

# Inference and Representation

David Sontag

New York University

Lecture 6, Oct. 14, 2015

# Approximate marginal inference

- Given the joint  $p(x_1, \dots, x_n)$  represented as a graphical model, how do we perform **marginal inference**, e.g. to compute  $p(x_1 | e)$ ?
- We showed in Lecture 4 that doing this exactly is NP-hard
- Nearly all *approximate inference* algorithms are either:
  - 1 **Monte-carlo methods (e.g., Gibbs sampling, likelihood reweighting, MCMC)**
  - 2 Variational algorithms (e.g., mean-field, loopy belief propagation)

---

## Algorithm 12.1 Forward Sampling in a Bayesian network

---

**Procedure** Forward-Sample (  
     $\mathcal{B}$  // Bayesian network over  $\mathcal{X}$   
)

- 1 Let  $X_1, \dots, X_n$  be a topological ordering of  $\mathcal{X}$
  - 2 **for**  $i = 1, \dots, n$
  - 3      $\mathbf{u}_i \leftarrow \mathbf{x}\langle \text{Pa}_{X_i} \rangle$  // Assignment to  $\text{Pa}_{X_i}$  in  $x_1, \dots, x_{i-1}$
  - 4     Sample  $x_i$  from  $P(X_i \mid \mathbf{u}_i)$
  - 5 **return**  $(x_1, \dots, x_n)$
- 

(Koller & Friedman, *Probabilistic Graphical Models*, MIT Press 2009)

# Monte-Carlo algorithms

- Given a joint distribution  $p(x_1, \dots, x_n)$ , how do we compute marginals?

$$\begin{aligned} p[X_1 = x_1] &= E_{\mathbf{x} \sim p}[f(\mathbf{x})], \text{ where } f(\mathbf{x}) = 1[X_1 = x_1] \\ &= \sum_{\mathbf{x}} p(\mathbf{x}) f(\mathbf{x}). \end{aligned}$$

- Rather than explicitly enumerating *all* assignments, consider the following Monte-Carlo estimate of the expectation:

$$\begin{aligned} \mathbf{x}^1 &\sim p(\mathbf{x}) \\ \mathbf{x}^2 &\sim p(\mathbf{x}) \\ &\vdots \\ \mathbf{x}^M &\sim p(\mathbf{x}) \end{aligned}$$

- Then, our *estimate* is  $\hat{E}_p[f(\mathbf{x})] = \frac{1}{M} \sum_{m=1}^M f(\mathbf{x}^m)$ . **How good is it?**

# Monte-Carlo algorithms

- Let  $\mathcal{D} = \{\mathbf{x}^1, \dots, \mathbf{x}^M\}$ . Since  $\mathcal{D}$  was drawn randomly from  $p(\mathbf{x})$ , the estimate is itself a random variable
- The estimate is *unbiased* because

$$\begin{aligned} E_{\mathbf{x}^1, \dots, \mathbf{x}^M \sim p(\mathbf{x})} \left[ \hat{E}[f(\mathbf{x})] \right] &= E_{\mathbf{x}^1, \dots, \mathbf{x}^M \sim p(\mathbf{x})} \left[ \frac{1}{M} \sum_{m=1}^M f(\mathbf{x}^m) \right] \\ &= \frac{1}{M} \sum_{m=1}^M E_{\mathbf{x}^m \sim p(\mathbf{x})} \left[ f(\mathbf{x}^m) \right] \\ &= E_{\mathbf{x} \sim p(\mathbf{x})} [f(\mathbf{x})]. \end{aligned}$$

- How quickly does the estimate converge to the true expectation?

