All of Graphical Models

Xiaojin Zhu

Department of Computer Sciences University of Wisconsin–Madison, USA

Tutorial at ICMLA 2011

◆□▶ ◆□▶ ◆三▶ ◆三▶ 三三 のへぐ

- Given $GM = joint distribution p(x_1, \ldots, x_n)$
- ▶ Do inference = $p(X_Q | X_E)$, in general $X_Q \cup X_E \subset \{x_1 \dots x_n\}$
- If $p(x_1, \ldots, x_n)$ not given, estimate it from data

・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・

Outline

Life without Graphical Models

Representation

Directed Graphical Models (Bayesian Networks) Undirected Graphical Models (Markov Random Fields)

(日)、(型)、(E)、(E)、(E)、(Q)

Inference

Exact Inference Markov Chain Monte Carlo Variational Inference Loopy Belief Propagation Mean Field Algorithm Exponential Family Maximizing Problems

Parameter Learning

Structure Learning

Outline

Life without Graphical Models

Representation

Directed Graphical Models (Bayesian Networks) Undirected Graphical Models (Markov Random Fields)

▲ロト ▲帰ト ▲ヨト ▲ヨト - ヨ - の々ぐ

Inference

Exact Inference Markov Chain Monte Carlo Variational Inference Loopy Belief Propagation Mean Field Algorithm Exponential Family Maximizing Problems

Parameter Learning

Structure Learning

Life without Graphical Models

... is fine mathematically:

- The universe is reduced to a set of random variables x_1, \ldots, x_n
 - e.g., x_1, \ldots, x_{n-1} can be the discrete or continuous features
 - e.g., $x_n \equiv y$ can be the discrete class label
- ▶ The joint $p(x_1, ..., x_n)$ completely describes how the universe works
- ▶ "Machine learning": estimate $p(x_1, ..., x_n)$ from training data $X^{(1)}, ..., X^{(N)}$, where $X^{(i)} = (x_1^{(i)}, ..., x_n^{(i)})$
- ▶ "Prediction": $y^* = \operatorname{argmax} p(x_n \mid x_1^*, \dots, x_n^*)$, a.k.a. inference
 - by the definition of conditional probability

$$p(x_n \mid x_1^*, \dots, x_n^*) = \frac{p(x_1^*, \dots, x_n^*, x_n)}{\sum_v p(x_1^*, \dots, x_n^*, x_n = v)}$$

・ロト・西ト・モン・モー うへぐ

Life without graphical models is just fine

◆□▶ ◆□▶ ◆臣▶ ◆臣▶ 臣 のへぐ

So why are we still here?

- Given $GM = joint distribution p(x_1, \ldots, x_n)$
 - exponential naïve storage $(2^n \text{ for binary r.v.})$
 - hard to interpret (conditional independence)
- ▶ Do inference = $p(X_Q | X_E)$, in general $X_Q \cup X_E \subset \{x_1 \dots x_n\}$
 - Often can't do it computationally
- If $p(x_1, \ldots, x_n)$ not given, estimate it from data

Can't do it either

Much of this tutorial is based on

- ► Koller & Friedman, Probabilistic Graphical Models. MIT 2009
- Wainwright & Jordan, Graphical Models, Exponential Families, and Variational Inference. FTML 2008
- Bishop, Pattern Recognition and Machine Learning. Springer 2006.

・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・

Outline

Life without Graphical Models

Representation

Directed Graphical Models (Bayesian Networks) Undirected Graphical Models (Markov Random Fields)

▲ロト ▲帰ト ▲ヨト ▲ヨト - ヨ - の々ぐ

Inference

Exact Inference Markov Chain Monte Carlo Variational Inference Loopy Belief Propagation Mean Field Algorithm Exponential Family Maximizing Problems

Parameter Learning

Structure Learning

- "Graphical model" is the study of probabilistic models
- Just because there is a graph with nodes and edges doesn't mean it's GM

These are not graphical models









neural network

decision tree

network flow

HMM template

・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・

Outline

Life without Graphical Models

Representation

Directed Graphical Models (Bayesian Networks)

Undirected Graphical Models (Markov Random Fields)

◆□▶ ◆□▶ ◆三▶ ◆三▶ 三三 のへで

Inference

Exact Inference Markov Chain Monte Carlo Variational Inference Loopy Belief Propagation Mean Field Algorithm Exponential Family Maximizing Problems

Parameter Learning

Structure Learning

Bayesian Network

- A directed graph has nodes X = (x₁,...,x_n), some of them connected by directed edges x_i → x_j
- A cycle is a directed path $x_1 \rightarrow \ldots \rightarrow x_k$ where $x_1 = x_k$
- A directed acyclic graph (DAG) contains no cycles
- A Bayesian network on the DAG is a family of distributions satisfying

$$\{p \mid p(X) = \prod_{i} p(x_i \mid Pa(x_i))\}$$

where $Pa(x_i)$ is the set of parents of x_i .

