

CS 589 Fall 2021 Lecture 6

Topic model

**Monday 6:30-9:00
Babbio 122**

All zoom links in Canvas



photo: <https://www.scubedstudios.com/information-retrieval/>

Review of Lecture 5

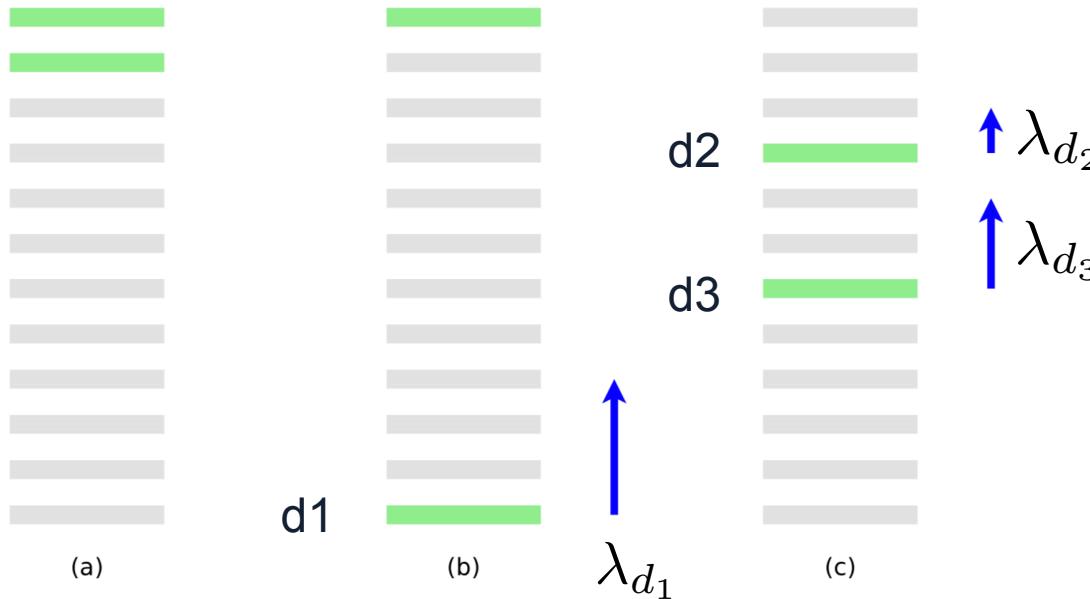
- The idea of machine learning for ranking has been around for a long time
- Learning to rank can be turned into a classification or regression problem
- LTR algorithms:
 - Multiple additive regression tree (MART)
 - RankNet
- LTR categories:
 - Pointwise
 - Pairwise
 - Listwise
- Yahoo! Learning to rank challenge (Tree ensemble algorithm wins!)

Review of Lecture 5

- Multiple additive regression tree (MART)
 - Minimizes the square loss (i.e., just like linear regression and regression tree)
 - First iteration: $f_0(x) = \text{avg of all } y_i$
 - From iteration 1 through k:
 - Update the residual of each y_i
 - Fit a regression tree $f_k(x)$ for the residuals
 - Update the sum $F(x)$
- The name “gradient boosting regression tree” comes from the fact that the residual = gradient of the objective function

Review of Lecture 5

- RankNet
 - Convert the pairwise loss function to the summation of λ_i for each document d_i
 - That is, λ_i is a sort of “gradient” with respect to d_i



$$w_k = w_k - \alpha \frac{\partial \Sigma}{\partial w_k}$$

$$\frac{\partial \Sigma}{\partial w_k} = \uparrow \cdot \frac{\partial s_1}{\partial w_k} + \uparrow \cdot \frac{\partial s_2}{\partial w_k} + \uparrow \cdot \frac{\partial s_3}{\partial w_k}$$

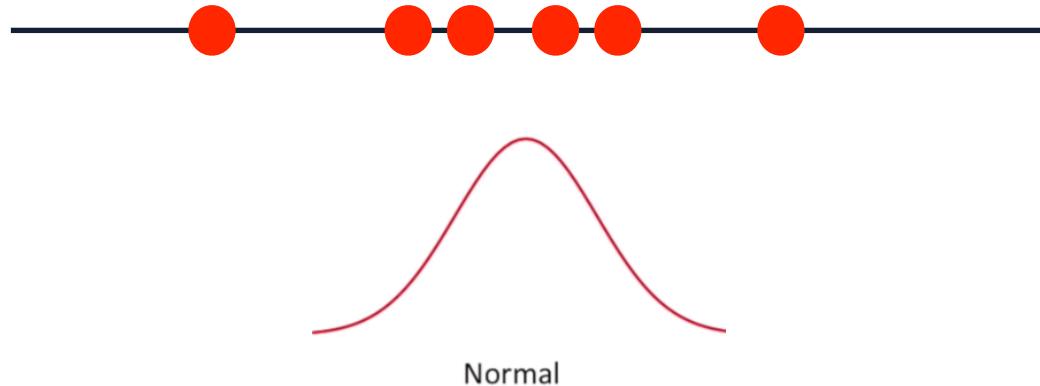
$$s_i = w x_i + b$$

Lecture 6: Topic model

- Pop quiz
- Expectation maximization
- Probabilistic topic model
 - The toy coin topic problem
 - Probabilistic latent semantic analysis
- Latent Dirichlet allocation
 - Bayesian inference of topic model
 - Variational inference for LDA
 - Gibbs sampling, Markov chain Monte-Carlo

Maximum likelihood estimation

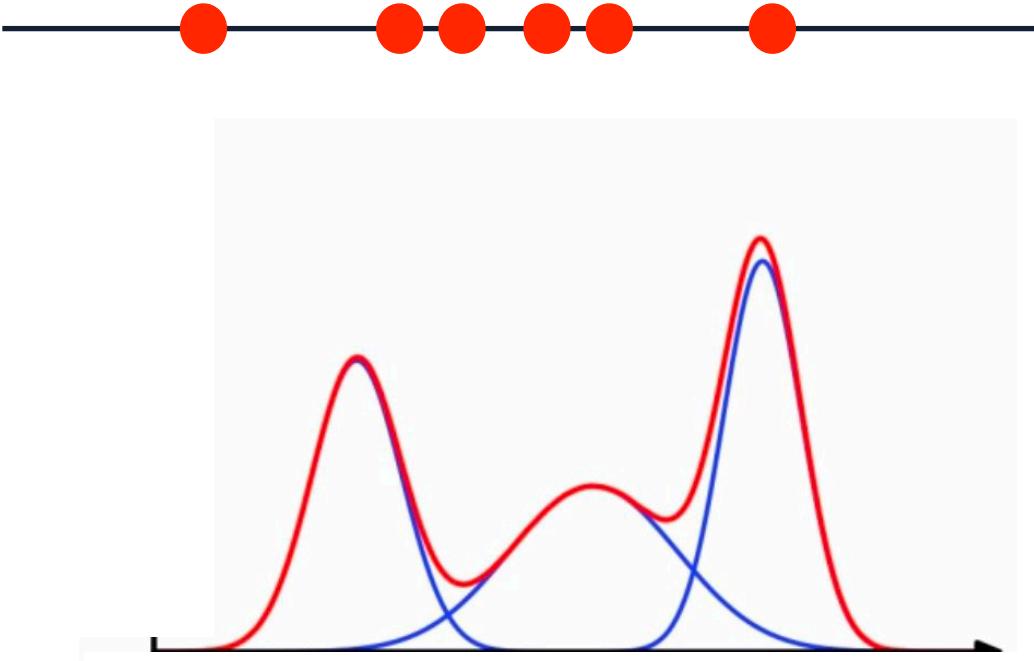
- Estimating the underlying distribution behind the observations
 - i.e., by building a model with parameter θ and find the θ that maximizes the likelihood of the observation
- Example 1: Estimating the weights of mice from observations x_1, \dots, x_n



$$\begin{aligned}\theta^* &= \operatorname{argmax}_{\mu, \sigma} \sum_{i=1}^n f_X(x_i; \mu, \sigma^2) \\ &= \operatorname{argmax}_{\mu, \sigma} \sum_{i=1}^n \frac{1}{\sigma \sqrt{2\pi}} e^{-\frac{(x_i - \mu)^2}{2\sigma^2}}\end{aligned}$$

