

Case Study 2: Document Retrieval

MAP EM, Latent Dirichlet Allocation, Gibbs Sampling

Machine Learning/Statistics for Big Data
CSE599C1/STAT592, University of Washington

Emily Fox

February 5th, 2013

©Emily Fox 2013

1

Gaussian Mixture Model

- Most commonly used mixture model
- Observations: x^1, \dots, x^N

- Parameters: $\theta = \{\pi, \phi\}$

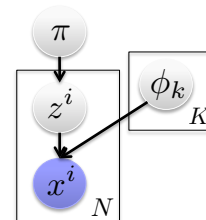
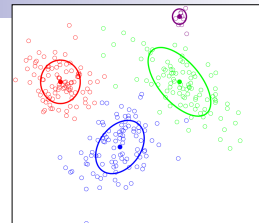
$$\pi = [\pi_1, \dots, \pi_K]$$

$$\phi = \{\phi_k\} = \{\mu_k, \Sigma_k\}$$

- Likelihood:

$$p(x^i | \theta) = \sum_k \pi_k p(x^i | \phi_k)$$

- Ex. z^i = country of origin, x^i = height of i^{th} person
 - k^{th} mixture component = distribution of heights in country k



©Emily Fox 2013

2

Motivates EM Algorithm

- Initial guess: $\hat{\theta}^{(0)}$
- Estimate at iteration t : $\hat{\theta}^{(t)}$
- **E-Step**
Compute $U(\theta, \hat{\theta}^{(t)}) = E[\log p(y | \theta) | x, \hat{\theta}^{(t)}]$
- **M-Step**
Compute $\hat{\theta}^{(t+1)} = \arg \max_{\theta} U(\theta, \hat{\theta}^{(t)})$

©Emily Fox 2013

3

MAP Estimation

- Bayesian approach:
 - Place **prior** $p(\theta)$ on parameters
 - Infer **posterior** $p(\theta | x)$
- Many, many, many motivations and implications
 - For the sake of this class, simplest motivation is to think of this as akin to regularization

$$\hat{\theta}^{MAP} = \arg \max_{\theta} \log p(\theta | x)$$

- Saw importance of regularization in logistic regression (ML estimate can overfit data and lead to poor generalization)

©Emily Fox 2013

4

EM Algorithm – MAP Case

- Re-derive EM algorithm for $p(\theta | x)$
- Add $\log p(\theta)$ to $U(\theta, \hat{\theta}^{(t)})$
 - What must be computed in E-Step remains unchanged because this term does not depend on y .
 - M-Step becomes:

$$\hat{\theta}^{(t+1)} = \arg \max_{\theta} U(\theta, \hat{\theta}^{(t)})$$

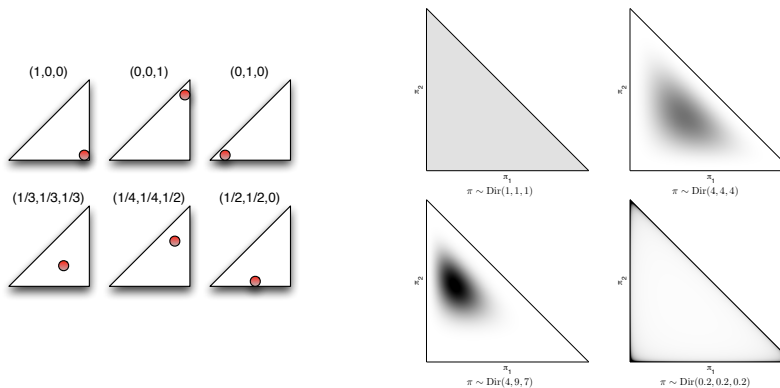
©Emily Fox 2013

5

MAP EM Example – MoG

- For mixture of Gaussians, conjugate priors are:

$$\pi \sim \text{Dir}(\alpha_1, \dots, \alpha_K)$$



©Emily Fox 2013

6

MAP EM Example – MoG

- For mixture of Gaussians, conjugate priors are:

$$\pi \sim \text{Dir}(\alpha_1, \dots, \alpha_K) \quad p(\pi | \alpha) = \frac{\Gamma(\sum_k \alpha_k)}{\prod_k \Gamma(\alpha_k)} \prod_k \pi_k^{\alpha_k - 1}$$

