

Gaussian Mixture Model

- Most commonly used mixture model

- Observations: x^1, \dots, x^N

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

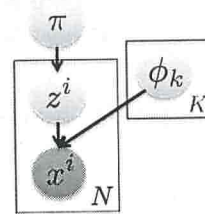
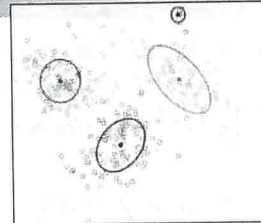
$$\pi = [\pi_1, \dots, \pi_K] \quad \text{mixture weights}$$

$$\phi = \{\phi_k\} = \{\mu_k, \Sigma_k\} \quad \text{params for each cluster } k$$

- 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

attempting to compute $\hat{\theta}^{ML}$

- Initial guess: $\hat{\theta}^{(0)}$

- Estimate at iteration t : $\hat{\theta}^{(t)}$

- E-Step

$$\text{Compute } U(\theta, \hat{\theta}^{(t)}) = E[\log p(y | \theta) | x, \hat{\theta}^{(t)}]$$

- M-Step

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

→ $\hat{\theta}^{(t)}$ converges to local mode

©Emily Fox 2013

3

MAP Estimation

- Bayesian approach:

- Place **prior** $p(\theta)$ on parameters
- Infer **posterior** $p(\theta | x) = \frac{p(x|\theta)p(\theta)}{p(x)}$

- Many, many, many motivations and implications

- For the sake of this class, simplest motivation is to think of this as akin to regularization

$$\hat{\theta}^{MAP} = \arg \max_{\theta} \log p(\theta | x) = \arg \max_{\theta} \underbrace{\log p(x|\theta)}_{\text{ML term}} + \log p(\theta) \leftarrow \text{reg.}$$

- 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:

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

affects max
w.r.t. θ

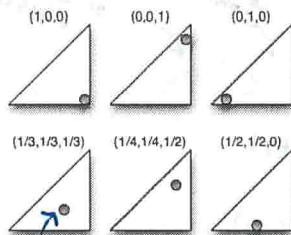
©Emily Fox 2013

5

MAP EM Example – MoG

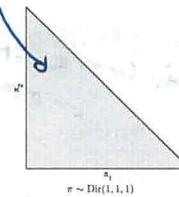
- For mixture of Gaussians, conjugate priors are:

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



(π_1, π_2, π_3) is a pt on the simplex

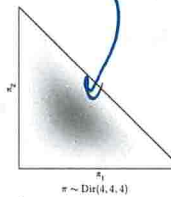
uniform on simplex



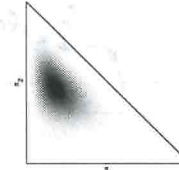
$$\pi \sim \text{Dir}(1, 1, 1)$$

$p(\theta|x)$ in same family as $p(\theta)$

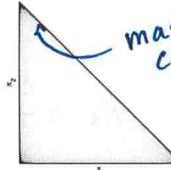
$\text{Dir}(\alpha_1, \alpha_2, \alpha_3)$ is prior on simplex



$$\pi \sim \text{Dir}(4, 4, 4)$$



$$\pi \sim \text{Dir}(4, 9, 7)$$



$$\pi \sim \text{Dir}(0.2, 0.2, 0.2)$$

mass at corners

hyperparams α determine form

©Emily Fox 2013

MAP EM Example – MoG

- For mixture of Gaussians, conjugate priors are:

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

- Dirichlet posterior

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

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

$$\propto \prod_k \left(\prod_{i: z^i = k} \pi_k \right) \cdot \pi_k^{\alpha_k - 1} \propto \prod_k \pi_k^{n_k + \alpha_k - 1}$$

- 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

$$p(x^{N+1} | x^1, \dots, x^N) \leftarrow \text{int. over uncertain model params} \\ = \int p(x^{N+1} | \theta) p(\theta | x^1, \dots, x^N) d\theta$$

- 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)

$$\hat{f}(\theta) = \frac{1}{M} \sum_{i=1}^M f(\theta^{(i)}) \\ \theta^{(i)} \sim \pi(\theta)$$

