**Case Study 5: Mixed Membership Modeling**

## Latent Dirichlet Allocation Collapsed Gibbs Sampler,

## Variational Methods

Machine Learning for Big Data
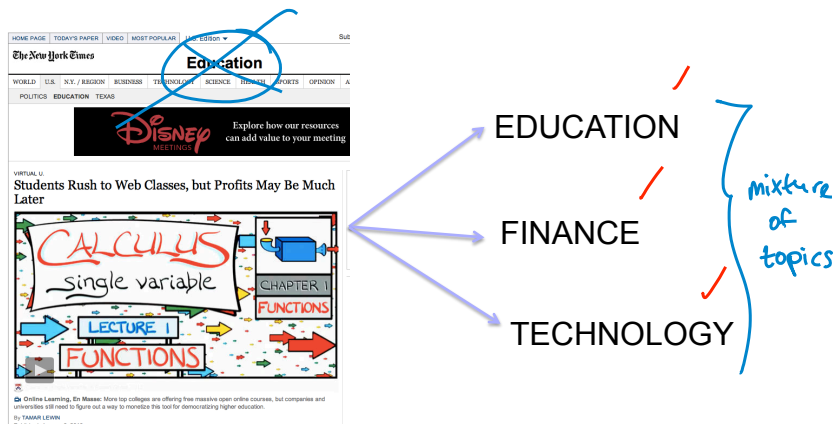CSE547/STAT548, University of Washington