- Dirichlet posterior

- Assume we condition on observations $z^i \sim \pi$
- Count occurrences of $z^i = k$
- Then,

$$p(\pi | \alpha, z^1, \dots, z^N) \propto$$

- Conjugacy: This **posterior** has same form as **prior**

©Emily Fox 2013

7

MAP EM Example – MoG

- For mixture of Gaussians, conjugate priors are:

$$\pi \sim \text{Dir}(\alpha_1, \dots, \alpha_K) \quad \{\mu_k, \Sigma_k\} \sim \text{NIW}(m_0, \kappa_0, \nu_0, S_0)$$

- Results in following M-Step:

$$\hat{\mu}_k = \frac{r_k \bar{x}_k + \kappa_0 m_0}{r_k + \kappa_0} \quad \hat{\pi}_k = \frac{r_k + \alpha_k - 1}{N + \sum_k \alpha_k - K}$$

$$\hat{\Sigma}_k = \frac{S_0 + r_k S_k + \frac{\kappa_0 r_k}{\kappa_0 + r_k} (\bar{x}_k - m_0)(\bar{x}_k - m_0)'}{\nu_0 + r_k + d + 2}$$

©Emily Fox 2013

8

Posterior Computations

- MAP EM focuses on point estimation:

$$\hat{\theta}^{MAP} = \arg \max_{\theta} p(\theta | x)$$

- What if we want a full characterization of the posterior?
 - Maintain a measure of uncertainty
 - Estimators other than posterior mode (different loss functions)
 - Predictive distributions for future observations
- Often no closed-form characterization (e.g., mixture models)
- Alternatives:
 - Monte Carlo based estimates using samples from posterior
 - Variational approximations to posterior (more next time)

©Emily Fox 2013

9

Gibb Sampling

- Want draws:
 - Construct Markov chain whose steady state distribution is
 - Simplest case:

©Emily Fox 2013

10

Example – Mixture of Gaussians

- Recall model

- Observations: x^1, \dots, x^N
- Cluster indicators: z^1, \dots, z^N
- Parameters: $\theta = \{\pi, \phi\}$ $\pi = [\pi_1, \dots, \pi_K]$
 $\phi = \{\phi_k\} = \{\mu_k, \Sigma_k\}$

- Generative model:

$$\pi \sim \text{Dir}(\alpha_1, \dots, \alpha_K) \qquad z^i \sim \pi$$

$$\{\mu_k, \Sigma_k\} \sim F(\phi) \qquad x^i | z^i \sim N(x^i; \mu_{z^i}, \Sigma_{z^i})$$

- Want to draw posterior samples of model parameters

$$\pi \sim p(\pi | \phi, x^1, \dots, x^N)$$

$$\phi \sim p(\phi | \pi, x^1, \dots, x^N)$$

©Emily Fox 2013

11

Auxiliary Variable Samplers

- Augment variables of interest θ with variables z to allow closed-form for sampling, just like in EM

- In both cases, simply looking at subchain $\{\theta^{(t)}\}$ converges to draws from marginal distribution $\pi(\theta)$

©Emily Fox 2013

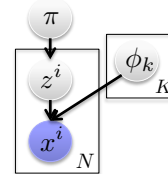
12

Example – Mixture of Gaussians

$$\pi \sim \text{Dir}(\alpha_1, \dots, \alpha_K) \quad z^i \sim \pi$$

$$\{\mu_k, \Sigma_k\} \sim F(\phi) \quad x^i | z^i \sim N(x^i; \mu_{z^i}, \Sigma_{z^i})$$