- ▶ p(x_i | Pa(x_i)) is the conditional probability distribution (CPD) at x_i
- By specifying the CPDs for all i, we specify a particular distribution p(X)

Binary variables



$$P(B, \sim E, A, J, \sim M)$$

$$= P(B)P(\sim E)P(A \mid B, \sim E)P(J \mid A)P(\sim M \mid A)$$

$$= 0.001 \times (1 - 0.002) \times 0.94 \times 0.9 \times (1 - 0.7)$$

$$\approx .000253$$

Example: Naive Bayes





- $p(y, x_1, \dots, x_d) = p(y) \prod_{i=1}^d p(x_i \mid y)$
- Used extensively in natural language processing
- Plate representation on the right



The two BNs are equivalent in all respects

- Bayesian networks imply no causality at all
- They only encode the joint probability distribution (hence correlation)
- However, people tend to design BNs based on causal relations



A generative model for $p(\phi, \theta, z, w \mid \alpha, \beta)$: For each topic t

```
\phi_t \sim \mathsf{Dirichlet}(\beta)
```

For each document d

 $\begin{array}{l} \theta \sim \mathsf{Dirichlet}(\alpha) \\ \mathsf{For \ each \ word \ position \ in \ } d \\ \mathsf{topic} \ z \sim \mathsf{Multinomial}(\theta) \\ \mathsf{word} \ w \sim \mathsf{Multinomial}(\phi_z) \\ \mathsf{Inference \ goals:} \ p(z \mid w, \alpha, \beta), \mathrm{argmax}_{\phi, \theta} \ p(\phi, \theta \mid w, \alpha, \beta) \end{array}$



A generative model for $p(\phi, \theta, z, w \mid \alpha, \beta)$: For each topic t

 $\phi_t \sim \mathsf{Dirichlet}(\beta)$

For each document d

$$\begin{split} \theta &\sim \mathsf{Dirichlet}(\alpha) \\ \mathsf{For each word position in } d \\ \mathsf{topic} \ z &\sim \mathsf{Multinomial}(\theta) \\ \mathsf{word} \ w &\sim \mathsf{Multinomial}(\phi_z) \\ \mathsf{Inference goals:} \ p(z \mid w, \alpha, \beta), \mathrm{argmax}_{\phi, \theta} \ p(\phi, \theta \mid w, \alpha, \beta) \end{split}$$

・ロト ・ 理 ト ・ ヨ ト ・ ヨ ト ・ ヨ



A generative model for $p(\phi, \theta, z, w \mid \alpha, \beta)$: For each topic t

```
\phi_t \sim \mathsf{Dirichlet}(\beta)
```

For each document d

$$\begin{split} \theta &\sim \mathsf{Dirichlet}(\alpha) \\ \mathsf{For each word position in } d \\ \mathsf{topic} \ z &\sim \mathsf{Multinomial}(\theta) \\ \mathsf{word} \ w &\sim \mathsf{Multinomial}(\phi_z) \\ \mathsf{Inference goals:} \ p(z \mid w, \alpha, \beta), \mathrm{argmax}_{\phi \mid \theta} \, p(\phi, \theta \mid w, \alpha, \beta) \end{split}$$

・ロト ・ 理 ト ・ ヨ ト ・ ヨ ト ・ ヨ



A generative model for $p(\phi, \theta, z, w \mid \alpha, \beta)$: For each topic t

```
\phi_t \sim \mathsf{Dirichlet}(\beta)
```

For each document d

 $\begin{array}{l} \theta \sim \mathsf{Dirichlet}(\alpha) \\ \mathsf{For each word position in } d \\ \mathsf{topic} \ z \sim \mathsf{Multinomial}(\theta) \\ \mathsf{word} \ w \sim \mathsf{Multinomial}(\phi_z) \\ \mathsf{Inference goals:} \ p(z \mid w, \alpha, \beta), \mathrm{argmax}_{\phi, \theta} \ p(\phi, \theta \mid w, \alpha, \beta) \end{array}$