Emily Fox

March 11th, 2014

©Emily Fox 2014                     1

---

# Task 3: Mixed Membership Models

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



EDUCATION

FINANCE

TECHNOLOGY

mixture of topics

©Emily Fox 2014                     2

# Latent Dirichlet Allocation (LDA)

*"Global" params*

*each topic k*

*is a distribution over words in vocab*

$\beta_k$, just as before

$\beta_k$

1 2 ... V

$\beta_k$

*previously, each doc had one topic*

*now, each is a mixture of topics*

**Topics**

```
gene     0.04
dna      0.02
genetic  0.01
...
```

```
life     0.02
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**

### 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

**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

**Topic proportions and assignments**

$z_{d,i}$

$\pi_d$

*topic 1 topic 2 topic 3*

*every word is assigned to a topic*

*each doc has its own prevalence of topics in that doc*

©Emily Fox 2014

3

---

# Latent Dirichlet Allocation (LDA)

**Topics**

**Documents**

### 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

**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

**Topic proportions and assignments**

**Obs:** All we see are words — $\beta_k$'s

**Want:** posterior $p($topics, doc prop. of topics, assign. vars. | words in docs$)$

©Emily Fox 2014

4

2

# LDA Generative Model

- Observations: $w_1^d, \ldots, w_{N_d}^d$
- Associated topics: $z_1^d, \ldots, z_{N_d}^d$ ← *topic per word in doc d*
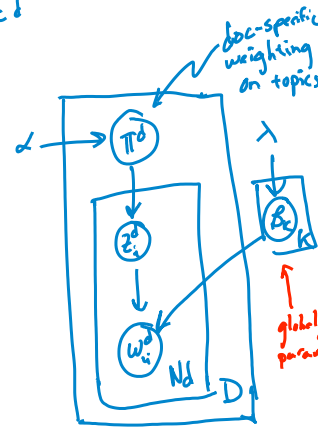- Parameters: $\theta = \{\{\pi^d\}, \{\beta_k\}\}$
- Generative model:

$$z_i^d \sim \pi^d \quad d=1,\ldots,D$$
$$\qquad i=1,\ldots,N_d$$

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

Priors:
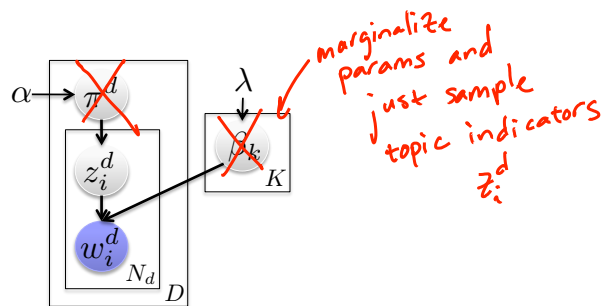$$\pi^d \sim Dir(\alpha_1,\ldots,\alpha_K) \quad d=1,\ldots,D$$
$$\beta_k \sim Dir(\lambda_1,\ldots,\lambda_V) \quad k=1,\ldots,K$$

*doc-specific weighting on topics*

*global params*

---

# LDA Joint Probability

*marginalize params and just sample topic indicators $z_i^d$*

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

3

# Collapsed LDA Sampling

- Marginalize parameters
  - Document-specific topic weights
  - Corpus-wide topic-specific word distributions

$$p(z_i^d = k \mid z_{\backslash id}, \{w_i^d\}, \alpha, \lambda)$$
$$\propto p(z_i^d = k \mid z_{\backslash id}, \alpha) p(w_i^d \mid z_i^d = k, z_{\backslash id}, w_{\backslash id}, \lambda)$$

- Unplate to see dependencies induced

*(handwritten annotations)* separately for each doc ($\pi^d$ = local var) · global var · All $z_\cdot^d = k$: (in all of corpus) · marg. $\pi^d$'s · $\beta_k$ · & marg. $\beta_k$ introducing coupling in all of above

7

# Collapsed LDA Sampling

- Sample topic indicators for each word
  - Algorithm:

$$p(z_i^d = k \mid z_{\backslash id}, \{w_i^d\}, \alpha, \lambda)$$
$$\propto p(z_i^d = k \mid \{z_j^d, j \neq i\}, \alpha) p(w_i^d \mid \{w_j^c : z_j^c = k, (j,c) \neq (i,d)\}, \lambda)$$

*(handwritten annotations)*

"prior" · "likelihood" · only dependence within doc $d$

$$\propto \frac{n_k^{-id} + \alpha_k}{N_d - 1 + \sum \alpha_k} \cdot \frac{m_{w_i^d, k}^{-id} + \lambda_{w_i^d}}{\sum_\gamma \left( m_{\gamma, k}^{-id} + \lambda_\gamma \right)}$$

# words assigned to topic k in doc d not counting $i^{th}$ word

normalize within doc

over whole corpus

$m_{\gamma, k}^{-id}$ = # of times word $\gamma$ appears in topic k (not counting $w_i^d$)

8

4

# Select a Document

| | | | | |
|---|---|---|---|---|
| Etruscan | trade | price | temple | market |

*all words in doc d*

# Randomly Assign Topics

$z_i^d$

$w_i^d$

| 3 | 2 | 1 | 3 | 1 |
|---|---|---|---|---|
| Etruscan | trade | price | temple | market |

*one approach to initialize sampler*

# Randomly Assign Topics

| 3 | 2 | 1 | 3 | 1 |
|---|---|---|---|---|
| Etruscan | trade | price | temple | market |

$z_i^d$

$w_i^d$

| 2 | 3 | 2 | 1 | 1 |
|---|---|---|---|---|
