probs}_d \sim \text{Dirichlet}(\beta)$

# LDA Step Two

- Generate the topic proportions for each document
  - Each topic “controls” a subset of the words in a document
  - $\text{topic\_probs}_d[t]$  is the probability that an arbitrary word in document  $d$  will be controlled by topic  $t$
  - Vector is sampled from a Dirichlet (beta) distribution
  - So, for each  $d$  in  $\{0 \dots n - 1\}$ ,  $\text{topic\_probs}_d \sim \text{Dirichlet}(\text{beta})$
- Ex:  $n = 4$ 
  - $\text{topic\_probs}_0 = (.98, 0.01, 0.01)$
  - $\text{topic\_probs}_1 = (0.01, .98, 0.01)$
  - $\text{topic\_probs}_2 = (0.02, .49, .49)$
  - $\text{topic\_probs}_3 = (.98, 0.01, 0.01)$

# LDA Step Three

- Generate the words in each document
  - Each topic “controls” a subset of the words in a document
  - $\text{words\_in\_doc}_d[w]$  is the number of occurrences of word  $w$  in document  $d$
  - To get this vector, generate the words one-at-a-time
  - For a given word in doc  $d$ :
    - (1) Figure out the topic  $t$  that controls it by sampling from a Multinomial ( $\text{topic\_probs}_d, 1$ ) distribution
    - (2) Generate the word by sampling from a Multinomial ( $\text{word\_probs}_t, 1$ ) distribution



# LDA Step Three

- Generate the words in each document
  - Each topic “controls” a subset of the words in a document
  - $\text{words\_in\_doc}_d[w]$  is the number of occurrences of word  $w$  in document  $d$
  - To get this vector, generate the words one-at-a-time
  - For a given word in doc  $d$ :
    - (1) Figure out the topic  $t$  that controls it by sampling from a Multinomial ( $\text{topic\_probs}_d, 1$ ) distribution
    - (2) Generate the word by sampling from a Multinomial ( $\text{word\_probs}_t, 1$ ) distribution
- Ex: doc 0...  $\text{topic\_probs}_0 = (.98, 0.01, 0.01)$ 
  - $t$  for word zero is...

# LDA Step Three

- Generate the words in each document
  - Each topic “controls” a subset of the words in a document
  - $\text{words\_in\_doc}_d[w]$  is the number of occurrences of word  $w$  in document  $d$
  - To get this vector, generate the words one-at-a-time
  - For a given word in doc  $d$ :
    - (1) Figure out the topic  $t$  that controls it by sampling from a Multinomial ( $\text{topic\_probs}_d, 1$ ) distribution
    - (2) Generate the word by sampling from a Multinomial ( $\text{word\_probs}_t, 1$ ) distribution
- Ex: doc 0...  $\text{topic\_probs}_0 = (.98, 0.01, 0.01)$ 
  - $t$  for word zero is zero, since we sampled  $(1, 0, 0)$  [there is a 1 in the zeroth entry]
  - So we generate the word using  $\text{word\_probs}_0 = (.2, .2, .2, .2, 0, 0, 0, 0, .2, 0, 0)$

# LDA Step Three

- Generate the words in each document
  - Each topic “controls” a subset of the words in a document
  - $\text{words\_in\_doc}_d[w]$  is the number of occurrences of word  $w$  in document  $d$
  - To get this vector, generate the words one-at-a-time
  - For a given word in doc  $d$ :
    - (1) Figure out the topic  $t$  that controls it by sampling from a Multinomial ( $\text{topic\_probs}_d, 1$ ) distribution
    - (2) Generate the word by sampling from a Multinomial ( $\text{word\_probs}_t, 1$ ) distribution
- Ex: doc 0...  $\text{topic\_probs}_0 = (.98, 0.01, 0.01)$  “I”
  - $t$  for word zero is zero, since we sampled  $(1, 0, 0)$  [there is a 1 in the zeroth entry]
  - So we generate the word using  $\text{word\_probs}_0 = (.2, .2, .2, .2, 0, 0, 0, 0, .2, 0, 0)$
  - And we get  $(0, 1, 0, 0, 0, 0, 0, 0, 0, 0, 0)$ , which is equivalent to “I”