A generative model for $p(\phi, \theta, z, w \mid \alpha, \beta)$: For each topic t

 $\phi_t \sim \mathsf{Dirichlet}(\beta)$

For each document d

 $\theta \sim \mathsf{Dirichlet}(\alpha)$

```
\begin{array}{l} \mbox{For each word position in } d \\ \mbox{topic } z \sim \mbox{Multinomial}(\theta) \\ \mbox{word } w \sim \mbox{Multinomial}(\phi_z) \\ \mbox{Inference goals: } p(z \mid w, \alpha, \beta), \mbox{argmax}_{\phi, \theta} \, p(\phi, \theta \mid w, \alpha, \beta) \end{array}
```



A generative model for $p(\phi, \theta, z, w \mid \alpha, \beta)$: For each topic t

```
\phi_t \sim \text{Dirichlet}(eta)
For each document d
	heta \sim \text{Dirichlet}(lpha)
```

```
For each word position in d
```

```
topic z \sim \text{Multinomial}(\theta)
word w \sim \text{Multinomial}(\phi_z)
Inference goals: p(z \mid w, \alpha, \beta), \operatorname{argmax}_{\phi, \theta} p(\phi, \theta \mid w, \alpha, \beta)
```

・ロト ・ 理 ト ・ ヨ ト ・ ヨ ト ・ ヨ



A generative model for $p(\phi, \theta, z, w \mid \alpha, \beta)$: For each topic t

```
\phi_t \sim \mathsf{Dirichlet}(\beta)
```

For each document d

 $\begin{array}{l} \theta \sim \mathsf{Dirichlet}(\alpha) \\ \mathsf{For each word position in } d \\ & \mathsf{topic} \ z \sim \mathsf{Multinomial}(\theta) \\ & \mathsf{word} \ w \sim \mathsf{Multinomial}(\phi_z) \\ \mathsf{Inference goals:} \ p(z \mid w, \alpha, \beta), \mathrm{argmax}_{\phi, \theta} \ p(\phi, \theta \mid w, \alpha, \beta) \end{array}$



A generative model for $p(\phi, \theta, z, w \mid \alpha, \beta)$: For each topic t

```
\phi_t \sim \mathsf{Dirichlet}(\beta)
```

For each document d

$$\begin{split} \theta &\sim \mathsf{Dirichlet}(\alpha) \\ \mathsf{For each word position in } d \\ \mathsf{topic } z &\sim \mathsf{Multinomial}(\theta) \\ \mathsf{word } w &\sim \mathsf{Multinomial}(\phi_z) \\ \mathsf{Inference goals: } p(z \mid w, \alpha, \beta), \mathrm{argmax}_{\phi, \theta} \, p(\phi, \theta \mid w, \alpha, \beta) \end{split}$$

Some Topics by LDA on the Wish Corpus

$p(\mathsf{word} \mid \mathsf{topic})$



"troops"



"election"



"love"

◆□▶ ◆□▶ ◆三▶ ◆三▶ 三三 のへぐ

- ► Two r.v.s A, B are independent if P(A, B) = P(A)P(B) or P(A|B) = P(A) (the two are equivalent)
- ► Two r.v.s A, B are conditionally independent given C if P(A, B | C) = P(A | C)P(B | C) or P(A | B, C) = P(A | C) (the two are equivalent)
- This extends to groups of r.v.s
- Conditional independence in a BN is precisely specified by d-separation ("directed separation")

(日) (同) (三) (三) (三) (○) (○)





- A, B in general dependent
- A, B conditionally independent given C
- C is a tail-to-tail node, blocks the undirected path A-B





- ► A, B in general dependent
- A, B conditionally independent given C
- C is a head-to-tail node, blocks the path A-B





- A, B in general independent
- A, B conditionally dependent given C, or any of C's descendants
- C is a head-to-head node, unblocks the path A-B

- Any groups of nodes A and B are conditionally independent given another group C, if all undirected paths from any node in A to any node in B are *blocked*
- \blacktriangleright A path is blocked if it includes a node x such that either
 - The path is head-to-tail or tail-to-tail at x and $x \in C$, or
 - The path is head-to-head at x, and neither x nor any of its descendants is in C.