Maximum likelihood estimation for mixture models

- Example 2: Estimating the weights of mice from observations x_1, \dots, x_n



$$\theta^* = \operatorname{argmax}_{\mu, \sigma} \sum_{i=1}^n \sum_{t=1}^3 p(z_i = t) \cdot f_X(x_i; \mu_t, \sigma_t^2)$$

$$= \operatorname{argmax}_{\mu, \sigma} \sum_{i=1}^n \sum_{t=1}^3 p(z_i = t) \frac{1}{\sigma \sqrt{2\pi}} e^{-\frac{(x_i - \mu_t)^2}{2\sigma_t^2}}$$



latent variable z_i : which topic x_i comes from

Expectation Maximization

- The expectation maximization algorithm is an algorithm for maximum likelihood estimation under the **latent space** assumption
 - Likelihood of the incomplete space: the observation

$$p(\mathbf{x}; \theta)$$

- Likelihood of the complete space: the observation + the latent variable, e.g., the latent topic

$$p(\mathbf{x}, z | \theta)$$

- Estimating the incomplete probability by marginalizing out the complete space

$$p(\mathbf{x}; \theta) = \sum_{k=1}^K p(\mathbf{x}|z; \theta)p(z = k)$$

Expectation maximization algorithm

- EM algorithm: repeat $n=1\dots N$:

E step:

$$Q(\theta|\hat{\theta}^{(n)}) = \mathbb{E}_{z|\mathbf{x},\hat{\theta}^{(n)}} [\log p(\mathbf{x}, z|\theta)]$$

M step:

$$\hat{\theta}^{(n+1)} = \arg \max_{\theta \in \mathcal{S}} Q(\theta|\hat{\theta}^{(n)})$$

- **Theorem:** the likelihood of observation, $\log p(\mathbf{x}; \theta^{(n)})$, monotonously increases with n (*Proof on whiteboard*)

An example problem: Coin-topic problem

- Author H and author T are co-authoring a paper in the following way:
 - At each time, they toss a coin to write the next word. If it's "head", author H writes the next word, if it's "tail", author T writes the next word. The probability for "head" is λ
 - The head author selects the next word by randomly sampling from $p(w|H)$, the tail author selects the next word by randomly sampling from $p(w|T)$



Problem 1: $p(w|H)$ and $p(w|T)$ are known, estimate λ

$$l_{cd}(\theta) = \sum_i \log p(w_i|z_i)p(z_i)$$

$$\begin{aligned}\mathbb{E}_{z|d,\lambda}[l_{cd}(\theta)] &= \max_{\lambda} \sum_i p(z_i = H|d, \lambda) \log p(w_i|H)\lambda + p(z_i = T|d, \lambda) \log p(w_i|T)(1 - \lambda) \\ &= \sum_i \sum_v \mathbf{1}[w_i = v] (p(z = H|v) \log p(v|H)\lambda + p(z = T|v) \log p(v|T)(1 - \lambda))\end{aligned}$$

- where both head and tail distributions are known, e.g.:

	the	computer	data	baseball	game	interesting
$p(w T)$	0.2	0.05	0.05	0.4	0.2	0.1
$p(w H)$	0.25	0.2	0.2	0.15	0.1	0.1

Expectation maximization

- Find θ that maximizes $\mathbb{E}_{z|d,\lambda}[l_{cd}(\theta)]$ (take the derivative):

$$\sum_i \sum_v \mathbf{1}[w_i = v] (p(z = H|v) \frac{1}{\lambda} + p(z = T|v) \frac{1}{1-\lambda}) = 0$$

$$\Rightarrow \lambda \propto \sum_v p(z = H|v) (\sum_i \mathbf{1}[w_i = v])$$

$$\Rightarrow \lambda^{(n+1)} \propto \frac{1}{|d|} \sum_v p^{(n)}(z = H|v) count(d, v) \quad \text{(M step)}$$

Expectation maximization

- E step: apply Bayes' theorem:

$$p^{(n+1)}(z = H|v) \propto p(v|z = H)p(z = H) = p(v|H)\lambda^{(n)}$$

$$p^{(n+1)}(z = T|v) \propto p(v|z = T)p(z = T) = p(v|T)(1 - \lambda^{(n)})$$

Problem 2: λ and $p(w|H)$ known, estimating $p(w|T)$

$$\begin{aligned}\mathbb{E}_{z|d,\lambda}[l_{cd}(\theta)] &= \max_{\lambda} \sum_i p(z_i = H|d, \lambda) \log p(w_i|H)\lambda + p(z_i = T|d, \lambda) \log \textcolor{red}{p(w_i|T)}(1 - \lambda) \\ &\quad - \eta \left(\sum_v \textcolor{red}{p(v|T)} - 1 \right) \\ &= \sum_i \sum_v \mathbb{1}[w_i = v] p(z = H|v) \log p(v|H)\lambda + p(z = T|v) \log \textcolor{red}{p(v|T)}(1 - \lambda) \\ &\quad - \eta \left(\sum_v \textcolor{red}{p(v|T)} - 1 \right)\end{aligned}$$

- Take the derivative for $p(v|T)$ and set to 0, we can get (M step):

$$p^{(n+1)}(v|T) = \frac{\sum_i \mathbb{1}[w_i = v] p^{(n)}(z = H|v)}{\sum_v \sum_i \mathbb{1}[w_i = v] p^{(n)}(z = H|v)}$$

Problem 2: λ and $p(w|H)$ known, estimating $p(w|T)$

$$\begin{aligned}\mathbb{E}_{z|d,\lambda}[l_{cd}(\theta)] &= \max_{\lambda} \sum_i p(z_i = H|d, \lambda) \log p(w_i|H)\lambda + p(z_i = T|d, \lambda) \log \textcolor{red}{p(w_i|T)}(1 - \lambda) \\ &\quad - \eta \left(\sum_v \textcolor{red}{p(v|T)} - 1 \right) \\ &= \sum_i \sum_v \mathbb{1}[w_i = v] p(z = H|v) \log p(v|H)\lambda + p(z = T|v) \log \textcolor{red}{p(v|T)}(1 - \lambda) \\ &\quad - \eta \left(\sum_v \textcolor{red}{p(v|T)} - 1 \right)\end{aligned}$$