# LDA Step Three

- Generate the words in each document
  - Each topic “controls” a subset of the words in a document
  - $\text{words\_in\_doc}_d[w]$  is the number of occurrences of word  $w$  in document  $d$
  - To get this vector, generate the words one-at-a-time
  - For a given word in doc  $d$ :
    - (1) Figure out the topic  $t$  that controls it by sampling from a Multinomial ( $\text{topic\_probs}_d, 1$ ) distribution
    - (2) Generate the word by sampling from a Multinomial ( $\text{word\_probs}_t, 1$ ) distribution
- Ex: doc 0...  $\text{topic\_probs}_0 = (.98, 0.01, 0.01)$  “T”
  - Now onto the next word

# LDA Step Three

- Generate the words in each document
  - Each topic “controls” a subset of the words in a document
  - $\text{words\_in\_doc}_d[w]$  is the number of occurrences of word  $w$  in document  $d$
  - To get this vector, generate the words one-at-a-time
  - For a given word in doc  $d$ :
    - (1) Figure out the topic  $t$  that controls it by sampling from a Multinomial ( $\text{topic\_probs}_d, 1$ ) distribution
    - (2) Generate the word by sampling from a Multinomial ( $\text{word\_probs}_t, 1$ ) distribution
- Ex: doc 0...  $\text{topic\_probs}_0 = (.98, 0.01, 0.01)$  “T”
  - $t$  for word one is zero, since we sampled  $(1, 0, 0)$  [there is a 1 in the zeroth entry]

# LDA Step Three

- Generate the words in each document
  - Each topic “controls” a subset of the words in a document
  - $\text{words\_in\_doc}_d[w]$  is the number of occurrences of word  $w$  in document  $d$
  - To get this vector, generate the words one-at-a-time
  - For a given word in doc  $d$ :
    - (1) Figure out the topic  $t$  that controls it by sampling from a Multinomial ( $\text{topic\_probs}_d, 1$ ) distribution
    - (2) Generate the word by sampling from a Multinomial ( $\text{word\_probs}_t, 1$ ) distribution
- Ex: doc 0...  $\text{topic\_probs}_0 = (.98, 0.01, 0.01)$  “I can’t”
  - $t$  for word one is zero, since we sampled  $(1, 0, 0)$  [there is a 1 in the zeroth entry]
  - So we generate the word using  $\text{word\_probs}_0 = (.2, .2, .2, .2, 0, 0, 0, 0, .2, 0, 0)$
  - And we get  $(0, 0, 1, 0, 0, 0, 0, 0, 0, 0, 0)$ , which is equivalent to “can’t”

# LDA Step Three

- Generate the words in each document
  - Each topic “controls” a subset of the words in a document
  - $\text{words\_in\_doc}_d[w]$  is the number of occurrences of word  $w$  in document  $d$
  - To get this vector, generate the words one-at-a-time
  - For a given word in doc  $d$ :
    - (1) Figure out the topic  $t$  that controls it by sampling from a Multinomial ( $\text{topic\_probs}_d, 1$ ) distribution
    - (2) Generate the word by sampling from a Multinomial ( $\text{word\_probs}_t, 1$ ) distribution
- Ex: doc 0...  $\text{topic\_probs}_0 = (.98, 0.01, 0.01)$  “I can’t”
  - Now onto the next word

# LDA Step Three

- Generate the words in each document
  - Each topic “controls” a subset of the words in a document
  - $\text{words\_in\_doc}_d[w]$  is the number of occurrences of word  $w$  in document  $d$
  - To get this vector, generate the words one-at-a-time
  - For a given word in doc  $d$ :
    - (1) Figure out the topic  $t$  that controls it by sampling from a Multinomial ( $\text{topic\_probs}_d, 1$ ) distribution
    - (2) Generate the word by sampling from a Multinomial ( $\text{word\_probs}_t, 1$ ) distribution
- Ex: doc 0...  $\text{topic\_probs}_0 = (.98, 0.01, 0.01)$  “I can’t”
  - $t$  for word two is zero, since we sampled  $(1, 0, 0)$  [there is a 1 in the zeroth entry]