# Law of large numbers

- There are two general results we can use, depending on whether we care about additive or multiplicative error
- **Hoeffding bound** says that:

$$\Pr_{\mathcal{D} \sim p(\mathbf{x})} \left[ E_p[f(\mathbf{x})] - \epsilon \leq \hat{E}_{\mathcal{D}}[f(\mathbf{x})] \leq E_p[f(\mathbf{x})] + \epsilon \right] \geq 1 - 2e^{-2M\epsilon^2}$$

- **Chernoff bound** says that (assuming  $f(\mathbf{x}) \in [0, 1]$ ):

$$\Pr_{\mathcal{D} \sim p(\mathbf{x})} \left[ E_p[f(\mathbf{x})](1 - \epsilon) \leq \hat{E}_{\mathcal{D}}[f(\mathbf{x})] \leq E_p[f(\mathbf{x})](1 + \epsilon) \right] \geq 1 - 2e^{-\frac{M\epsilon^2}{3} E_p[f(\mathbf{x})]}$$

- Estimating *single-variable* marginals for a BN is easy: just forward sample!
- What about computing *conditional* queries such as  $p(\mathbf{X} = \mathbf{x} \mid \mathbf{E} = \mathbf{e})$ ?
- Computing denominator of  $p(\mathbf{X} = \mathbf{x}, \mathbf{E} = \mathbf{e})/p(\mathbf{E} = \mathbf{e})$  needs  $\Omega(1/p(\mathbf{E} = \mathbf{e}))$  samples, by Chernoff bound.

# “Normalized” Importance Sampling

- If we could instead directly sample from  $p(\mathbf{X} \mid \mathbf{E} = \mathbf{e})$ , we would be in business – but this is hard!
- For the same reason, sampling from an undirected graphical model  $p(\mathbf{x}) = \frac{1}{Z} \prod_{c \in C} \phi_c(\mathbf{x}_c)$  – even without evidence – is hard, because we don't know  $Z$
- Suppose we instead had a simpler-to-sample-from distribution  $q(\mathbf{x})$ , called the “proposal distribution”
- Let  $\tilde{p}(\mathbf{x})$  be an unnormalized version of the distribution, e.g.

$$\tilde{p}(\mathbf{x}) = p(\mathbf{x}, E = \mathbf{e}) \quad (\text{BN with evidence})$$

$$\tilde{p}(\mathbf{x}) = \prod_{c \in C} \phi_c(\mathbf{x}_c) \quad (\text{MRF})$$

Note that we can efficiently *evaluate*  $\tilde{p}(\mathbf{x})$  for any  $\mathbf{x}$

# “Normalized” Importance Sampling

- Consider the following estimate (now using  $\mathbf{x}^1, \dots, \mathbf{x}^M \sim q(\mathbf{x})$ ):

$$\hat{E}_{\mathcal{D}}[f(x)] = \frac{\frac{1}{M} \sum_{m=1}^M f(\mathbf{x}^m) \tilde{w}(\mathbf{x}^m)}{\frac{1}{M} \sum_{m=1}^M \tilde{w}(\mathbf{x}^m)}, \quad \text{where } \tilde{w}(\mathbf{x}) = \frac{\tilde{p}(\mathbf{x})}{q(\mathbf{x})}$$

- This is *not* an unbiased estimate! E.g., for  $M = 1$ , we have

$$\begin{aligned} E_{\mathbf{x}^1 \sim q(\mathbf{x})} \left[ \hat{E}_{\mathcal{D}}[f(x)] \right] &= E_{\mathbf{x}^1 \sim q(\mathbf{x})} \left[ \frac{f(\mathbf{x}^1) \tilde{w}(\mathbf{x}^1)}{\tilde{w}(\mathbf{x}^1)} \right] = E_{\mathbf{x} \sim q(\mathbf{x})} [f(\mathbf{x})] \\ &\neq E_{\mathbf{x} \sim p(\mathbf{x})} [f(\mathbf{x})]. \end{aligned}$$

- However, the estimate is asymptotically correct (i.e., as  $M \rightarrow \infty$ )