- Try auxiliary variable sampler
 - Introduce cluster indicators into sampler



©Emily Fox 2013

13

Example – Clustering Results I

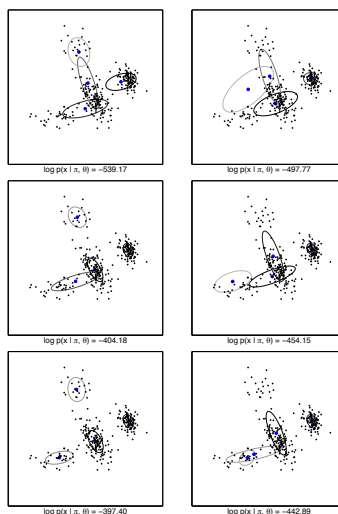


Figure courtesy of Erik Sudderth

©Emily Fox 2013

14

Collapsed Gibbs Samplers

- Marginalize a set of latent variables or parameters
 - Sometimes marginalized variables are nuisance parameters
 - Other times what gets marginalized are the variables
 - Make post-facto inferences on variables of interest based on sampled variables

- Can improve efficiency if marginalized variables are high-dim
 - Reduced dimension of search space
 - But, often introduces dependences!

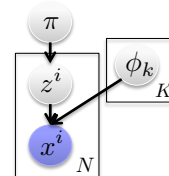
©Emily Fox 2013

15

Example – Collapsed MoG Sampling

$$\pi \sim \text{Dir}(\alpha_1, \dots, \alpha_K) \quad z^i \sim \pi$$
$$\{\mu_k, \Sigma_k\} \sim F(\phi) \quad x^i | z^i \sim N(x^i; \mu_{z^i}, \Sigma_{z^i})$$

- Collapsed sampler



©Emily Fox 2013

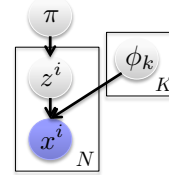
16

Example – Collapsed MoG Sampling

$$\pi \sim \text{Dir}(\alpha_1, \dots, \alpha_K) \quad z^i \sim \pi$$

$$\{\mu_k, \Sigma_k\} \sim F(\phi) \quad x^i | z^i \sim N(x^i; \mu_{z^i}, \Sigma_{z^i})$$

- Derivation



- Important facts:

$$p(z_{1:N} | \alpha) = \frac{\Gamma(\sum_k \alpha_k) \prod_k \Gamma(n_k + \alpha_k)}{\prod_k \Gamma(\alpha_k) \Gamma(\sum_k n_k + \alpha_k)} \quad \frac{\Gamma(m+1)}{\Gamma(m)} = m$$

©Emily Fox 2013

17

Example – Clustering Results II

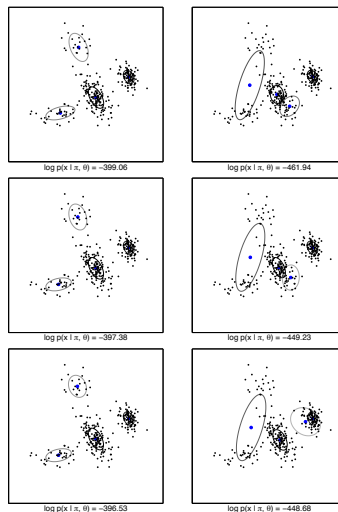


Figure courtesy of Erik Sudderth

©Emily Fox 2013

18

Comparing Collapsed vs. Regular

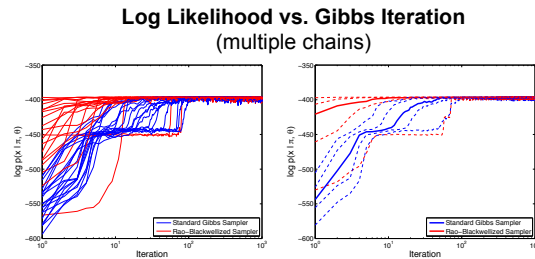


Figure courtesy of Erik Sudderth

©Emily Fox 2013

19

Task 2: Cluster Documents

■ Previously:

- Cluster documents based on topic

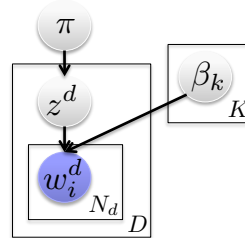


©Emily Fox 2013

20

A Generative Model

- Documents: x^1, \dots, x^D
- Associated topics: z^1, \dots, z^D
- Parameters: $\theta = \{\pi, \beta\}$
- Generative model:

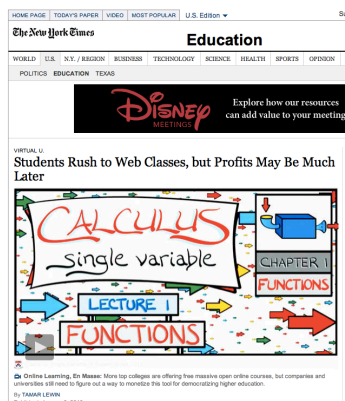


©Emily Fox 2013

21

Task 2: Cluster Documents

- **Now:** Document may belong to multiple clusters



EDUCATION

FINANCE

TECHNOLOGY

©Emily Fox 2013

22

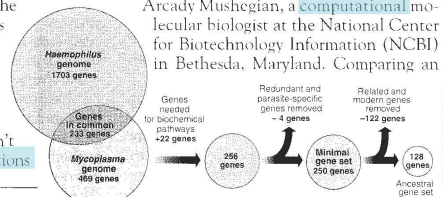
Latent Dirichlet Allocation (LDA)

COLD SPRING HARBOR, NEW YORK—How many genes does an organism need to survive? Last week at the genome meeting here,* two genome researchers with radically different approaches presented complementary views of the basic genes needed for life. One research team, using computer analyses to compare known genomes, concluded that today's organisms can be sustained with just 250 genes, and that the earliest life forms required a mere 128 genes. The other researcher mapped genes in a simple parasite and estimated that for this organism, 800 genes are plenty to do the job—but that anything short of 100 wouldn't be enough.

Although the numbers don't match precisely, those predictions

* Genome Mapping and Sequencing, Cold Spring Harbor, New York, May 8 to 12.

"are not all that far apart," especially in comparison to the 75,000 genes in the human genome, notes Siv Andersson of Uppsala University in Sweden, who arrived at the 800 number. But coming up with a consensus answer may be more than just a genetic numbers game, particularly as more and more genomes are completely mapped and sequenced. "It may be a way of organizing any newly sequenced genome," explains Arcady Mushegian, a computational molecular biologist at the National Center for Biotechnology Information (NCBI) in Bethesda, Maryland. Comparing an



Stripping down. Computer analysis yields an estimate of the minimum modern and ancient genomes.

Latent Dirichlet Allocation (LDA)

Topics

gene	0.94
dna	0.92
genetic	0.01
...	
life	0.62
evolve	0.01
organism	0.01
...	
brain	0.04
neuron	0.02
nerve	0.01
...	
data	0.02
number	0.02
computer	0.01
...	