| Italy | temple | ship | trade | market |

do for all docs in corpus

©Emily Fox 2014                                        11

# Maintain Local Statistics

| 3 | 2 | 1 – | 3 | 1 |
|---|---|---|---|---|
| Etruscan | trade | price | temple | market |

$z_i^d$

$w_i^d$

| | Topic 1 | Topic 2 | Topic 3 |
|---|---|---|---|
| Doc d | 2 | 1 | 2 |

$n_1^d$     $n_2^d$     $n_3^d$

©Emily Fox 2014                                        12

6

# Maintain Global Statistics

$z_i^d$

$w_i^d$

| 3 | 2 | 1 | 3 | 1 |
|---|---|---|---|---|
| Etruscan | trade | price | temple | market |

|  | Topic 1 | Topic 2 | Topic 3 |
|---|---|---|---|
| Doc d | 2 | 1 | 2 |

$M_{trade,2}$ ← word ← topic

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

Total counts from **all** docs

13

---

# Resample Assignments

$z_i^d$

$w_i^d$

$z_2^d$

| 3 | ~~2~~ | 1 | 3 | 1 |
|---|---|---|---|---|
| Etruscan | trade | price | temple | market |

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

$n_k^{-id}$

|  | Topic 1 | Topic 2 | Topic 3 |
|---|---|---|---|
| Etruscan | 1 | 0 | 35 |
| market | 50 | 0 | 1 |
| price | 42 | 1 | 0 |
| temple | 0 | 0 | 20 |
| trade | 10 | 7 ~~8~~ | 1 |
| ... |  |  |  |

$M_{word,topic}^{-id}$

14

7

# What is the conditional distribution for this topic?

$z_i^d$

| 3 | ? | 1 | 3 | 1 |
|---|---|---|---|---|
| Etruscan | trade | price | temple | market |

$w_i^d$

$$p(z_i^d \mid \text{everything else})$$

---

# What is the conditional distribution for this topic?

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

$z_i^d$

| 3 | ? | 1 | 3 | 1 |
|---|---|---|---|---|
| Etruscan | trade | price | temple | market |

$w_i^d$

Topic 1          Topic 2          Topic 3

$$\frac{n_k^{-id} + d_k}{N_d - 1 + \sum d_k}$$

"prior" term

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

# 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?

$z_i^d$

$w_i^d$

| 3 | ? | 1 | 3 | 1 |
|---|---|---|---|---|
| Etruscan | trade | price | temple | market |

Topic 1            Topic 2            Topic 3

row from global table

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

©Emily Fox 2014                                                                17

---

# 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?

$z_i^d$

$w_i^d$

| 3 | ? | 1 | 3 | 1 |
|---|---|---|---|---|
| Etruscan | trade | price | temple | market |

Topic 1            Topic 2            Topic 3

$$\frac{m_{trade,k}^{-id} + \lambda_{trade}}{\sum_\gamma \left( m_{\gamma,k} + \lambda_\gamma \right)}$$

"likelihood" term

©Emily Fox 2014                                                                18

9

## 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?

$z_i^d$

$w_i^d$

| 3 | ? | 1 | 3 | 1 |
|---|---|---|---|---|
| Etruscan | trade | price | temple | market |

Topic 1　　　　　　Topic 2　　　　　　Topic 3

$$\frac{n_k^{-id} + \alpha_k}{N_d - 1 + \sum_k \alpha_k} \frac{m_{trade,k}^{-id} + \lambda_{trade}}{\sum_{\gamma=1}^{V} m_{\gamma,k}^{-id} + \lambda_\gamma}$$

©Emily Fox 2014　　　　19


## Sample a New Topic Indicator

$z_i^d$

$w_i^d$

| 3 | 1 | 1 | 3 | 1 |
|---|---|---|---|---|
| Etruscan | trade | price | temple | market |

Topic 1　　　　　　Topic 2　　　　　　Topic 3

*happened to sample topic 1*

©Emily Fox 2014　　　　20

10

# Update Counts

$z_i^d$

$w_i^d$

| 3 | 1 ~~X~~ | 1 | 3 | 1 |
|---|---|---|---|---|
| Etruscan | trade | price | temple | market |

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

|  | Topic 1 | Topic 2 | Topic 3 |
|---|---|---|---|
| Etruscan | 1 | 0 | 35 |
| market | 50 | 0 | 1 |
| price | 42 | 1 | 0 |
| temple | 0 | 0 | 20 |
| trade | 11 ~~10~~ | 7 | 1 |
| … |  |  |  |

21

# Geometrically…

inc. popularity of topic 1 in doc d
and word prevalence for topic 1 in corpus

$z_i^d$

$w_i^d$

| 3 | (1) | 1 | 3 | 1 |
|---|---|---|---|---|
| Etruscan | trade | price | temple | market |

Topic 1          Topic 2          Topic 3

22