# “Normalized” Importance Sampling

- Consider the following estimate (now using  $\mathbf{x}^1, \dots, \mathbf{x}^M \sim q(\mathbf{x})$ ):

$$\hat{E}_{\mathcal{D}}[f(\mathbf{x})] = \frac{\frac{1}{M} \sum_{m=1}^M f(\mathbf{x}^m) \tilde{w}(\mathbf{x}^m)}{\frac{1}{M} \sum_{m=1}^M \tilde{w}(\mathbf{x}^m)}, \quad \text{where } \tilde{w}(\mathbf{x}) = \frac{\tilde{p}(\mathbf{x})}{q(\mathbf{x})}$$

- Letting  $\tilde{p}(\mathbf{x}) = p(\mathbf{x})Z$ , the expectation of the numerator is:

$$\begin{aligned} E_{\mathcal{D} \sim q(\mathbf{x})} \left[ \frac{1}{M} \sum_{m=1}^M f(\mathbf{x}^m) \tilde{w}(\mathbf{x}^m) \right] &= \frac{1}{M} \sum_{m=1}^M E_{\mathbf{x}^m \sim q(\mathbf{x})} [f(\mathbf{x}^m) \tilde{w}(\mathbf{x}^m)] \\ &= \frac{1}{M} \sum_{m=1}^M \sum_{\mathbf{x}} q(\mathbf{x}) \left[ f(\mathbf{x}) \frac{\tilde{p}(\mathbf{x})}{q(\mathbf{x})} \right] \\ &= \frac{1}{M} \sum_{m=1}^M \sum_{\mathbf{x}} \tilde{p}(\mathbf{x}) f(\mathbf{x}) = Z E_p[f(\mathbf{x})]. \end{aligned}$$

# “Normalized” Importance Sampling

- Consider the following estimate (now using  $\mathbf{x}^1, \dots, \mathbf{x}^M \sim q(\mathbf{x})$ ):

$$\hat{E}_{\mathcal{D}}[f(\mathbf{x})] = \frac{\frac{1}{M} \sum_{m=1}^M f(\mathbf{x}^m) \tilde{w}(\mathbf{x}^m)}{\frac{1}{M} \sum_{m=1}^M \tilde{w}(\mathbf{x}^m)}, \quad \text{where } \tilde{w}(\mathbf{x}) = \frac{\tilde{p}(\mathbf{x})}{q(\mathbf{x})}$$

- Letting  $\tilde{p}(\mathbf{x}) = p(\mathbf{x})Z$ , the expectation of the numerator is  $ZE_p[f(\mathbf{x})]$ .
- The expectation of the denominator is  $Z$ !

$$\begin{aligned} E_{\mathcal{D} \sim q(\mathbf{x})} \left[ \frac{1}{M} \sum_{m=1}^M \tilde{w}(\mathbf{x}^m) \right] &= \frac{1}{M} \sum_{m=1}^M E_{\mathbf{x}^m \sim q(\mathbf{x})} [\tilde{w}(\mathbf{x}^m)] \\ &= \frac{1}{M} \sum_{m=1}^M \sum_{\mathbf{x}} q(\mathbf{x}) \left[ \frac{\tilde{p}(\mathbf{x})}{q(\mathbf{x})} \right] \\ &= \frac{1}{M} \sum_{m=1}^M \sum_{\mathbf{x}} \tilde{p}(\mathbf{x}) = Z. \end{aligned}$$

# Likelihood weighting

- What should we use for  $q(\mathbf{x})$ ? For a Bayesian network, we can sample from the latent variables, keeping the evidence fixed