・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・

- The path from A to B not blocked by either E or F
- ► A, B dependent given C



▲ロト ▲御 ト ▲ 臣 ト ▲ 臣 ト の Q @

- The path from A to B is blocked both at E and F
- ► A, B conditionally independent given F



▲ロト ▲御 ト ▲ 臣 ト ▲ 臣 ト の Q @

Outline

Life without Graphical Models

Representation

Directed Graphical Models (Bayesian Networks) Undirected Graphical Models (Markov Random Fields)

◆□▶ ◆□▶ ◆三▶ ◆三▶ 三三 のへで

Inference

Exact Inference Markov Chain Monte Carlo Variational Inference Loopy Belief Propagation Mean Field Algorithm Exponential Family Maximizing Problems

Parameter Learning

Structure Learning

- The efficiency of directed graphical model (acyclic graph, locally normalized CPDs) also makes it restrictive
- A clique C in an undirected graph is a fully connected set of nodes (note: full of loops!)
- Define a nonnegative potential function $\psi_C: X_C \mapsto \mathbb{R}_+$
- An undirected graphical model (aka Markov Random Field) on the graph is a family of distributions satisfying

$$\left\{ p \mid p(X) = \frac{1}{Z} \prod_{C} \psi_{C}(X_{C}) \right\}$$

• $Z = \int \prod_{C} \psi_{C}(X_{C}) dX$ is the partition function

Example: A Tiny Markov Random Field



- ▶ $x_1, x_2 \in \{-1, 1\}$
- A single clique $\psi_C(x_1, x_2) = e^{ax_1x_2}$

▶
$$p(x_1, x_2) = \frac{1}{Z} e^{ax_1 x_2}$$

$$\blacktriangleright \ Z = (e^a + e^{-a} + e^{-a} + e^a)$$

- $\blacktriangleright \ p(1,1) = p(-1,-1) = e^a/(2e^a + 2e^{-a})$
- $\blacktriangleright \ p(-1,1) = p(1,-1) = e^{-a}/(2e^a + 2e^{-a})$
- When the parameter a > 0, favor homogeneous chains
- When the parameter a < 0, favor inhomogeneous chains

- ▶ Real-valued feature functions $f_1(X), \ldots, f_k(X)$
- Real-valued weights w_1, \ldots, w_k

$$p(X) = \frac{1}{Z} \exp\left(-\sum_{i=1}^{k} w_i f_i(X)\right)$$

◆□▶ ◆□▶ ◆臣▶ ◆臣▶ 臣 のへぐ

Example: The Ising Model



This is an undirected model with $x \in \{0, 1\}$.

$$p_{\theta}(x) = \frac{1}{Z} \exp\left(\sum_{s \in V} \theta_s x_s + \sum_{(s,t) \in E} \theta_{st} x_s x_t\right)$$

・ロト ・西ト ・ヨト ・ヨー うらぐ

►
$$f_s(X) = x_s$$
, $f_{st}(X) = x_s x_t$
► $w_s = -\theta_s$, $w_{st} = -\theta_{st}$
Example: Image Denoising



[From Bishop PRML]

$$p(X) \sim N(\mu, \Sigma) = \frac{1}{(2\pi)^{n/2} |\Sigma|^{1/2}} \exp\left(-\frac{1}{2}(X-\mu)^{\top} \Sigma^{-1}(X-\mu)\right)$$

- Multivariate Gaussian
- The $n \times n$ covariance matrix Σ positive semi-definite
- Let $\Omega = \Sigma^{-1}$ be the precision matrix
- x_i, x_j are conditionally independent given all other variables, if and only if Ω_{ij} = 0

• When $\Omega_{ij} \neq 0$, there is an edge between x_i, x_j

Conditional Independence in Markov Random Fields

- Two group of variables A, B are conditionally independent given another group C, if
 - Remove C and all edges involving C
 - A, B beome disconnected



Factor Graph

- For both directed and undirected graphical models
- Bipartite: edges between a variable node and a factor node
- Factors represent computation



Outline

Life without Graphical Models

Representation

Directed Graphical Models (Bayesian Networks) Undirected Graphical Models (Markov Random Fields)

Inference

Exact Inference Markov Chain Monte Carlo Variational Inference Loopy Belief Propagation Mean Field Algorithm Exponential Family Maximizing Problems

Parameter Learning

Structure Learning

Outline

Life without Graphical Models

Representation

Directed Graphical Models (Bayesian Networks) Undirected Graphical Models (Markov Random Fields)