# LDA Step Three

- Generate the words in each document
  - Each topic “controls” a subset of the words in a document
  - $\text{words\_in\_doc}_d[w]$  is the number of occurrences of word  $w$  in document  $d$
  - To get this vector, generate the words one-at-a-time
  - For a given word in doc  $d$ :
    - (1) Figure out the topic  $t$  that controls it by sampling from a Multinomial ( $\text{topic\_probs}_d, 1$ ) distribution
    - (2) Generate the word by sampling from a Multinomial ( $\text{word\_probs}_t, 1$ ) distribution
- Ex: doc 0...  $\text{topic\_probs}_0 = (.98, 0.01, 0.01)$  “I can’t stand”
  - $t$  for word two is zero, since we sampled  $(1, 0, 0)$  [there is a 1 in the zeroth entry]
  - So we generate the word using  $\text{word\_probs}_0 = (.2, .2, .2, .2, 0, 0, 0, 0, .2, 0, 0)$
  - And we get  $(0, 0, 0, 1, 0, 0, 0, 0, 0, 0, 0)$ , which is equivalent to “stand”

# LDA Step Three

- Generate the words in each document
  - Each topic “controls” a subset of the words in a document
  - $\text{words\_in\_doc}_d[w]$  is the number of occurrences of word  $w$  in document  $d$
  - To get this vector, generate the words one-at-a-time
  - For a given word in doc  $d$ :
    - (1) Figure out the topic  $t$  that controls it by sampling from a Multinomial ( $\text{topic\_probs}_d, 1$ ) distribution
    - (2) Generate the word by sampling from a Multinomial ( $\text{word\_probs}_t, 1$ ) distribution
- Ex: doc 0...  $\text{topic\_probs}_0 = (.98, 0.01, 0.01)$  “I can’t stand”
  - Onto next word

# LDA Step Three

- Generate the words in each document
  - Each topic “controls” a subset of the words in a document
  - $\text{words\_in\_doc}_d[w]$  is the number of occurrences of word  $w$  in document  $d$
  - To get this vector, generate the words one-at-a-time
  - For a given word in doc  $d$ :
    - (1) Figure out the topic  $t$  that controls it by sampling from a Multinomial ( $\text{topic\_probs}_d, 1$ ) distribution
    - (2) Generate the word by sampling from a Multinomial ( $\text{word\_probs}_t, 1$ ) distribution
- Ex: doc 0...  $\text{topic\_probs}_0 = (.98, 0.01, 0.01)$  “I can’t stand”
  - $t$  for word three is zero, since we sampled  $(1, 0, 0)$  [there is a 1 in the zeroth entry]

# LDA Step Three

- Generate the words in each document
  - Each topic “controls” a subset of the words in a document
  - $\text{words\_in\_doc}_d[w]$  is the number of occurrences of word  $w$  in document  $d$
  - To get this vector, generate the words one-at-a-time
  - For a given word in doc  $d$ :
    - (1) Figure out the topic  $t$  that controls it by sampling from a Multinomial ( $\text{topic\_probs}_d, 1$ ) distribution
    - (2) Generate the word by sampling from a Multinomial ( $\text{word\_probs}_t, 1$ ) distribution
- Ex: doc 0...  $\text{topic\_probs}_0 = (.98, 0.01, 0.01)$  “I can’t stand bad”
  - $t$  for word three is zero, since we sampled  $(1, 0, 0)$  [there is a 1 in the zeroth entry]
  - So we generate the word using  $\text{word\_probs}_0 = (.2, .2, .2, .2, 0, 0, 0, 0, .2, 0, 0)$
  - And we get  $(1, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0)$ , which is equivalent to “bad”