Documents

Topic proportions and assignments

Seeking Life's Bare (Genetic) Necessities

COLD SPRING HARBOR, NEW YORK—How many genes does an organism need to survive? Last week at the genome meeting here,* two genome researchers with radically different approaches presented complementary views of the basic genes needed for life. One research team, using computer analyses to compare known genomes, concluded that today's organisms can be sustained with just 250 genes, and that the earliest life forms required a mere 128 genes. The other researcher mapped genes in a simple parasite and estimated that for this organism, 800 genes are plenty to do the job—but that anything short of 100 wouldn't be enough. Although the numbers don't match precisely, those predictions

"are not all that far apart," especially in comparison to the 75,000 genes in the human genome, notes Siv Andersson of Uppsala University in Sweden, who arrived at the 800 number. But coming up with a consensus answer may be more than just a genetic numbers game, particularly as more and more genomes are completely mapped and sequenced. "It may be a way of organizing any newly sequenced genome," explains Arcady Mushegian, a computational molecular biologist at the National Center for Biotechnology Information (NCBI) in Bethesda, Maryland. Comparing an

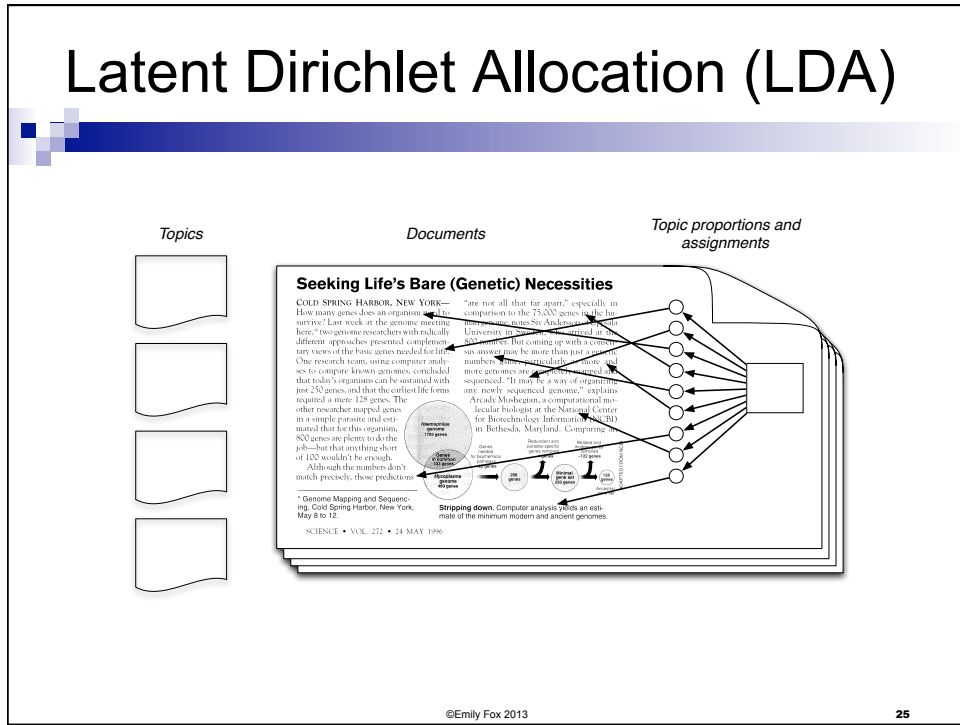
ADAPTED FROM NCBI

* Genome Mapping and Sequencing, Cold Spring Harbor, New York, May 8 to 12.

Stripping down. Computer analysis yields an estimate of the minimum modern and ancient genomes.

SCIENCE • VOL. 311 • 24 MAY 1996

Latent Dirichlet Allocation (LDA)



©Emily Fox 2013

25

Example Inference – Topic Weights

- **Data:** The OCR'd collection of *Science* from 1990-2000
 - 17K documents
 - 11M words
 - 20K unique terms (stop words and rare words removed)
- **Model:** 100-topic LDA model