- Take the derivative for $p(v|T)$ and set to 0, we can get (M step):

$$p^{(n+1)}(v|T) = \frac{\sum_i \mathbb{1}[w_i = v] p^{(n)}(z = H|v)}{\sum_v \sum_i \mathbb{1}[w_i = v] p^{(n)}(z = H|v)}$$

(E step same as Problem 1)

Problem 2: λ and $p(w|H)$ known, estimating $p(w|T)$

- How can Problem 2 be useful?
 - Suppose $p(w|T)$ is the main topic (Python)
 - $p(w|H)$ is the background topic (StackOverflow general question)
 - The mixture of head and tail topic is dominated by background words:
 - After stop words removal, the true topic $p(w|T)$ is “revealed”:

Problem 3: Unknown topic and λ

- Suppose $p(w|H)$ is unknown, but $p(w|T)$ and λ are both unknown:
 - Use Problem 1 to estimate λ
 - Use Problem 2 to estimate $p(w|T)$
 - Same E step

Applications of Coin Topic Problem for Text Mining

- Application Scenarios:

- $p(w|H)$ & $p(w|T)$ are known; estimate λ ↪

how much percent of the document is about computer game?
- $p(w|H)$ & λ are known; estimate $p(w|T)$ ↪

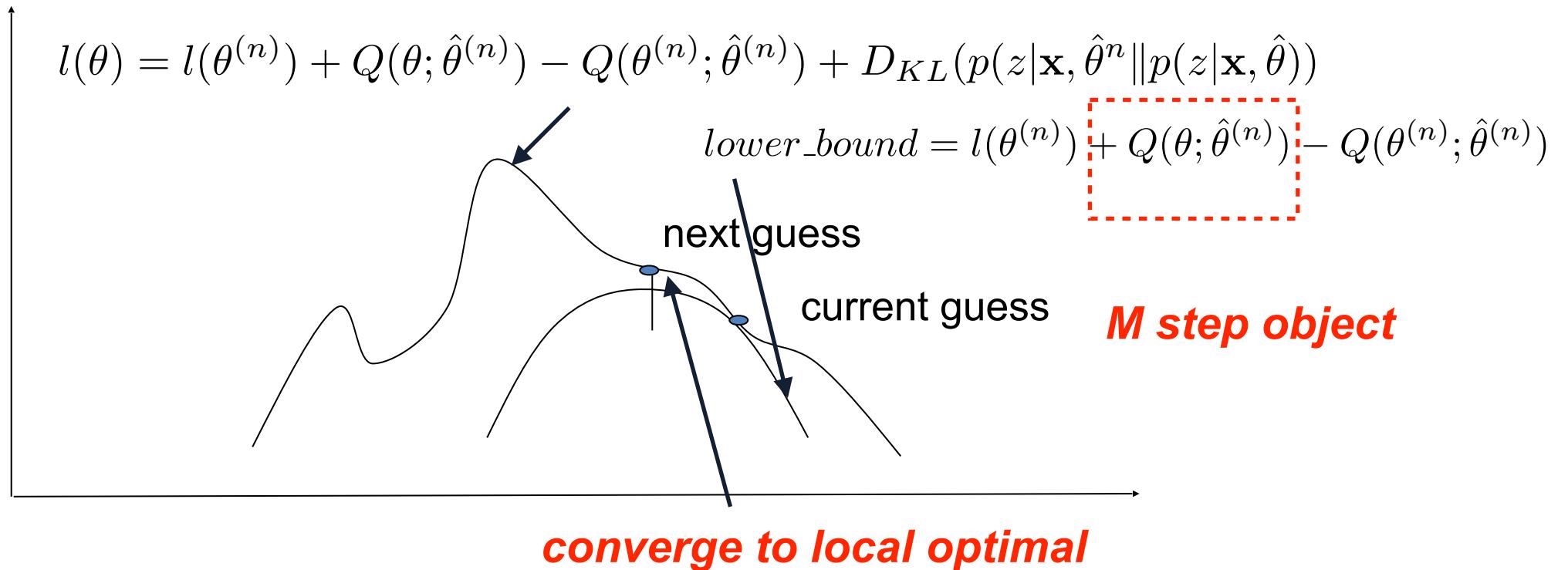
30% of the doc is about computer game, what's the other topic about?
- $p(w|H)$ is known; estimate λ & $p(w|T)$ ↪

The doc is about computer game, is it also about some other topic, and if so to what extent?
- λ is known; estimate $p(w|H)$ & $p(w|T)$ ↪

30% of the doc is about one topic and 70% is about another, what are these two topics?
- Estimate λ , $p(w|H)$, $p(w|T)$ ↪

The doc is about two subtopics, find out what these two subtopics are and to what extent the doc covers each.

Expectation maximization as hill climbing



EM algorithm in action

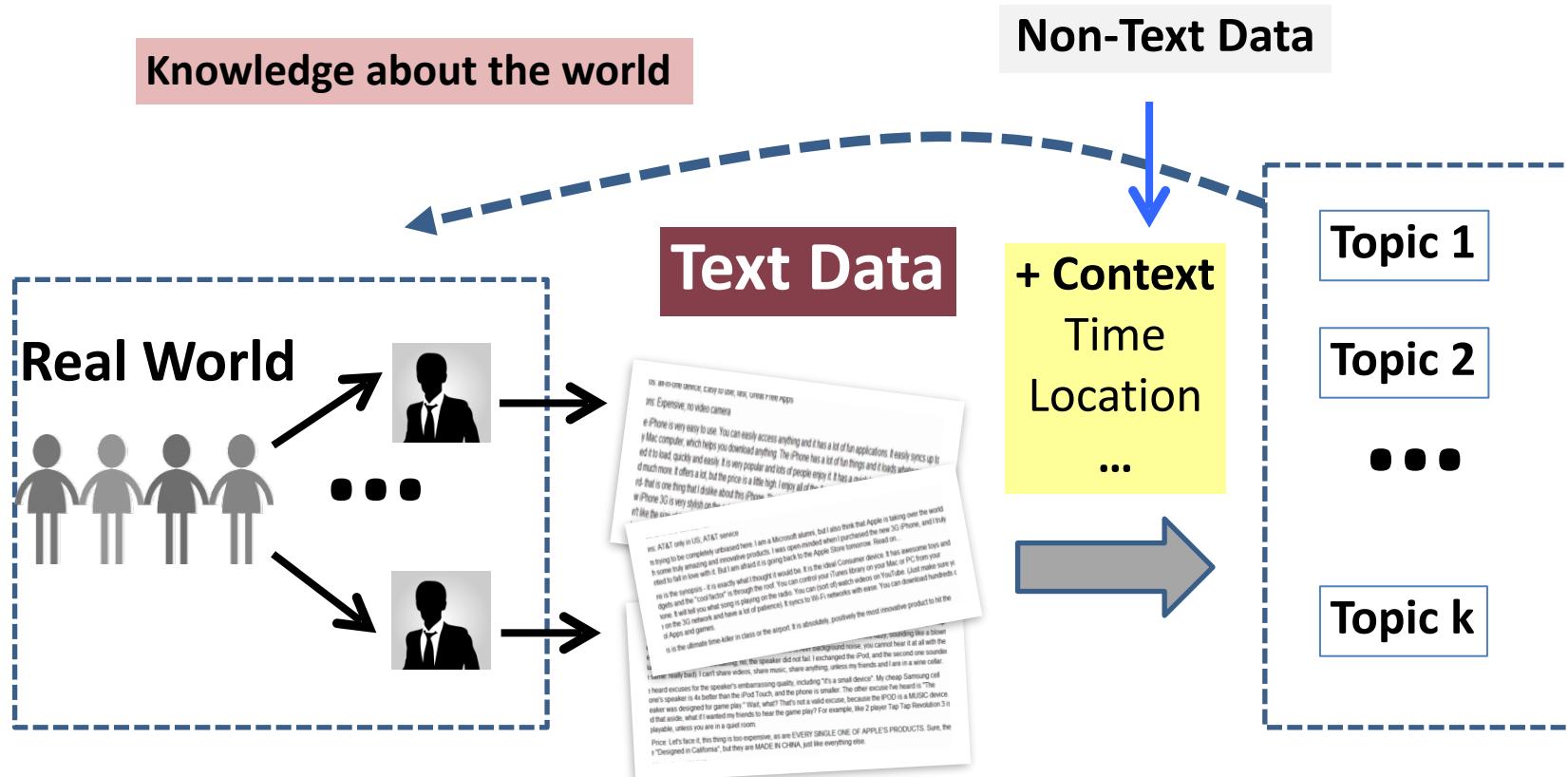
- Log likelihood increases:

Word	#	$p(w \theta_B)$	Iteration 1		Iteration 2		Iteration 3	
			$P(w \theta)$	$p(z=0 w)$	$P(w \theta)$	$p(z=0 w)$	$P(w \theta)$	$p(z=0 w)$
The	4	0.5	0.25	0.33	0.20	0.29	0.18	0.26
Paper	2	0.3	0.25	0.45	0.14	0.32	0.10	0.25
Text	4	0.1	0.25	0.71	0.44	0.81	0.50	0.93
Mining	2	0.1	0.25	0.71	0.22	0.69	0.22	0.69
Log-Likelihood			-16.96		-16.13		-16.02	

Topic models and analysis

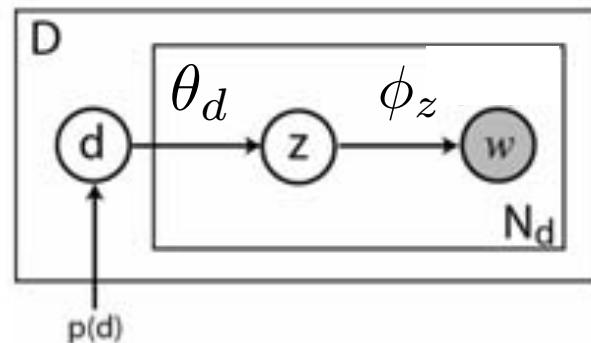
- Topic ≈ main idea discussed in text data
 - Theme/subject of a discussion or conversation
 - Different granularities (e.g., topic of a sentence, an article, etc.)
- Many applications require discovery of topics in text
 - What are Twitter users talking about today?
 - What are the current research topics in data mining? How are they different from those 5 years ago?
 - What do people like about the iPhone 6? What do they dislike?
 - What were the major topics debated in 2012 presidential election?

Lifecycle of topic and text data



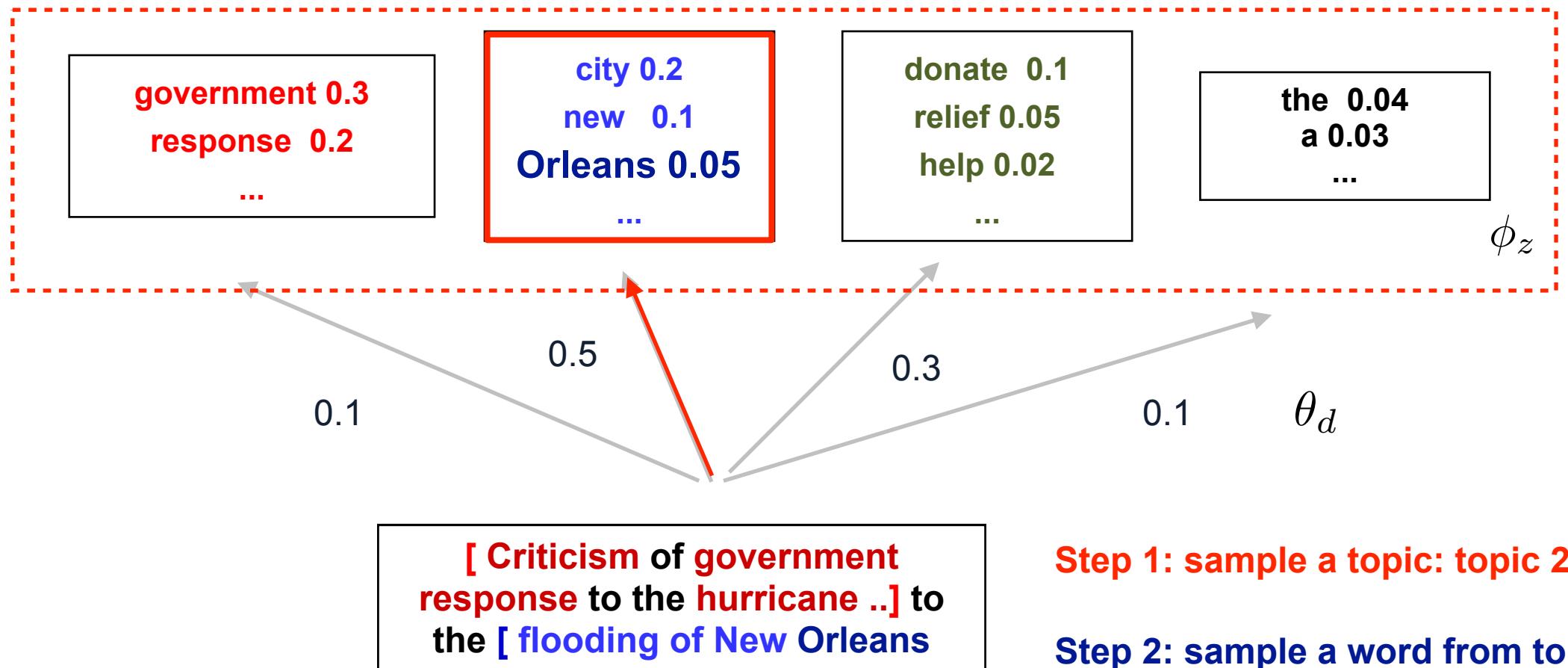
Probabilistic latent semantic analysis

- Assume documents are generated by sampling words from k latent topics
- For each document d:
 - For each token position i
 - Choose a topic $z \sim \text{Multinomial}(\theta_d)$
 - Choose a term $w \sim \text{Multinomial}(\phi_z)$



Review of PLSA

Phi $T \times V$, V : vocabulary size ~50,000, T : #topics, $T=20$
theta $D \times T$, D : #documents: ~10,000, T : #topics, $T=20$



Document as a Sample of Mixed Topics

Topic θ_1

government
0.3
response 0.2

Topic θ_2

...

Topic θ_k

city 0.2
new 0.1
orleans
0.05
...
donate 0.1
relief 0.05
help 0.02
...

Background θ_B

the 0.04
a 0.03
...