# LDA Step Three

- Generate the words in each document
  - Each topic “controls” a subset of the words in a document
  - $\text{words\_in\_doc}_d[w]$  is the number of occurrences of word  $w$  in document  $d$
  - To get this vector, generate the words one-at-a-time
  - For a given word in doc  $d$ :
    - (1) Figure out the topic  $t$  that controls it by sampling from a Multinomial ( $\text{topic\_probs}_d, 1$ ) distribution
    - (2) Generate the word by sampling from a Multinomial ( $\text{word\_probs}_t, 1$ ) distribution
- Ex: doc 0...  $\text{topic\_probs}_0 = (.98, 0.01, 0.01)$  “I can’t stand bad”
  - Onto the last word in the document

# LDA Step Three

- Generate the words in each document
  - Each topic “controls” a subset of the words in a document
  - $\text{words\_in\_doc}_d[w]$  is the number of occurrences of word  $w$  in document  $d$
  - To get this vector, generate the words one-at-a-time
  - For a given word in doc  $d$ :
    - (1) Figure out the topic  $t$  that controls it by sampling from a Multinomial ( $\text{topic\_probs}_d, 1$ ) distribution
    - (2) Generate the word by sampling from a Multinomial ( $\text{word\_probs}_t, 1$ ) distribution
- Ex: doc 0...  $\text{topic\_probs}_0 = (.98, 0.01, 0.01)$  “I can’t stand bad”
  - $t$  for word three is zero, since we sampled  $(1, 0, 0)$  [there is a 1 in the zeroth entry]

# LDA Step Three

- Generate the words in each document
  - Each topic “controls” a subset of the words in a document
  - $\text{words\_in\_doc}_d[w]$  is the number of occurrences of word  $w$  in document  $d$
  - To get this vector, generate the words one-at-a-time
  - For a given word in doc  $d$ :
    - (1) Figure out the topic  $t$  that controls it by sampling from a Multinomial ( $\text{topic\_probs}_d, 1$ ) distribution
    - (2) Generate the word by sampling from a Multinomial ( $\text{word\_probs}_t, 1$ ) distribution
- Ex: doc 0...  $\text{topic\_probs}_0 = (.98, 0.01, 0.01)$  “I can’t stand bad beer”
  - $t$  for word three is zero, since we sampled  $(1, 0, 0)$  [there is a 1 in the zeroth entry]
  - So we generate the word using  $\text{word\_probs}_0 = (.2, .2, .2, .2, 0, 0, 0, 0, .2, 0, 0)$
  - And we get  $(0, 0, 0, 0, 0, 0, 0, 0, 1, 0, 0)$ , which is equivalent to “beer”

## In The End... For Doc 0...

- text is “I can’t stand bad beer” (equiv. to “1 2 3 0 8”)
- $\text{topic\_probs}_0 = (.98, 0.01, 0.01)$
- $\text{words\_in\_doc}_0 = (1, 1, 1, 1, 0, 0, 0, 0, 1, 0, 0)$ 
  - Why? Word 0 appears once, word 1 appears once, word 4 zero times, etc.
- $\text{produced}_0 = (1, 1, 1, 1, 0, 0, 0, 0, 1, 0, 0)$ 
  - $(0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0)$
  - $(0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0)$
  - Why? Topic 0 (associated with first line) produced 5 words  
Those words were  $(1, 1, 1, 1, 0, 0, 0, 0, 1, 0, 0)$
  - Topic 1, topic 2 produced no words
  - “produced” always a matrix with  $\text{num\_words}$  cols,  $k$  rows



Repeat For Each Doc in the Corpus!