Seeking Life's Bare (Genetic) Necessities

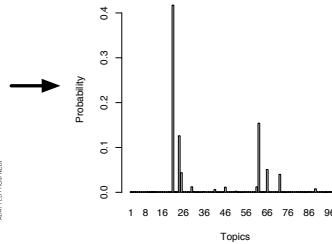
COLD SPRING HARBOR, NEW YORK—How many genes does an organism need to survive? Last week at the genome meeting here,* two genome researchers with radically different approaches presented complementary views of the basic genes needed for life. One research team, using computer analysis to compare known genomes, concluded that today's organisms can be sustained with just 250 genes, and that the earliest life forms required a mere 128 genes. The other researcher mapped genes in a simple parasite and estimated that for this organism, 800 genes are plenty to do the job—but that anything short of 100 wouldn't be enough. Although the numbers don't match precisely, those predictions

"are not all that far apart," especially in comparison to the 25,000 genes in the human genome, notes Siv Andersson of Uppsala University in Sweden, who arrived at the 800 number. But coming up with a consensus answer may be more than just a genetic numbers game, particularly as more and more genomes are completely mapped and sequenced. "It may be a way of organizing any newly sequenced genome," explains Arcady Mushegian, a computational molecular biologist at the National Center for Biotechnology Information (NCBI) in Bethesda, Maryland. Comparing an

Stripping down. Computer analysis yields an estimate of the minimum modern and ancient genomes.

* Genome Mapping and Sequencing, Cold Spring Harbor, New York, May 8 to 12.

SCIENCE • VOL. 272 • 24 MAY 1996



©Emily Fox 2013

26

Example Inference – Topic Words

human	evolution	disease	computer
genome	evolutionary	host	models
dna	species	bacteria	information
genetic	organisms	diseases	data
genes	life	resistance	computers
sequence	origin	bacterial	system
gene	biology	new	network
molecular	groups	strains	systems
sequencing	phylogenetic	control	model
map	living	infectious	parallel
information	diversity	malaria	methods
genetics	group	parasite	networks
mapping	new	parasites	software
project	two	united	new
sequences	common	tuberculosis	simulations

©Emily Fox 2013

27

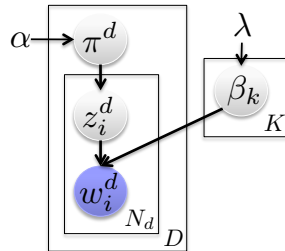
LDA Generative Model

- Observations: $w_1^d, \dots, w_{N_d}^d$
- Associated topics: $z_1^d, \dots, z_{N_d}^d$
- Parameters: $\theta = \{\{\pi^d\}, \{\beta_k\}\}$
- Generative model:

©Emily Fox 2013

28

LDA Generative Model



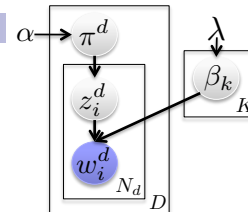
$$p(\cdot) = \prod_{k=1}^K p(\beta_k | \lambda) \prod_{d=1}^D p(\pi^d | \alpha) \left(\prod_{i=1}^{N_d} p(z_i^d | \pi^d) p(w_i^d | z_i^d, \beta) \right)$$

©Emily Fox 2013

29

Collapsed LDA Sampling

- Marginalize parameters
 - Document-specific topic weights
 - Corpus-wide topic-specific word distributions
- Sample topic indicators for each word
 - Derivation:



$$p(z_{1:N_d}^d | \alpha) = \frac{\Gamma(\sum_k \alpha_k) \prod_k \Gamma(n_k^d + \alpha_k)}{\prod_k \Gamma(\alpha_k) \Gamma(\sum_k n_k^d + \alpha_k)}$$

$$p(\{w_i^d | z_i^d = k\}, \lambda) = \frac{\Gamma(\sum_\nu \lambda_\nu) \prod_\nu \Gamma(v_\nu^k + \lambda_\nu)}{\prod_\nu \Gamma(\lambda_\nu) \Gamma(\sum_\nu v_\nu^k + \lambda_\nu)}$$