Blog article about “Hurricane Katrina”

[Criticism of government response to the hurricane primarily consisted of criticism of its response to the approach of the storm and its aftermath, specifically in the delayed response] to the [flooding of New Orleans. ... 80% of the 1.3 million residents of the greater New Orleans metropolitan area evacuated] ...[Over seventy countries pledged monetary donations or other assistance]. ...

proportion of topics

Probabilistic latent semantic analysis

$$p(d_i = w | \Phi, \theta_d) = \sum_{z=1}^T \phi_{z,w} \theta_{d,z}$$

$$p(\mathcal{W} | \Phi, \Theta)$$

$$= \prod_{d=1}^D \prod_{d_i=1}^{N_d} \sum_{z=1}^T \phi_{z,w} \theta_{d,z}$$

$$= \prod_{d=1}^D \prod_{w=1}^V \left(\sum_{z=1}^T \phi_{z,w} \theta_{d,z} \right)^{count(d,w)}$$

$$\arg \max_{\Phi, \Theta} [\log p(\mathcal{W} | \Phi, \Theta) + \sum_{d=1}^D \lambda_d (1 - \sum_{z=1}^T \theta_{(d,z)}) + \sum_{z=1}^T \sigma_k (1 - \sum_{w=1}^V \phi_{z,w})]$$

Probabilistic latent semantic analysis

$$\operatorname{argmax}_{\Theta, \Phi} \sum_d \sum_v \operatorname{count}(d, v) \sum_t^T R_{d, v, t} (\log \phi_{t, v} + \log \theta_{d, t})$$

$$+ \left(\sum_d \lambda_d \left(1 - \sum_t^T \theta_{d, t} \right) + \sum_t^T \sigma_t \left(1 - \sum_v \phi_{t, v} \right) \right)$$

$R_{d, v, t}$: latent variable's conditional probability, i.e., $p(t|v, d)$

- **M step: take the derivative of the above formula for $\theta_{d, t}, \phi_{t, v}$:**

$$\theta_{d, t} \propto \sum_v R_{d, v, t} \operatorname{count}(d, v)$$

$$\Rightarrow \theta_{d, t} = \frac{\sum_v R_{d, v, t} \operatorname{count}(d, v)}{\sum_t \sum_v R_{d, v, t} \operatorname{count}(d, v)}$$

$$\phi_{t, v} \propto \sum_d R_{d, v, t} \operatorname{count}(d, v)$$

$$\Rightarrow \phi_{t, v} = \frac{\sum_d R_{d, v, t} \operatorname{count}(d, v)}{\sum_v \sum_d R_{d, v, t} \operatorname{count}(d, v)}$$

- **E step: apply Bayes theorem:** $R_{d, v, t} \propto \phi_{d, t} \theta_{t, v} \Rightarrow R_{d, v, t} = \frac{\phi_{d, t} \theta_{t, v}}{\sum_t \phi_{d, t} \theta_{t, v}}$

Probabilistic latent semantic analysis: partially available labels

- Generalized topic modeling:
 - Each document can contain just one topic, e.g., short documents
 - That is, topic inference = topic classification
- If we already know the document tags for a part of the documents, does the partial labels help us make better predictions for the entire corpus? (**Homework 3**)
- Example: news tagging, StackOverflow question tagging

Probabilistic latent semantic analysis: partially available labels

Trends for you · Change

#RAF100

The RAF celebrates 100th anniversary

#TuesdayThoughts

@NWMCblog is Tweeting about this

#ThailandCaveRescue

237K Tweets

Thai Navy Seal

120K Tweets

All 12

All 12 boys and coach rescued from Thai cave

#IgniteB2B

1,850 Tweets

#WildBoars

4,934 Tweets

Science

159K Tweets

George Clooney

George Clooney injured in motorcycle accident in Italy

#NationalPinaColadaDay

1,870 Tweets

T Tagger News tags

1. Intel AMT Checker for Linux (github.com) **Security** **Linux**
116 points by laamalif 4 hours ago | 34 comments
2. Fwaf – Machine Learning Driven Web Application Firewall (fsecurify.com) **AI/Machine Learning** **Security** **Data Science**
46 points by Faizann20 7 hours ago | 9 comments
3. How to Spot a Spook (1974) (cryptome.org) **Politics**
30 points by mercer 3 hours ago | 6 comments
4. LittleTable: A Relational Time-Series Database at Cisco Meraki (acm.org) **Databases**
26 points by rodionos 7 hours ago | 2 comments
5. Wcry ransomware is reborn without its killswitch, starts spreading anew (boingboing.net) **Security**
34 points by rbanffy 4 hours ago | 4 comments
6. Using Deep Learning at Scale in Twitter's Timelines (twitter.com) **AI/Machine Learning**
39 points by hunglee2 8 hours ago | 23 comments
7. Cyberattacks in 12 Nations Said to Use Leaked N.S.A. Hacking Tool (nytimes.com) **Security**
1200 points by ghosh 21 hours ago | 468 comments
8. What it's like to be in the most automated job in the United States (qz.com) **AI/Machine Learning**
28 points by billifuduo 5 hours ago | 17 comments
9. The Right to Read (1997) (gnu.org) **Blockchain**
196 points by tobysullivan 21 hours ago | 62 comments
10. Spends 10 Years Mastering Microsoft Paint to Illustrate His Book (boredpanda.com) **Microsoft**
59 points by prncpinto 8 hours ago | 12 comments
11. Visual Studio 2017 now fully supports Python and Django (visualstudio.com) **Python** **Microsoft**
91 points by vanflymen 14 hours ago | 38 comments
12. Your tl;dr by an ai: a deep reinforced model for abstractive summarization (metamind.io) **AI/Machine Learning**
99 points by etiam 18 hours ago | 18 comments
13. Lessons scaling from 10 to 20 people (josephwalla.com) **Startups**
60 points by gkop 12 hours ago | 15 comments
14. Happy nations don't focus on growth (bloomberg.com) **Politics**
133 points by smollett 21 hours ago | 45 comments
15. Ublock vs. ublock origin (reddit.com) **Web Development**
39 points by based2 4 hours ago | 12 comments
16. Rejection Letter (antipope.org) **Security**
475 points by cstruss 2 hours ago | 63 comments
17. Video Solves Mystery of How Narwhals Use Their Tusks (nationalgeographic.com) **Science**
105 points by clouddrover 8 hours ago | 7 comments
18. Netflix confirms it is blocking rooted/unlocked Android devices (androidpolice.com) **Mobile**
184 points by msq 14 hours ago | 167 comments
19. Fuzzing Irrsi (irssi.org) **Security**
136 points by jbisch 2 hours ago | 13 comments

Homework 3

- Suppose each document has only 1 topic. We have two document set: S1 (100 documents) contains all the tagged documents; S2 (10,000 documents) contains all the untagged documents. Each tag is a topic, there are only 2 topics
- (Part 1): derive the EM algorithm using pLSA that maximizes the probability of the observed document, given the known topics from S1
- (Part 2): implement your EM algorithm, output the predicted topic for each document in S2