---

**Algorithm 12.2 Likelihood-weighted particle generation**

---

```
Procedure LW-Sample (  
     $\mathcal{B}$ , // Bayesian network over  $\mathcal{X}$   
     $Z = \mathbf{z}$  // Event in the network  
)  
1 Let  $X_1, \dots, X_n$  be a topological ordering of  $\mathcal{X}$   
2  $w \leftarrow 1$   
3 for  $i = 1, \dots, n$   
4    $\mathbf{u}_i \leftarrow \mathbf{x}(\text{Pa}_{X_i})$  // Assignment to  $\text{Pa}_{X_i}$  in  $x_1, \dots, x_{i-1}$   
5   if  $X_i \notin Z$  then  
6     Sample  $x_i$  from  $P(X_i | \mathbf{u}_i)$   
7   else  
8      $x_i \leftarrow \mathbf{z}(X_i)$  // Assignment to  $X_i$  in  $\mathbf{z}$   
9      $w \leftarrow w \cdot P(x_i | \mathbf{u}_i)$  // Multiply weight by probability of desired value  
10  return  $(x_1, \dots, x_n), w$ 
```

---

(Koller & Friedman, *Probabilistic Graphical Models*, MIT Press 2009)

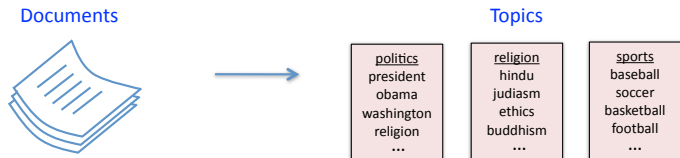
- Corresponds to importance sampling using:

$$q(\mathbf{x}) = \prod_{t \notin \mathbf{E}} p(x_t | \mathbf{x}_{pa(t)}) \prod_{t \in \mathbf{E}} 1[x_t = \mathbf{e}_t], \text{ so } \tilde{w}(\mathbf{x}) = \frac{\tilde{p}(\mathbf{x})}{q(\mathbf{x})} = \prod_{t \in \mathbf{E}} p(x_t | \mathbf{x}_{pa(t)}).$$

Problem Set 4 will explore Gibbs sampling for topic models

# Latent Dirichlet allocation (LDA)

- **Topic models** are powerful tools for exploring large data sets and for making inferences about the content of documents



- Many applications in information retrieval, document summarization, and classification



- LDA is one of the simplest and most widely used topic models

# Generative model for a document in LDA

- 1 Sample the document's **topic distribution**  $\theta$  (aka topic vector)

$$\theta \sim \text{Dirichlet}(\alpha_1:T)$$

where the  $\{\alpha_t\}_{t=1}^T$  are fixed hyperparameters. Thus  $\theta$  is a distribution over  $T$  topics with mean  $\theta_t = \alpha_t / \sum_{t'} \alpha_{t'}$

- 2 For  $i = 1$  to  $N$ , sample the **topic**  $z_i$  of the  $i$ 'th word

$$z_i | \theta \sim \theta$$

- 3 ... and then sample the actual **word**  $w_i$  from the  $z_i$ 'th topic

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

where  $\{\beta_t\}_{t=1}^T$  are the *topics* (a fixed collection of distributions on words)

# Generative model for a document in LDA

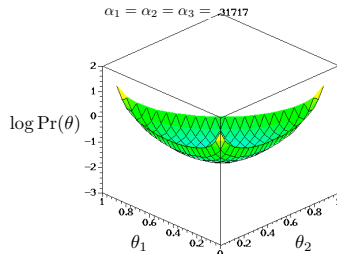
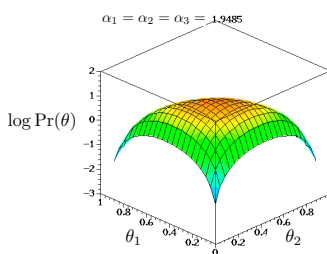
- 1 Sample the document's **topic distribution**  $\theta$  (aka topic vector)

$$\theta \sim \text{Dirichlet}(\alpha_{1:T})$$

where the  $\{\alpha_t\}_{t=1}^T$  are hyperparameters. The Dirichlet density, defined over  $\Delta = \{\vec{\theta} \in \mathbb{R}^T : \forall t \theta_t \geq 0, \sum_{t=1}^T \theta_t = 1\}$ , is:

$$p(\theta_1, \dots, \theta_T) \propto \prod_{t=1}^T \theta_t^{\alpha_t - 1}$$

For example, for  $T=3$  ( $\theta_3 = 1 - \theta_1 - \theta_2$ ):