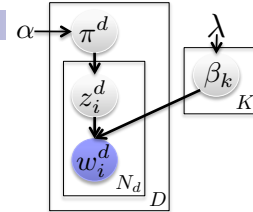
$$p(z | \alpha) = \prod_{d=1}^D p(z_{1:N_d}^d | \alpha) \quad p(w | z, \lambda) = \prod_{k=1}^K p(\{w_i^d | z_i^d = k\}, \lambda)$$

©Emily Fox 2013

30

Collapsed LDA Sampling

- Marginalize parameters
 - Document-specific topic weights
 - Corpus-wide topic-specific word distributions
- Sample topic indicators for each word
 - Algorithm:



©Emily Fox 2013

31

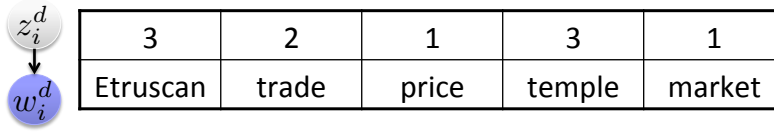
Sample Document

Etruscan	trade	price	temple	market

©Emily Fox 2013

32

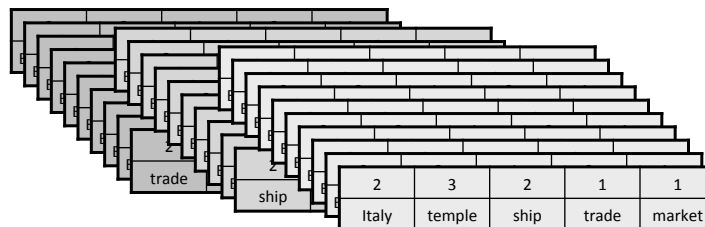
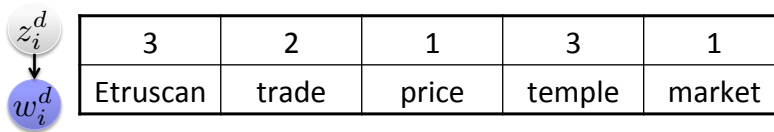
Randomly Assign Topics



©Emily Fox 2013

33

Randomly Assign Topics



©Emily Fox 2013

34

Maintain Global Statistics

z_i^d	3	2	1	3	1
w_i^d	Etruscan	trade	price	temple	market

Total counts from **all** docs

	1	2	3
Etruscan	1	0	35
market	50	0	1
price	42	1	0
temple	0	0	20
trade	10	8	1
...			

©Emily Fox 2013

35

Resample Assignments

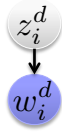
z_i^d	3	2	1	3	1
w_i^d	Etruscan	trade	price	temple	market

	1	2	3
Etruscan	1	0	35
market	50	0	1
price	42	1	0
temple	0	0	20
trade	10	8	1
...			

©Emily Fox 2013

36

What is the conditional distribution for this topic?




3	?	1	3	1
Etruscan	trade	price	temple	market

©Emily Fox 2013

37

What is the conditional distribution for this topic?

- Part I: How much does this document like each topic?



3	?	1	3	1
Etruscan	trade	price	temple	market

Topic 1

Topic 2

Topic 3



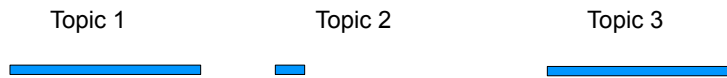
©Emily Fox 2013

38

What is the conditional distribution for this topic?

- Part I: How much does this document like each topic?
- Part II: How much does each topic like this word?

z_i^d	3	?	1	3	1
w_i^d	Etruscan	trade	price	temple	market



	1	2	3
trade	10	7	1

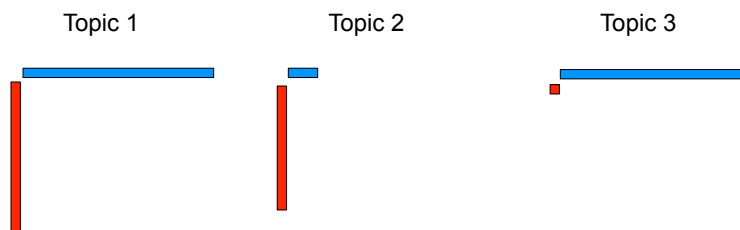
©Emily Fox 2013

39

What is the conditional distribution for this topic?