Inference

Exact Inference

Markov Chain Monte Carl Variational Inference Loopy Belief Propagation Mean Field Algorithm Exponential Family Maximizing Problems

Parameter Learning

Structure Learning

- ▶ Let X = (X_Q, X_E, X_O) for query, evidence, and other variables.
- Infer $P(X_Q \mid X_E)$
- By definition

$$P(X_Q \mid X_E) = \frac{P(X_Q, X_E)}{P(X_E)} = \frac{\sum_{X_O} P(X_Q, X_E, X_O)}{\sum_{X_Q, X_O} P(X_Q, X_E, X_O)}$$

Summing exponential number of terms: with k variables in X_O each taking r values, there are r^k terms

- There are a bunch of "other" variables x_1, \ldots, x_k
- We sum over r values each variable can take $\sum_{x_i=v_1}^{v_r}$

• This is exponential
$$(r^k)$$
: $\sum_{x_1...x_k}$

• We want
$$\sum_{x_1...x_k} p(X)$$

▶ For a graphical model, the joint probability factors $p(X) = \prod_{j=1}^{m} f_j(X_{(j)})$

・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・

• Each factor f_j operates on $X_{(j)} \subseteq X$

Eliminating a Variable

- ▶ Rearrange factors $\sum_{x_1...x_k} f_1^- \dots f_l^- f_{l+1}^+ \dots f_m^+$ by whether $x_1 \in X_{(j)}$
- Obviously equivalent: $\sum_{\boldsymbol{x}_2...\boldsymbol{x}_k} f_1^- \dots f_l^- \left(\sum_{\boldsymbol{x}_1} f_{l+1}^+ \dots f_m^+ \right)$
- Introduce a new factor $f_{m+1}^- = \left(\sum_{x_1} f_{l+1}^+ \dots f_m^+\right)$
- f_{m+1}^- contains the union of variables in $f_{l+1}^+ \dots f_m^+$ except x_1
- ▶ In fact, x_1 disappears altogether in $\sum_{x_2...x_k} f_1^- \dots f_l^- f_{m+1}^-$
- ▶ Dynamic programming: compute f_{m+1}^- once, use it thereafter

- Hope: f_{m+1}^- contains very few variables
- Recursively eliminate other variables in turn

Example: Chain Graph



- Binary variables
- ▶ Say we want $P(D) = \sum_{A,B,C} P(A)P(B|A)P(C|B)P(D|C)$

Let f₁(A) = P(A). Note f₁ is an array of size two:
 P(A = 0)
 P(A = 1)

$$P(B = 1|A = 0)$$

 $P(B = 1|A = 1))$

 $\sum_{A,B,C} f_1(A) f_2(A,B) f_3(B,C) f_4(C,D) = \\ \sum_{B,C} f_3(B,C) f_4(C,D) (\sum_A f_1(A) f_2(A,B))$

Example: Chain Graph



- ▶ $f_1(A)f_2(A, B)$ an array of size four: match A values P(A = 0)P(B = 0|A = 0) P(A = 1)P(B = 0|A = 1) P(A = 0)P(B = 1|A = 0)P(A = 1)P(B = 1|A = 1)
- ▶ $f_5(B) \equiv \sum_A f_1(A)f_2(A, B)$ an array of size two P(A = 0)P(B = 0|A = 0) + P(A = 1)P(B = 0|A = 1)P(A = 0)P(B = 1|A = 0) + P(A = 1)P(B = 1|A = 1)
- For this example, $f_5(B)$ happens to be P(B)
- ► $\sum_{B,C} f_3(B,C) f_4(C,D) f_5(B) =$ $\sum_C f_4(C,D) (\sum_B f_3(B,C) f_5(B))$, and so on
- ▶ In the end, $f_7(D) = (P(D = 0), P(D = 1))$



- Computation for P(D): 12 ×, 6 +
- Enumeration: 48 \times , 14 +
- Saving depends on elimination order. Finding optimal order NP-hard; there are heuristic methods.
- Saving depends more critically on the graph structure (tree width), can be intractable

- ▶ For evidence variables X_E , simply plug in their value e
- Eliminate variables $X_O = X X_E X_Q$
- The final factor will be the *joint* $f(X_Q) = P(X_Q, X_E = e)$
- Normalize to answer query:

$$P(X_Q \mid X_E = e) = \frac{f(X_Q)}{\sum_{X_Q} f(X_Q)}$$

(日)、(型)、(E)、(E)、(E)、(Q)