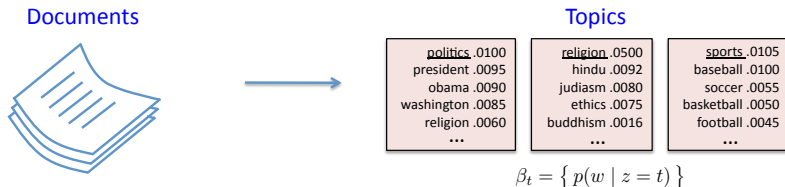


# Generative model for a document in LDA

- ③ ... and then sample the actual **word**  $w_i$  from the  $z_i$ 'th topic

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

where  $\{\beta_t\}_{t=1}^T$  are the *topics* (a fixed collection of distributions on words)





# Example of using LDA

 $\beta_1$ 

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
...	

 $\beta_T$ 

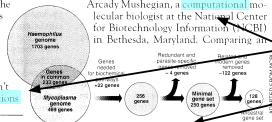
data	0.02
number	0.02
computer	0.01
...	

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 at 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.

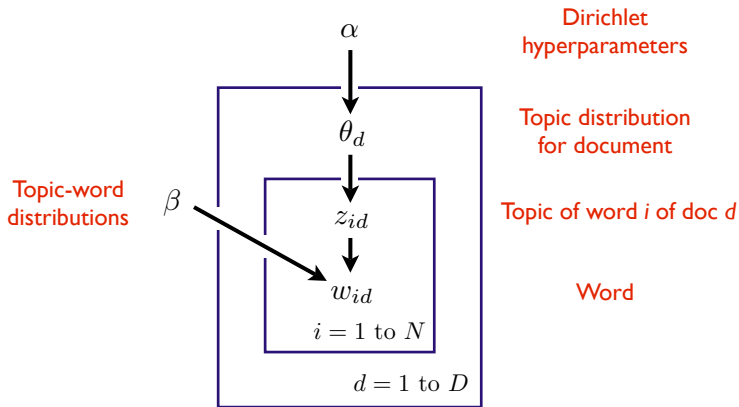
SCIENCE • VOL. 272 • 24 MAY 1996

Topic proportions and assignments

 $z_{1d}$ 
 $\theta_d$ 
 $z_{Nd}$ 

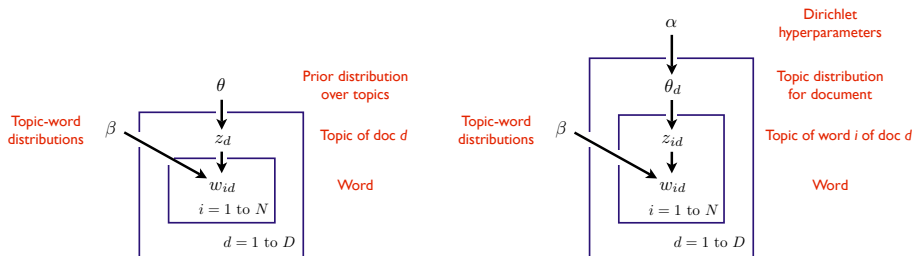
(Blei, *Introduction to Probabilistic Topic Models*, 2011)

# “Plate” notation for LDA model



Variables within a plate are replicated in a conditionally independent manner

# Comparison of mixture and admixture models



- Model on left is a **mixture model**
  - Called *multinomial* naive Bayes (a word can appear multiple times)
  - Document is generated from a single topic
- Model on right (LDA) is an **admixture model**
  - Document is generated from a distribution over topics