# 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 include:
  - 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
  - And many, many more!

- Alternative: Variational methods instead of sampling
  - Approximate posterior with an optimized variational distribution

©Emily Fox 2014    23

---

# Case Study 5: Mixed Membership Modeling

# Variational Methods

Machine Learning for Big Data
CSE547/STAT548, University of Washington
Emily Fox
March 11th, 2014

©Emily Fox 2014    24

# Variational Methods Goal

- Recall task: Characterize the posterior $p(\theta, z \mid x)$ obs.

  params ↗ ↖ latent vars

- Turn posterior inference into an optimization task
- Introduce a "tractable" family of distributions over parameters and latent variables
  - Family is indexed by a set of "free parameters"
  - Find member of the family closest to: $p(\theta, z \mid x)$

Call the family $Q$ and want $q \in Q$
that is closest to $p(\theta, z \mid x)$

# Variational Methods Cartoon

- Cartoon of goal:

best $p(\theta, z \mid x)$    one member $q$ in family $Q$

$\sigma^2$    $\sigma^2$    $\sigma^2$

$\mu$    $\mu$    $\mu$

e.g., $Q$: all Gaussians

- Questions:
  - ① How do we measure "closeness"?
  - ② If the posterior is intractable, how can we approximate something we do not have to begin with?

# A Measure of Closeness

- Kullback-Leibler (KL) divergence
  - Measures "distance" between two distributions $p$ and $q$

$$KL(p\|q) \triangleq D(p\|q) = E_p\left[\log \frac{p}{q}\right] = \int_\theta p(\theta) \log \frac{p(\theta)}{q(\theta)} d\theta$$

- If $p = q$ for all $\theta$

$$D(p\|q) = \int p(\theta) \log 1 \, d\theta = 0$$

- Otherwise, $\quad D(p\|q) > 0$

27

---

# A Measure of Closeness

$$\mathrm{KL}(p||q) \triangleq D(p||q) = \int_\theta p(\theta) \log \frac{p(\theta)}{q(\theta)} d\theta$$

- Not symmetric $\quad D(p\|q) \neq D(q\|p) \quad \longleftarrow \int q(\theta) \log \frac{q(\theta)}{p(\theta)} d\theta$

$\Longrightarrow$ not a true distance metric

- $p$ determines where the difference is important:
  - $p(\theta)=0$ and $q(\theta)\neq0$ $\quad 0 \log 0 = 0$
  - $p(\theta)\neq0$ and $q(\theta)=0$ $\quad \epsilon \log \frac{\epsilon}{0} = \infty$ ($=\epsilon$)

$$\text{If } D(p\|q) \text{ finite}, \quad \mathrm{supp}(q) \supseteq \mathrm{supp}(p)$$

- Want $\quad \hat{q} = \arg\min_{q \in Q} D(p\|q)$

- Just as hard as the original problem! $\quad E_p[\cdots]$

28

14

# Reverse Divergence

- Divergence D($p \| q$)
  - true distribution $p$ defines support of diff.
  - the "correct" direction
  - will be intractable to compute
- Reverse divergence D($q \| p$)
  - approximate distribution defines support
  - tends to give overconfident results
  - will be tractable

*What we have control over*

*q now less diffuse than p*

*p*

*divergence*

*reverse*

# Interpretations of Minimizing Reverse KL

$$D(q\|p) = E_q\left[\log \frac{q}{p}\right]$$

- Similarity measure:

$$D(q(\theta,z)\|p(\theta,z|x)) = E_q[\log q(\theta,z)] - E_q[\log p(\theta,z|x)]$$

$$= E_q[\log q(\theta,z)] - E_q[\log p(\theta,z,x)] + \log p(x)$$

$$-\mathcal{L}(q)$$

- Evidence lower bound (ELBO)

$$\log p(x) = D(q(z,\theta)\|p(\theta,z|x)) + \mathcal{L}(q) \geq \mathcal{L}(q)$$

$$\geq 0$$

*"ELBO"*

# Interpretations of Minimizing Reverse KL

- Evidence lower bound (ELBO)

  *log marginal likelihood or "evidence"*

  $$\log p(x) = D(q(z,\theta)||p(z,\theta|x)) + \mathcal{L}(q) \geq \mathcal{L}(q)$$

  *const.*  *add to a const.*  *"ELBO"*

- Therefore,
  - ELBO provides a lower bound on marginal likelihood
  - Maximizing ELBO is equivalent to minimizing KL

  $$\max \mathcal{L}(q) = \min D(q||p) = \max \text{ lower bound of } \log p(x)$$

  *what we can control*  *depends on what we don't know*

# Mean Field    $\mathcal{L}(q) = E_q[\log p(z, \theta, x)] - E_q[\log q(z, \theta)]$

- How do we choose a *Q* such that the following is tractable?

  $$\hat{q} = \arg\max_{q \in Q} \mathcal{L}(q) \quad \leftarrow \text{new objective}$$