©Emily Fox 2013

9

Gibb Sampling

- Want draws:

$$(\theta_1, \dots, \theta_n) \sim \pi(\theta)$$

n parameters or latent vars
can't sample directly

- Construct Markov chain whose steady state distribution is $\pi(\theta)$

- Simplest case:

For $t=1, \dots, N_{iter}$

For $i=1, \dots, n$

can use random ordering

$$\theta_i^{(t)} \sim p(\theta_i | \theta_1^{(t)}, \dots, \theta_{i-1}^{(t)}, \theta_{i+1}^{(t-1)}, \dots, \theta_n^{(t-1)})$$

Gibbs sampling assumes that this has a closed-form that we can sample from

©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:

Priors

$$\pi \sim \text{Dir}(\alpha_1, \dots, \alpha_K) \quad z^i \sim \pi \quad i=1, \dots, N$$

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

■ Want to draw posterior samples of model parameters

$$\begin{aligned} \pi^{(t)} &\sim p(\pi | \phi^{(t-1)}, x^1, \dots, x^N) \\ \phi^{(t)} &\sim p(\phi | \pi^{(t)}, x^1, \dots, x^N) \end{aligned} \quad \left. \vphantom{\begin{aligned} \pi^{(t)} \\ \phi^{(t)} \end{aligned}} \right\} \text{NO CLOSED FORM}$$

only cond. on obs.

©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

Ex. Assume just one var of interest θ

Sample $\theta^{(t)} \sim p(\theta | z^{(t-1)})$ \leftarrow each has

$z^{(t)} \sim p(z | \theta^{(t)})$ \leftarrow closed form

$\theta^{(1)}$	$\theta^{(2)}$	$\theta^{(3)}$	\dots
$z^{(1)}$	$z^{(2)}$	$z^{(3)}$	\dots

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

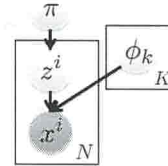
©Emily Fox 2013

12

Example – Mixture of Gaussians

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

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



■ Try auxiliary variable sampler

- Introduce cluster indicators into sampler

$$\pi^{(t)} \sim p(\pi | z_{1:N}^{(t-1)}, \{x_{1:N}, \lambda\}) = \text{Dir}(n_1 + \alpha_1, \dots, n_K + \alpha_K)$$

For $k=1, \dots, K$

$$\{\mu_k^{(t)}, \Sigma_k^{(t)}\} \sim p(\phi_k | z_{1:N}^{(t-1)}, \{x_{1:N}, \lambda\})$$

$$\text{what if high dim?} \quad = p(\phi_k | \{x^i : z^{(t-1)} = k\}, \lambda) = \text{NIW}(\text{post})$$