## For Example, Let's Look At Doc 2...

- $\text{topic\_probs}_2 = (.02, 0.49, 0.49)$
- Imagine that when we generate doc 2, we get:
  - Word 0: produced by topic 2, is 1 or “I”
  - Word 1: produced by topic 2, is 7 or “love”
  - Word 2: produced by topic 2, is 8 or “beer”
  - Word 3: produced by topic 1, is 1 or “I”
  - Word 4: produced by topic 1, is 2 or “can’t”
  - Word 5: produced by topic 2, is 7 or “love”
  - Word 6: produced by topic 1, is 9 or “humanities”
  - Word 7: produced by topic 1, is 10 or “classes”
- $\text{words\_in\_doc}_2 = (0, 2, 1, 0, 0, 0, 0, 2, 1, 1, 1)$
- $\text{produced}_2 = \begin{pmatrix} 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0 \\ 0, 1, 1, 0, 0, 0, 0, 0, 0, 1, 1 \\ 0, 1, 0, 0, 0, 0, 0, 2, 1, 0, 0 \end{pmatrix}$

# OK We've Got Our Generative Process

- Now, someone gives us an actual, real-life corpus
  - This means we have got a dictionary
  - And we have  $\text{words\_in\_doc}_d$  for all  $d$  in  $\{0 \dots n - 1\}$
- Our goal is to figure out a posterior distribution on
  - $\text{word\_probs}_t$  for all  $t$  in  $\{0 \dots k - 1\}$
  - $\text{topic\_probs}_d$  for all  $d$  in  $\{0 \dots n - 1\}$
  - $\text{produced}_d$  for all  $d$  in  $\{0 \dots n - 1\}$
- Why?
  - The hope is this will reveal something about the corpus
  - For example, what the different topics in the corpus are
  - And what topics are present in each document

# As In All Applications of Bayesian ML

- The posterior is derived using Bayes' Theorem...

$$p(\text{word\_probs}_{all}, \text{topic\_probs}_{all}, \text{produced}_{all} \mid \text{words\_in\_doc}_{all}) = \\ p(\text{words\_in\_doc}_{all} \mid \text{word\_probs}_{all}, \text{topic\_probs}_{all}, \text{produced}_{all}) \times \\ p(\text{word\_probs}_{all}, \text{topic\_probs}_{all}, \text{produced}_{all}) / p(\text{words\_in\_doc}_{all})$$

# We Can Write This Formula

- But what does this do for us?
- In the simple “guess the test score average” case...
  - We had a couple of posterior distributions that we could plot
  - It was quite intuitive
- Not the case here!
  - Got all kinds of weird, multi-dimensions distributions
  - How to proceed?

# With More Complicated Models Such as LDA

- It is common to resort to something called “MCMC”
  - stands for “Markov Chain Monte Carlo”
- MCMC is a class of algorithms
  - That can often be used to draw samples from even the most complex distributions
  - Many of the ideas behind MCMC came out the the Manhattan project
- Using MCMC, we can actually draw samples from
$$p(\text{word\_probs}_{all}, \text{topic\_probs}_{all}, \text{produced}_{all} \mid \text{words\_in\_doc}_{all})$$
- Each “sample” will contain everything in the model!
  - $\text{word\_probs}_t$  for all  $t$  in  $\{0 \dots k - 1\}$
  - $\text{topic\_probs}_d$  for all  $d$  in  $\{0 \dots n - 1\}$
  - $\text{produced}_d$  for all  $d$  in  $\{0 \dots n - 1\}$

# Why Useful?

- Can take one sample, use it as a representative set of values
- Or take many samples, use to get a summary of the distribution

# How Does MCMC Work?

- Could spend an entire semester on MCMC alone
- Many flavors of MCMC
- For LDA, we'll employ a simple MCMC methodology
  - A “Gibbs Sampler”



# In the Case of LDA

- Here is pseudo-code for the Gibbs sampler:

Choose consistent, initial values for  $\text{word\_probs}_{all}$ ,  $\text{topic\_probs}_{all}$ ,  $\text{produced}_{all}$

For  $i = 1$  to  $\text{num\_iters}$  do:

For all  $t$  in  $\{0 \dots k - 1\}$ :

$\text{word\_probs}_t \sim p(\text{word\_probs}_t \mid \text{topic\_probs}_{all}, \text{produced}_{all}, \text{words\_in\_doc}_{all})$

For all  $d$  in  $\{0 \dots n - 1\}$ :