- Simplest case = mean field approximation  $\theta, z = \{z^1, \dots, z^N\}$
  - Assume each parameter and latent variable is conditionally independent given the set of free parameters

  $$q(z, \theta) = q(\theta|\gamma) \prod_{i=1}^{N} q(z^i|\phi^i)$$

  $\gamma, \{\phi^i\}$ are "free params" = control knobs in getting q close to p

  

  Original graph          Naïve mean field

  *can also look at "structured mean field approx (break only some depend.)*

# Mean Field

$$\mathcal{L}(q) = E_q[\log p(z, \theta, x)] - E_q[\log q(z, \theta)]$$

*entropy*

- Naïve mean field decomposition:

$$q(z, \theta) = q(\theta \mid \gamma) \prod_{i=1}^{N} q(z^i \mid \phi^i)$$

$q_\theta$    $q_{z^i}$

- Under this approximation, entropy term decomposes as

$$-E_q[\log q(z, \theta)] = -E_q[\log q(\theta \mid \gamma)] - \sum_i E_q[\log q(z^i \mid \phi^i)]$$

*decouples across $\gamma$, $\phi^i$*

- Can (always) rewrite joint term as   *full cond. of $\theta$*

$$E_q[\log p(\theta, z, x)] = E_q[\log p(\theta \mid z, x)] + E_q[\log p(z, x)]$$

OR    *full cond. of $z^i$*

$$E_q[\log p(\theta, z, x)] = E_q[\log p(z^i \mid z_{\backslash i}, \theta, x)] + E_q[\log p(z_{\backslash i}, \theta, x)]$$

---

# Mean Field – Optimize $\gamma$

- Examine one free parameter, *e.g.*, $\gamma$

$$\mathcal{L}(q) = E_q[\log p(\theta \mid z, x)] + E_q[\log p(z, x)] - E_q[\log q(\theta \mid \gamma)] - \sum_i E_q[\log q(z^i \mid \phi^i)]$$

*consider $\theta$-full-cond. form*

- Look at terms of ELBO just depending on $\gamma$

*don't depend on $\gamma$ because under $q$, $z^i \perp\!\!\!\perp \theta$!*

$$\mathcal{L}^\gamma = E_q[\log p(\theta \mid z, x)] - E_q[\log q(\theta \mid \gamma)] + const.$$

*w.r.t. $\gamma$*

*really just $q_\theta = q(\theta \mid \gamma)$ needed here*

# Mean Field – Optimize $\phi^i$

- Examine another free parameter, *e.g.*, $\phi^i$

$$\mathcal{L}(q) = E_q[\log p(z^i \mid z_{\setminus i}, \theta, x)] + E_q[\log p(z_{\setminus i}, \theta, x)] - E_q[\log q(\theta \mid \gamma)] - \sum_i E_q[\log q(z^i \mid \phi^i)]$$

*consider the $z^i$-full-cond. form*

*const. wrt $\phi^i$*

  □ Look at terms of ELBO just depending on $\phi^i$

$$\mathcal{L}^{\phi^i} = E_q\left[\log p(z^i \mid z_{\setminus i}, \theta, x)\right] - E_q\left[\log q(z^i \mid \phi^i)\right]$$

*really just $q_{z^i} = q(z^i \mid \phi^i)$ here*

- This motivates using a coordinate ascent algorithm for optimization
  □ Iteratively optimize each free parameter holding all others fixed

35

---

# Algorithm Outline

- **Initialization:** Randomly select starting distribution $q_\theta^{(0)}$
- **E-Step:** Given parameters, find posterior of hidden data

  *optimize $\phi$* → $q_z^{(t)} = \arg\max_{q_z} \mathcal{L}(q_z, q_\theta^{(t-1)})$  *latent vars $z$*

- **M-Step:** Given posterior distributions, find likely parameters $\theta$

  *optimize $\gamma$* → $q_\theta^{(t)} = \arg\max_{q_\theta} \mathcal{L}(q_z^{(t)}, q_\theta)$

- **Iteration:** Alternate E-step & M-step until convergence

36

18

# What you need to know…

- Latent Dirichlet allocation (LDA)
  - Motivation and generative model specification
  - Collapsed Gibbs sampler

- Variational methods
  - Overall goal
  - Interpretation in terms of minimizing (reverse) KL
  - Mean field approximation

37

---

# Reading

- **Mixed Membership Models: KM Sec. 27.3**
  - Basic LDA:
    Blei, David M., Andrew Y. Ng, and Michael I. Jordan. "Latent dirichlet allocation." the Journal of machine Learning research 3 (2003): 993-1022.
  - Introduction:
    Blei, David M. "Probabilistic topic models." Communications of the ACM, vol. 55, no. 4 (2012): 77-84.
  - Sampling:
    Griffith, Thomas L. and Mark Steyvers. "Finding scientific topics." Proceedings of the National Academy of Sciences of the United States of America, Volume: 101, Supplement: 1 (2004): Pages: 5228-5235

38

# Acknowledgements

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

**39**