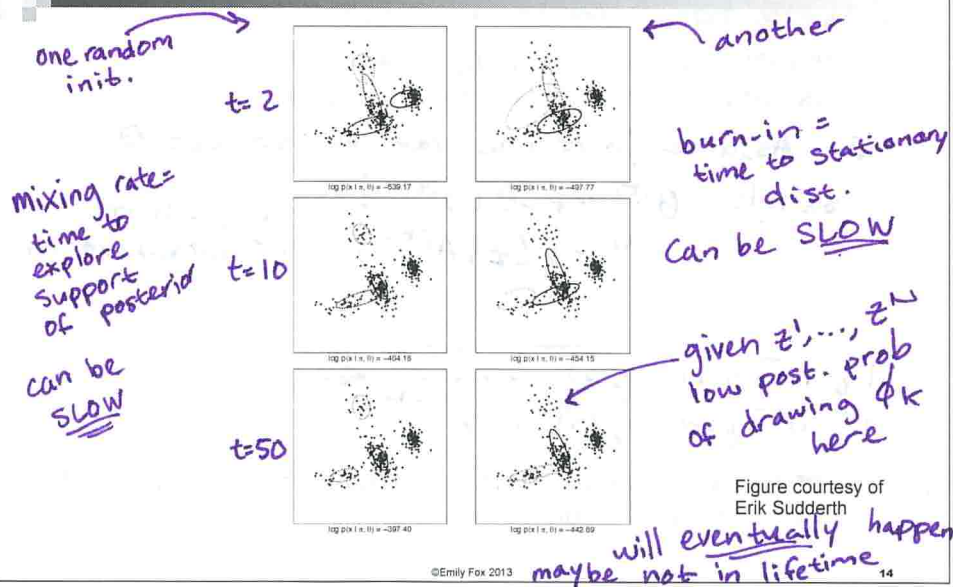
For $i=1, \dots, N$

$$z^i{}^{(t)} \sim p(z^i | x^i, \pi^{(t)}, \phi^{(t)}) \propto \pi_{z^i}^{(t)} N(x^i | \phi_{z^i}^{(t)}) \quad \text{as in EM}$$

©Emily Fox 2013

13

Example – Clustering Results I



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

what you don't care about

θ = param of interest, but high dim
 z = enables sampling

What about just sampling z 's? Feasible?

If so, can for $\hat{\theta}^{(t)}$ from sampled $z_{1:N}^{(t)}$

- Can improve efficiency if marginalized variables are high-dim

- Reduced dimension of search space
- But, often introduces dependences!

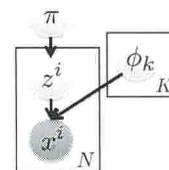
©Emily Fox 2013

15

Example – Collapsed MoG Sampling

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

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



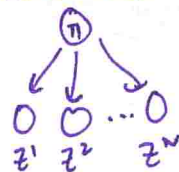
- Collapsed sampler

For $i=1, \dots, N$

$$z^{i(t)} \sim p(z^i | z^{1(t)}, \dots, z^{i-1(t)}, z^{i+1(t)}, \dots, z^{N(t)}, x_{1:N}, \alpha, \lambda)$$

$z_{\setminus i}^{(t)}$

cond. ind.



now tightly coupled



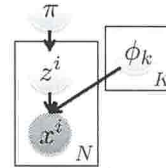
©Emily Fox 2013

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

$$p(z^i | z_{1:i}, x_{1:i}, \alpha, \lambda) \propto p(z^i | z_{1:i}, \alpha)$$

$$p(z^i | z_{1:i}, \alpha) = \int p(z^i = k | \pi) p(\pi | z_{1:i}) d\pi = \frac{n_{k-1}^i + \alpha_k}{N-1 + \sum \alpha_k}$$

$$p(x^i | z_{1:i}, x_{1:i}, \lambda) \propto \int \prod_{i: z^i = k} p(x^i | \phi_k) p(\phi_k | \lambda) d\phi_k$$

counts not including $z^i = k$

= student-t pred. like.

■ Important facts:

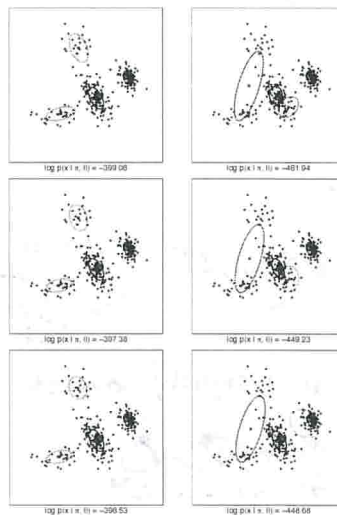
$$p(z_{1:N} | \alpha) = \frac{\Gamma(\sum_k \alpha_k) \prod_k \Gamma(n_k + \alpha_k)}{\prod_k \Gamma(\alpha_k) \Gamma(\sum_k n_k + \alpha_k)}$$

$$\frac{\Gamma(m+1)}{\Gamma(m)} = m$$

©Emily Fox 2013

17

Example – Clustering Results II



← worst case is still bad

seq. sampling z^i can make it hard to make global assign. changes (strong dependencies)

Trade-off

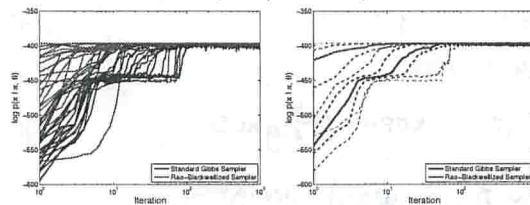
Figure courtesy of Erik Sudderth

©Emily Fox 2013

18

Comparing Collapsed vs. Regular

Log Likelihood vs. Gibbs Iteration
(multiple chains)