- Part I: How much does this document like each topic?
- Part II: How much does each topic like this word?

z_i^d	3	?	1	3	1
w_i^d	Etruscan	trade	price	temple	market



©Emily Fox 2013

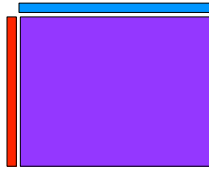
40

What is the conditional distribution for this topic?

- Part I: How much does this document like each topic?
- Part II: How much does each topic like this word?

z_i^d	3	?	1	3	1
w_i^d	Etruscan	trade	price	temple	market

Topic 1



Topic 2



Topic 3



$$\frac{n_k^d + \alpha_k}{\sum_{j=1}^K n_j^d + \alpha_j} \frac{v_{trade}^k + \lambda_k}{\sum_{j=1}^V v_j^k + \lambda_j}$$

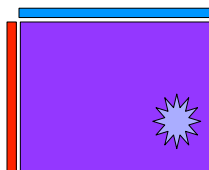
©Emily Fox 2013

41

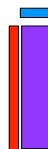
Sample a New Topic Indicator

z_i^d	3	?	1	3	1
w_i^d	Etruscan	trade	price	temple	market

Topic 1



Topic 2



Topic 3



©Emily Fox 2013

42

Update Counts

Diagram illustrating the relationship between document variables z_i^d and w_i^d and their corresponding word counts in a document.

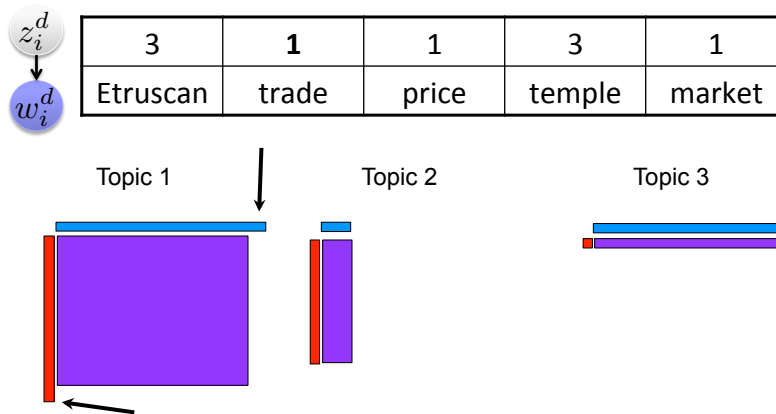
z_i^d	3	?	1	3	1
w_i^d	Etruscan	trade	price	temple	market

	1	2	3
Etruscan	1	0	35
market	50	0	1
price	42	1	0
temple	0	0	20
trade	10	7	1
...			

©Emily Fox 2013

43

Geometrically...



©Emily Fox 2013

44

Issues with Generic LDA Sampling

- Slow mixing rates → Need many iterations
- Each iteration cycles through sampling topic assignments for *all* words in *all* documents
- Modern approaches:
 - Large-scale LDA. For example, [Mimno, David, Matthew D. Hoffman and David M. Blei. "Sparse stochastic inference for latent Dirichlet allocation." International Conference on Machine Learning, 2012.](#)
 - Distributed LDA. For example, [Ahmed, Amr, et al. "Scalable inference in latent variable models." Proceedings of the fifth ACM international conference on Web search and data mining \(2012\): 123-132](#)
- Next time: Variational methods instead of sampling

Acknowledgements

- Thanks to Dave Blei, David Mimno, and Jordan Boyd-Graber for some material in this lecture relating to LDA