$\text{topic\_probs}_d \sim p(\text{topic\_probs}_d \mid \text{word\_probs}_{all}, \text{produced}_{all}, \text{words\_in\_doc}_{all})$

For all  $d$  in  $\{0 \dots n - 1\}$ :

$\text{produced}_d \sim p(\text{produced}_d \mid \text{word\_probs}_{all}, \text{topic\_probs}_{all}, \text{words\_in\_doc}_{all})$

- Run this loop for awhile

— Then at any point, stop

— Current  $\text{word\_probs}_{all}$ ,  $\text{topic\_probs}_{all}$ ,  $\text{produced}_{all}$  are a sample from posterior!

## So All That's Left

- Is to give pseudo-code for each of the three steps

# Word\_Probs

- To sample each  $\text{word\_probs}_t$ 
  - Create a vector called “counter”, where counter is  $\sum_d \text{produced}_{d,t}$
  - Note that  $\text{produced}_{d,t}$  denotes the  $t$ th row of the produced matrix for doc  $d$
  - This counts the number of times topic  $t$  produced each word  $w$
  - Then  $\text{word\_probs}_t \sim \text{Dirichlet}(\text{counter} + \text{alpha})$

↑  
Constant used in generative process  
We assume we know this,  
or can guess it

# Word\_Probs

- To sample each  $\text{word\_probs}_t$ 
    - Create a vector called “counter”, where counter is  $\sum_d \text{produced}_{d,t}$
    - Note that  $\text{produced}_{d,t}$  denotes the  $t$ th row of the produced matrix for doc  $d$
    - This counts the number of times topic  $t$  produced each word  $w$
    - Then  $\text{word\_probs}_t \sim \text{Dirichlet}(\text{counter} + \text{alpha})$
  - $\text{produced}_0 = (0, 0, 0, 0, 2, 0, 0, 0, 0, 0, 0)$   
 $(0, 1, 0, 0, 0, 0, 1, 0, 0, 1, 1)$   
 $(0, 1, 1, 0, 0, 0, 0, 0, 0, 0, 1)$
  - $\text{produced}_1 = (0, 1, 0, 0, 0, 0, 0, 0, 0, 0, 0)$   
 $(0, 0, 1, 0, 1, 0, 0, 2, 0, 1, 0)$   
 $(0, 1, 1, 0, 0, 1, 0, 0, 0, 1, 1)$
  - $\text{produced}_2 = (0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0)$   
 $(0, 1, 1, 0, 0, 0, 0, 0, 0, 0, 1)$   
 $(0, 1, 0, 0, 3, 0, 0, 2, 1, 0, 0)$
- For  $t = 1$ , counter is  $(0, 2, 2, 0, 1, 0, 1, 2, 0, 2, 2)$

# Word\_Probs

- To sample each  $\text{word\_probs}_t$ 
  - Create a vector called “counter”, where counter is  $\sum_d \text{produced}_{d,t}$
  - Note that  $\text{produced}_{d,t}$  denotes the  $t$ th row of the produced matrix for doc  $d$
  - This counts the number of times topic  $t$  produced each word  $w$
  - Then  $\text{word\_probs}_t \sim \text{Dirichlet}(\text{counter} + \text{alpha})$
- Intuitively
  - You are “guessing” the probability topic  $t$  will produce each word
  - If  $\text{counter}[w]$  is large, then topic  $t$  produced  $w$  quite often in practice...
  - And the Dirichlet is then likely to give that word a high probability

# Topic\_Probs

- To sample each  $\text{topic\_probs}_d$ 
  - Create a vector “counter”, where  $\text{counter}[t]$  is  $\sum_w \text{produced}_{d,t}[w]$
  - That is,  $\text{counter}[t]$  is the sum over the  $t$ th row in the  $\text{produced}_d$  matrix
  - This counts the number of times topic  $t$  was used to produce a word in  $d$
  - Then  $\text{topic\_probs}_d \sim \text{Dirichlet}(\text{counter} + \text{beta})$