Overall, in this case,
collapsed has faster burn-in
typically, but worst
case still the same

Figure courtesy of
Erik Sudderth

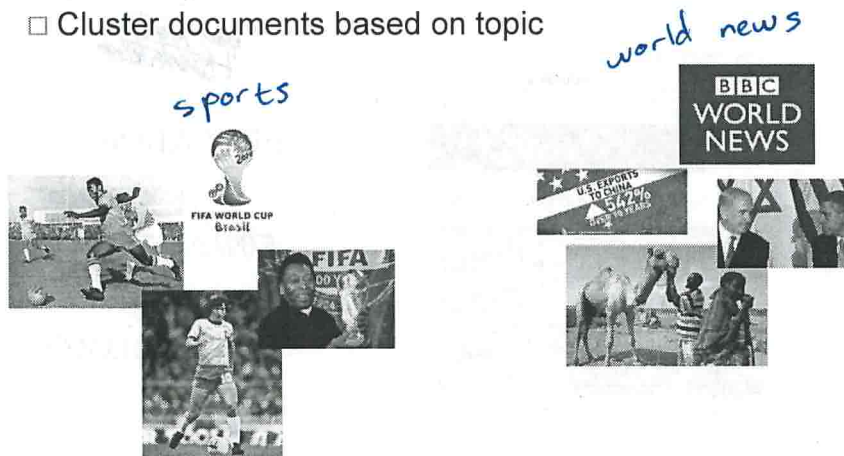
©Emily Fox 2013

18

Task 2: Cluster Documents

■ Previously:

- Cluster documents based on topic

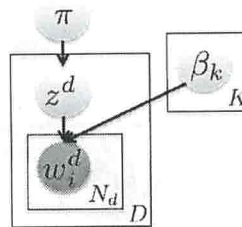


©Emily Fox 2013

19

A Generative Model

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



$z^d \sim \pi$ topic ~~weights~~
 $w_i^d | z^d \sim \beta_{z^d}$ word weights for topic z^d

Bayesian approach:

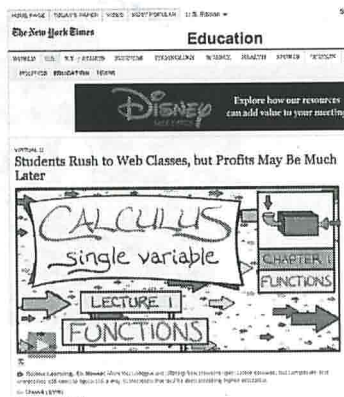
Priors $\pi \sim \text{Dir}(\alpha_1, \dots, \alpha_K)$ size of vocab
 $\beta_k \sim \text{Dir}(\lambda_1, \dots, \lambda_V)$ $k=1, \dots, K$

©Emily Fox 2013

20

Task 2: Cluster Documents

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



EDUCATION

FINANCE

TECHNOLOGY

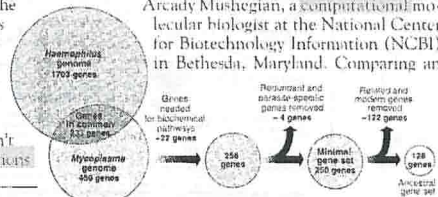
Latent Dirichlet Allocation (LDA)