CS 589 Fall 2021 Lecture 6

Latent Dirichlet allocation

**Monday 6:30-9:00
Babbio 122**

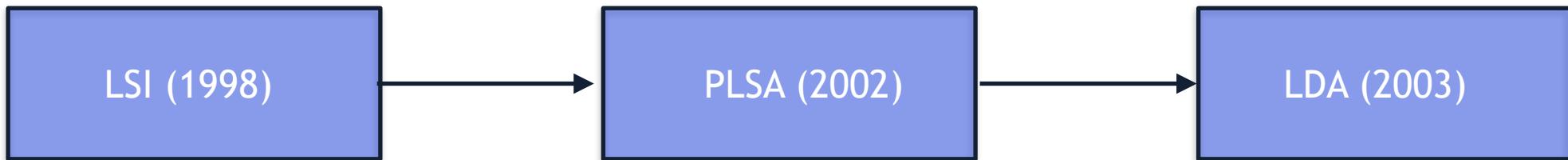
All zoom links in Canvas



photo: <https://www.scubedstudios.com/information-retrieval/>

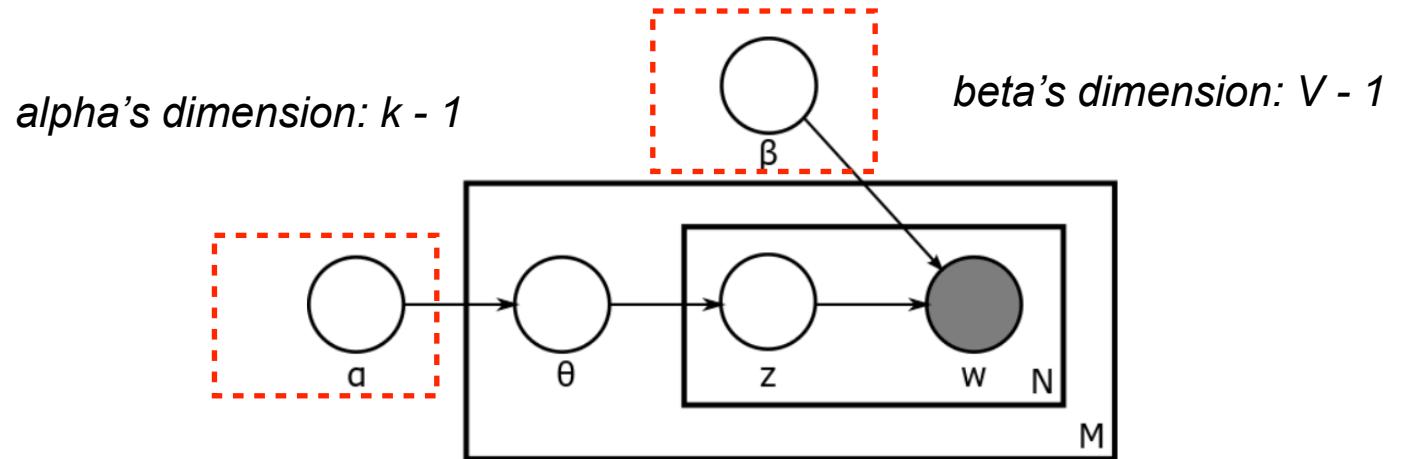
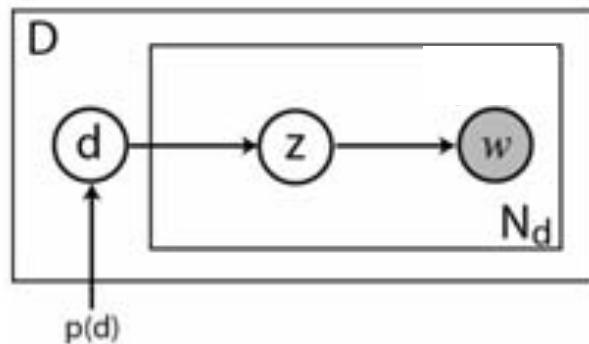
Latent Dirichlet allocation [Blei et al. 03]

- A Bayesian topic model by considering the prior distribution of $p(t|d)$
 - pLSA assumes $\theta_{d,t}$ is a variable by itself
 - LDA assumes $\theta_{d,t}$ is drawn from yet another distribution (i.e., "distribution of distribution")



Latent Dirichlet allocation [Blei et al. 2003]

- Adding a prior *hyper parameter* α to PLSA model
 - The document-topic probability is **randomly sampled** from the same prior distribution
 - That is, when computing the probability, we include $p(\theta_{d,t}|\alpha)$ in the model



Bayes' rules

Chain rule: joint distribution

$$P(A, B) = P(A \cap B) = P(A|B)P(B) = P(B|A)P(A)$$

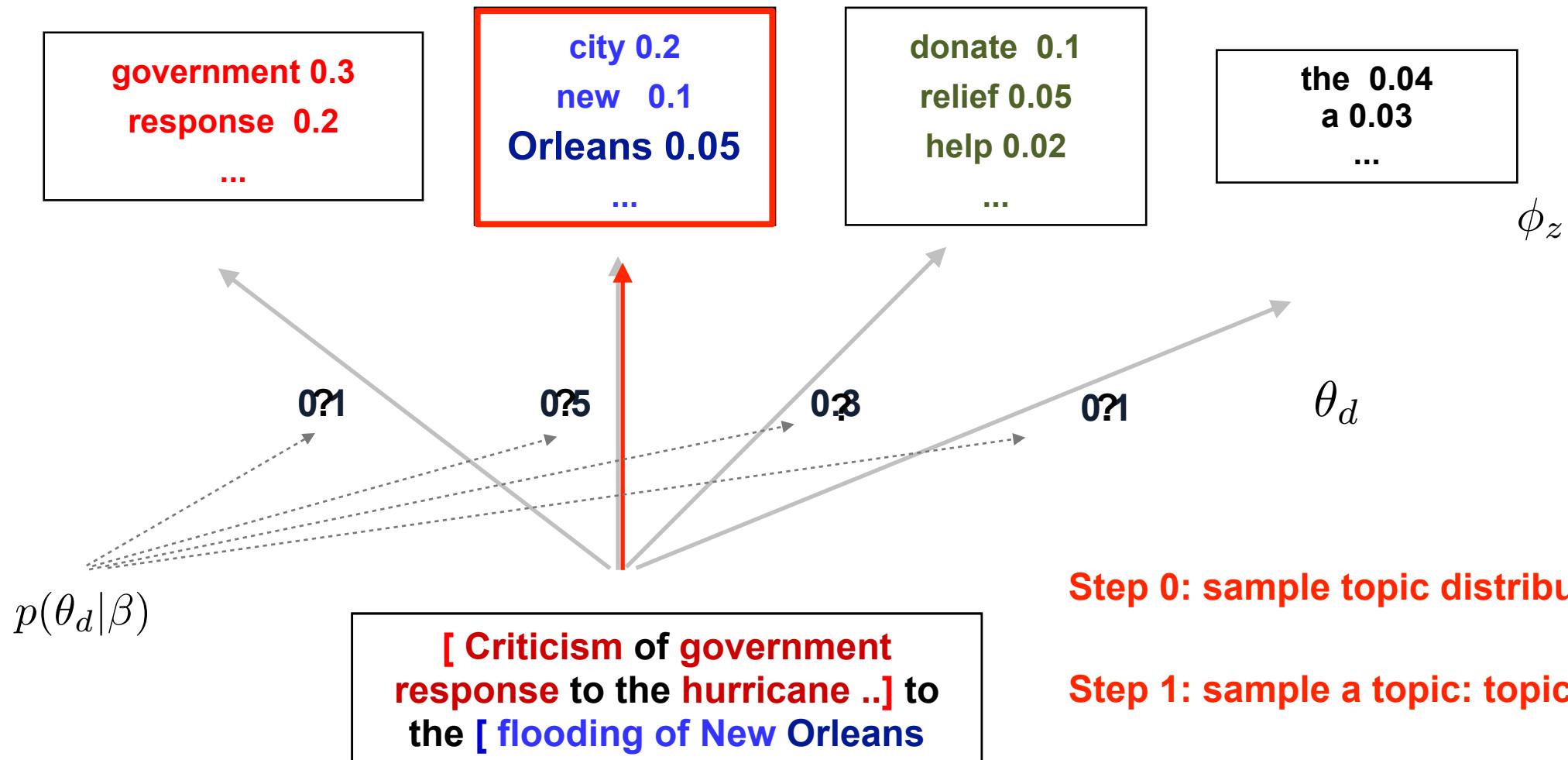
Bayes' rule:

posterior likelihood prior

$$P(A|B) = \frac{P(B|A)P(A)}{P(B)} = \left[\frac{P(B|A)}{\sum_{X \in \{A, \bar{A}\}} P(B|X)P(X)} \right] P(A)$$

$$P(A|B) \propto P(B|A)P(A)$$

PLSA \rightarrow LDA



Maximum a Posterior (MAP) estimation

- Maximum likelihood estimation:

$$\theta_{MLE} = \arg \max_{\theta} p(x_1, \dots, x_n | \theta)$$

- Choose the parameter with the maximum posterior probability

$$\begin{aligned}\theta_{MAP} &= \arg \max_{\theta} p(\theta | x_1, \dots, x_n) \\ &= \arg \max_{\theta} p(x_1, \dots, x_n | \theta) p(\theta)\end{aligned}$$

- LDA is the MAP version of pLSA (MLE)

PLSA -> LDA

- PLSA:

$$p_d(w \mid \{\theta_d\}, \{\phi_z\}) = \sum_z \theta_{d,z} \phi_{z,w}$$

$$\log p(D \mid \{\phi_z\}, \{\theta_d\}) = \sum_{d \in D} \sum_{w \in V} c(w, d) \log \left[\sum_z \theta_{d,z} \phi_{z,w} \right]$$

- LDA:

$$p_d(w \mid \{\theta_d\}, \{\phi_z\}) = \sum_z \theta_{d,z} \phi_{z,w}$$

$$\log p(d \mid \alpha, \{\phi_z\}) = \int \sum_{w \in V} c(w, d) \log \left[\sum_z \theta_{d,z} \phi_{z,w} \right] p(\theta_d \mid \alpha) d\theta_d$$

marginalized probability:

$$\log p(D \mid \alpha, \beta) = \int \sum_{d \in D} \log p(d \mid \alpha, \{\phi_z\}) \prod_z p(\phi_z \mid \beta) d\phi_1 \dots d\phi_k$$

Solving Maximum a Posteriori inference for LDA

- Exact inference is intractable:

$$p(Z, \Phi, \Theta | D, \alpha, \beta) = \frac{p(D, Z, \Phi, \Theta | \alpha, \beta)}{p(D | \alpha, \beta)}$$

$$\log p(D | \alpha, \beta) = \int \sum_{d \in D} \log p(d | \alpha, \{\phi_z\}) \prod_{z=1}^k p(\phi_z | \beta) d\phi_1 \dots d\phi_k$$

- Equation (1) is computationally intractable due to the coupling of beta and phi in the denominator

Variational inference

- Key idea: use a **surrogate distribution** to approximate the posterior distribution of latent variables
 - Surrogate distribution is simpler to estimate than the true posterior
- Goal: finding the “best” surrogate distribution from a certain parametric family by **minimizing** the KL-divergence between the surrogate function (Q) to the true posterior (P)

Evidence lower bound (ELBO)

- Given that $p(Z|D) = \frac{p(D, Z)}{p(D)}$, we want to minimize the KL divergence between $Q(Z)$ and $p(Z|D)$:

$$\begin{aligned} D_{\text{KL}}(q\|p) &= \int_Z q(Z) \left[\log \frac{q(Z)}{p(Z, D)} + \log p(D) \right] \\ &= \int_Z q(Z)[\log q(Z) - \log p(Z, D)] + \int_Z q(Z)[\log p(D)] \\ &= \int_Z q(Z)[\log q(Z) - \log p(Z, D)] + \log p(D) \\ \Rightarrow \log p(D) &= D_{\text{KL}}(q\|p) - \mathbb{E}_q[\log q(Z) - \log p(Z, D)] = D_{\text{KL}}(q\|p) + \mathcal{L}(q) \end{aligned}$$

Evidence lower bound (ELBO)

- $p(D)$ does not rely on the latent variable Z :

$$\Rightarrow \log p(D) = D_{\text{KL}}(q\|p) - \mathbb{E}_q[\log q(Z) - \log p(Z, D)] = D_{\text{KL}}(q\|p) + \mathcal{L}(q)$$

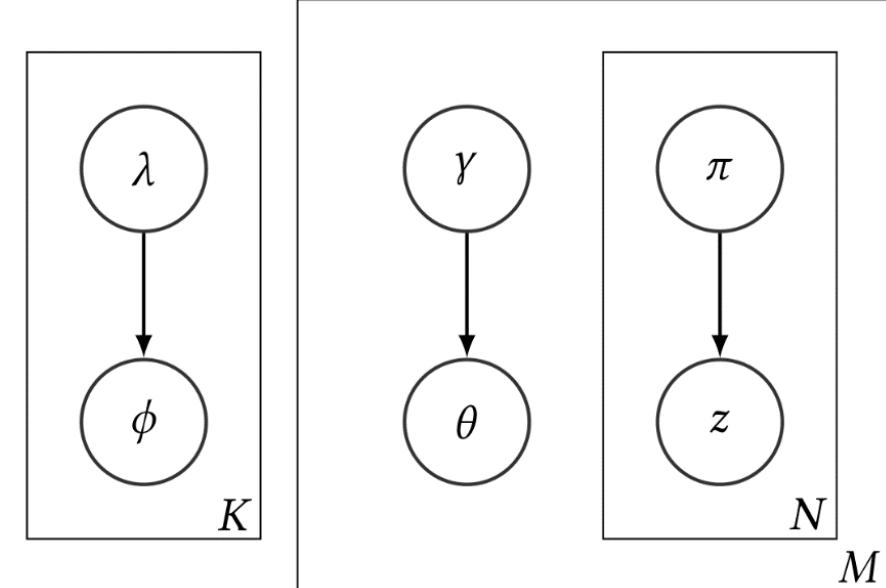
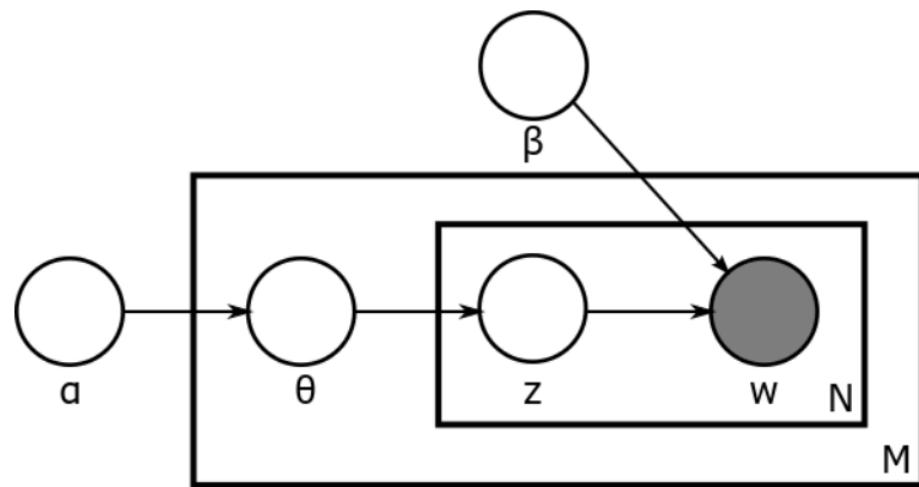
- Therefore minimizing KL divergence = maximizing $\mathcal{L}(q)$, which we call evidence lower bound, as KL divergence is non-negative
- The lower bound is more tractable to optimize than the evidence:

$$L(q) = \int_Z q(Z)[\log q(Z) - \log p(Z, D)]$$

LDA variational inference

- Assumption: q belongs to a family of distribution which can be factorized

$$q(\theta, \mathbf{z} \mid \gamma, \phi) = q(\theta \mid \gamma) \prod_{n=1}^N q(z_n \mid \phi_n)$$



LDA variational inference

$$L(q) = \int_Z q(Z)[\log q(Z) - \log p(Z, D)]$$

$$p(Z, \Phi, \Theta | D, \alpha, \beta) = \frac{p(D, Z, \Phi, \Theta | \alpha, \beta)}{p(D | \alpha, \beta)}$$

- ELBO of LDA under the factorization assumption:

$$\begin{aligned} L(\gamma, \phi; \alpha, \beta) &= \mathbb{E}_q[\log p(\theta | \alpha)] + \mathbb{E}_q[\log p(\mathbf{z} | \theta)] + \mathbb{E}_q[\log p(\mathbf{w} | \mathbf{z}, \beta)] \\ &\quad - \mathbb{E}_q[\log q(\theta)] - \mathbb{E}_q[\log q(\mathbf{z})] \end{aligned}$$

$$\begin{aligned} L(\gamma, \phi; \alpha, \beta) &= \log \Gamma \left(\sum_{j=1}^k \alpha_j \right) - \sum_{i=1}^k \log \Gamma (\alpha_i) + \sum_{i=1}^k (\alpha_i - 1) \left(\Psi (\gamma_i) - \Psi \left(\sum_{j=1}^k \gamma_j \right) \right) \\ &\quad + \sum_{n=1}^N \sum_{i=1}^k \phi_{ni} \left(\Psi (\gamma_i) - \Psi \left(\sum_{j=1}^k \gamma_j \right) \right) \\ &\quad + \sum_{n=1}^N \sum_{i=1}^k \sum_{j=1}^V \phi_{ni} w_n^j \log \beta_{ij} \\ &\quad - \log \Gamma \left(\sum_{j=1}^k \gamma_j \right) + \sum_{i=1}^k \log \Gamma (\gamma_i) - \sum_{i=1}^k (\gamma_i - 1) \left(\Psi (\gamma_i) - \Psi \left(\sum_{j=1}^k \gamma_j \right) \right) \\ &\quad - \sum_{n=1}^N \sum_{i=1}^k \phi_{ni} \log \phi_{ni} \end{aligned}$$

← *θ and z integrated out*

LDA variational inference

- LDA variational inference algorithm: optimizing γ , φ , β
1. Randomly initialize variational parameters (can't be uniform)
 2. For each iteration:
 1. For each document, update γ and φ
 2. For corpus, update β
 3. Compute L for diagnostics
 3. Return **expectation of variational parameters** for solution to latent variables

$$\phi_{ni} \propto \beta_{iv} \exp \left(\psi(\gamma_i) - \Psi \left(\sum_j \gamma_j \right) \right)$$
$$E_q [\log (\theta_i) | \gamma] = \Psi(\gamma_i) - \Psi \left(\sum_{j=1}^k \gamma_j \right)$$

Pros and Cons of Variational inference

- Pros
 - Deterministic algorithm - easy to tell when converged
 - Embarrassingly parallelizable
- Cons:
 - Speed: make many calls to transcendental functions (no close form solution)
 - Quality is questionable due to the factorization assumption
 - Memory usage: requires $O(MNK)$ to store the per-token variational distributions

Conjugate prior distribution

- That is, the **posterior distributions** are in the same probability distribution family as the prior probability distribution
- The prior is called the **conjugate prior** of the likelihood probability
- Conjugate prior makes denominator easy to compute because you can integrate out theta
- e.g., Beta distribution is the conjugate prior of Bernoulli distribution

Bernoulli: $p(x | \theta) = \theta^x(1 - \theta)^{1-x}.$

Beta: $p(\theta | \alpha) = K(\alpha)\theta^{\alpha_1-1}(1 - \theta)^{\alpha_2-1}.$

$$\begin{aligned} K(\alpha) &= \left(\int \theta^{\alpha_1-1}(1 - \theta)^{\alpha_2-1} d\theta \right)^{-1} \\ &= \frac{\Gamma(\alpha_1 + \alpha_2)}{\Gamma(\alpha_1)\Gamma(\alpha_2)} \end{aligned}$$

Collapsed variational inference [Teh et al. 07]

- If your priors are conjugate, they can be integrated out!
- Let n_{dzv} be the number of times that word v in document d has topic z
- Let $n_{.zv}$ be the number of times word v in any document has topic z

$$p(\mathbf{z}, \mathbf{x} \mid \alpha, \beta) = \prod_d \frac{\Gamma(k\alpha)}{\Gamma(k\alpha + n_{d..})} \prod_z \frac{\Gamma(\alpha + n_{dz.})}{\Gamma(\alpha)} \prod_z \frac{\Gamma(V\beta)}{\Gamma(V\beta + n_{.z.})} \prod_v \frac{\Gamma(\beta + n_{.zv})}{\Gamma(\beta)}$$

No theta or phi in the above formula

Collapsed variational inference [Teh et al. 07]

- How do we get back phi and theta?

$$\theta_{dz} = \frac{\alpha + n_{dz.}}{K\alpha + n_{d..}}$$

$$\phi_{zv} = \frac{\beta + n_{.zv}}{W\beta + n_{.z.}}$$

- How do we get n_{dzv} ?
 - Gibbs sampling: sample values of Z, count from those examples

Summary

- Expectation maximization
 - maximizes the lower bound of the likelihood function
 - guarantee to improve the likelihood
 - may stuck at local optimal
- Know how to derive EM algorithm for coin topic and pLSA
- LDA is the Bayesian version of pLSA
 - Direct inference is intractable
 - Variational inference (leveraging ELBO)
 - Markov chain Monte Carlo (Gibbs sampling)

Resources for LDA

- Mallet: <http://mallet.cs.umass.edu/>

```
bin/mallet import-dir --input /data/topic-input --output topic-input.mallet \
--keep-sequence --remove-stopwords
```

- gensim:

```
>>> from gensim.test.utils import common_texts
>>> from gensim.corpora.dictionary import Dictionary
>>>
>>> # Create a corpus from a list of texts
>>> common_dictionary = Dictionary(common_texts)
>>> common_corpus = [common_dictionary.doc2bow(text) for text in common_texts]
>>>
>>> # Train the model on the corpus.
>>> lda = LdaModel(common_corpus, num_topics=10)
```