↑  
Again, constant used in generative process  
We assume we know this,  
or can guess it

# Topic\_Probs

- To sample each  $\text{topic\_probs}_d$ 
  - Create a vector “counter”, where  $\text{counter}[t]$  is  $\sum_w \text{produced}_{d,t}[w]$
  - That is,  $\text{counter}[t]$  is the sum over the  $t$ th row in the  $\text{produced}_d$  matrix
  - This counts the number of times topic  $t$  was used to produce a word in  $d$
  - Then  $\text{topic\_probs}_d \sim \text{Dirichlet}(\text{counter} + \text{beta})$
- $\text{produced}_1 =$ 
  - (0, 1, 0, 0, 0, 0, 0, 0, 0, 0, 0)
  - (0, 0, 1, 0, 1, 0, 0, 2, 0, 1, 0)
  - (0, 1, 1, 0, 0, 1, 0, 0, 0, 1, 1)

For  $w = 6$ , counter is (1, 5, 5)

# Topic\_Probs

- To sample each  $\text{topic\_probs}_d$ 
  - Create a vector “counter”, where  $\text{counter}[t]$  is  $\sum_w \text{produced}_{d,t}[w]$
  - That is,  $\text{counter}[t]$  is the sum over the  $t$ th row in the  $\text{produced}_d$  matrix
  - This counts the number of times topic  $t$  was used to produce a word in  $d$
  - Then  $\text{topic\_probs}_d \sim \text{Dirichlet}(\text{counter} + \text{beta})$
- Intuitively
  - You are “guessing” the probability a word in  $d$  will come from each topic
  - If  $\text{counter}[t]$  is large, then topic  $t$  produced a lot of words in  $d$ ...
  - And the Dirichlet is then likely to give that topic a high probability



# Produced

- To sample each produced<sub>*d*</sub>
  - For each word *w*:
    - (1) Create a vector “probs”, where  $\text{probs}[t] = \text{word\_probs}_t[w] \times \text{topic\_probs}_d[t]$
    - (2) Normalize probs
    - (3) Then (*w*th column in produced<sub>*d*</sub>)  $\sim$  Multinomial (words\_in\_doc<sub>*d*</sub>[*w*], probs)

# Produced

- To sample each produced<sub>d</sub>

— For each word  $w$ :

(1) Create a vector “probs”, where  $\text{probs}[t] = \text{word\_probs}_t[w] \times \text{topic\_probs}_d[t]$

(2) Normalize probs

(3) Then ( $w$ th column in produced<sub>d</sub>)  $\sim$  Multinomial (words\_in\_doc<sub>d</sub>[ $w$ ], probs)

Ex:

word\_probs<sub>0</sub> = (.2, .2, .2, .2, 0, 0, 0, 0, .2, 0, 0)

word\_probs<sub>1</sub> = (0, .2, .2, .2, 0, 0, 0, 0, 0, .2, .2)

word\_probs<sub>2</sub> = (0, .2, .2, 0, 0, 0, .2, .2, 0, 0, 0)

topic\_probs<sub>0</sub> = (.98, 0.01, 0.01)

Let  $w = 4$ ... then  $\text{probs} = \text{normalize}((0, 0, 0.2 \times 0.01)) = (0, 0, 1)$

# Produced

- To sample each produced  $d$ 
  - For each word  $w$ :
    - (1) Create a vector “probs”, where  $\text{probs}[t] = \text{word\_probs}_t[w] \times \text{topic\_probs}_d[t]$
    - (2) Normalize probs
    - (3) Then (with column in produced  $d$ )  $\sim$  Multinomial ( $\text{words\_in\_doc}_d[w]$ , probs)
- Intuitively
  - “probs” is telling you how likely each topic is to be the one responsible for an occurrence of  $w$  in  $d$
  - Then you are using the multinomial to guess how many occurs of  $w$  each topic is responsible for
  - $\text{probs}[t]$  large means  $t$  expectedly responsible for more occurs of  $w$

This Completes LDA. Questions?