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, these predictions

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

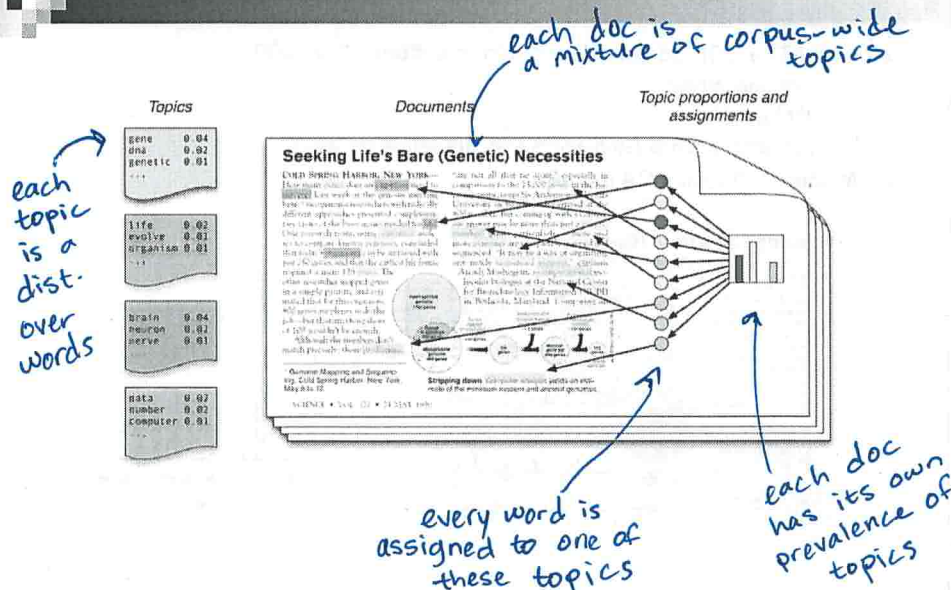
"are not all that far apart," especially in comparison to the 75,000 genes in the human genome, notes Siv Andersson of Uppsala University in Sweden, who arrived at the 800 number. But coming up with a consensus answer may be more than just a genetic numbers game, particularly as more and more genomes 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



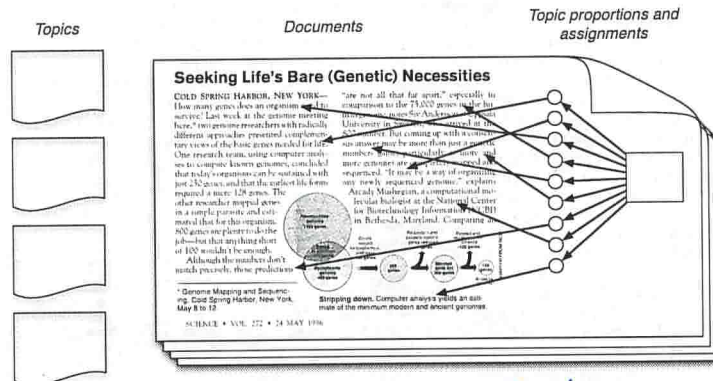
Stripping down. Computer analysis yields an estimate of the minimum modern and ancient genomes.

mixed membership model

Latent Dirichlet Allocation (LDA)



Latent Dirichlet Allocation (LDA)



Only observe the words!
Want posterior $p(\text{topics, proportions, assignments} | \text{docs})$

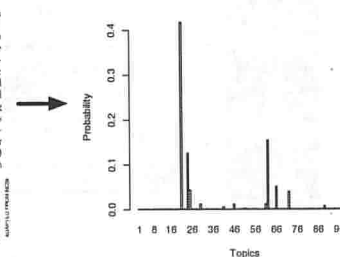
Example Inference – Topic Weights

- 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 genomic 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 Sir 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



Example Inference – Topic Words

highest prob. words in topic

topic ¹	2	3	4	...
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: $w_1^d, \dots, w_{N_d}^d$ $d=1, \dots, D$
- Associated topics: $z_1^d, \dots, z_{N_d}^d$ *corpus-wide topics that define word probabilities*
- Parameters: $\theta = \{\{\pi^d\}, \{\beta_k\}\}$
- Generative model: *doc-specific topic weights*

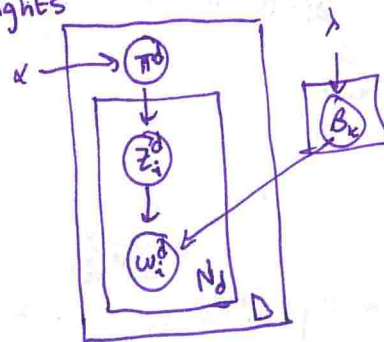
$$z_i^d \sim \pi^d$$

$$w_i^d | z_i^d \sim \beta_{z_i^d}$$

Priors:

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

$$\beta_k \sim \text{Dir}(\lambda_1, \dots, \lambda_V)$$



©Emily Fox 2013

28