- Enumeration
- Variable elimination
- Not covered: junction tree (aka clique tree)

◆□▶ ◆□▶ ◆三▶ ◆三▶ 三三 のへぐ

Exact, but intractable for large graphs

Outline

Life without Graphical Models

Representation

Directed Graphical Models (Bayesian Networks) Undirected Graphical Models (Markov Random Fields)

Inference

Exact Inference

Markov Chain Monte Carlo

Variational Inference Loopy Belief Propagation Mean Field Algorithm Exponential Family Maximizing Problems

Parameter Learning

Structure Learning

Inference by Monte Carlo

▶ Consider the inference problem $p(X_Q = c_Q \mid X_E)$ where $X_Q \cup X_E \subseteq \{x_1 \dots x_n\}$

$$p(X_Q = c_Q \mid X_E) = \int \mathbb{1}_{(x_Q = c_Q)} p(x_Q \mid X_E) dx_Q$$

► If we can draw samples x⁽¹⁾_Q,...x^(m)_Q ~ p(x_Q | X_E), an unbiased estimator is

$$p(X_Q = c_Q \mid X_E) \approx \frac{1}{m} \sum_{i=1}^m \mathbf{1}_{(x_Q^{(i)} = c_Q)}$$

- \blacktriangleright The variance of the estimator decreases as \mathbb{V}/m
- ▶ Inference reduces to sampling from $p(x_Q \mid X_E)$

Forward Sampling Example



To generate a sample X = (B, E, A, J, M):

- 1. Sample $B \sim \text{Ber}(0.001)$: $r \sim U(0,1)$. If (r < 0.001) then B = 1 else B = 0
- 2. Sample $E \sim Ber(0.002)$
- 3. If B = 1 and E = 1, sample $A \sim Ber(0.95)$, and so on
- 4. If A = 1 sample $J \sim Ber(0.9)$ else $J \sim Ber(0.05)$
- 5. If A = 1 sample $M \sim Ber(0.7)$ else $M \sim Ber(0.01)$

Works for Bayesian networks.

- ▶ Say the inference task is P(B = 1 | E = 1, M = 1)
- Throw away all samples except those with (E = 1, M = 1)

$$p(B = 1 \mid E = 1, M = 1) \approx \frac{1}{m} \sum_{i=1}^{m} \mathbf{1}_{(B^{(i)} = 1)}$$

・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・

where m is the number of surviving samples

- Can be highly inefficient (note P(E = 1) tiny)
- Does not work for Markov Random Fields

Gibbs Sampler Example: P(B = 1 | E = 1, M = 1)

- Gibbs sampler is a Markov Chain Monte Carlo (MCMC) method.
- Directly sample from $p(x_Q \mid X_E)$
- Works for both graphical models
- Initialization:
 - Fix evidence; randomly set other variables
 - e.g. $X^{(0)} = (B = 0, E = 1, A = 0, J = 0, M = 1)$



Gibbs Update

- For each non-evidence variable x_i , fixing all other nodes X_{-i} , resample its value $x_i \sim P(x_i \mid X_{-i})$
- This is equivalent to $x_i \sim P(x_i \mid \mathsf{MarkovBlanket}(x_i))$
- ▶ For a Bayesian network MarkovBlanket(x_i) includes x_i's parents, spouses, and children

$$P(x_i \mid \mathsf{MarkovBlanket}(x_i)) \propto P(x_i \mid Pa(x_i)) \prod_{y \in C(x_i)} P(y \mid Pa(y))$$

where Pa(x) are the parents of x, and C(x) the children of x.

- For many graphical models the Markov Blanket is small.
- For example,

$$B \sim P(B \mid E=1, A=0) \propto P(B)P(A=0 \mid B, E=1)$$



Gibbs Update

- Say we sampled B = 1. Then $X^{(1)} = (B = 1, E = 1, A = 0, J = 0, M = 1)$
- Starting from $X^{(1)}$, sample $A \sim P(A \mid B = 1, E = 1, J = 0, M = 1)$ to get $X^{(2)}$
- Move on to J, then repeat $B, A, J, B, A, J \dots$
- ▶ Keep all *later* samples. P(B = 1 | E = 1, M = 1) is the fraction of samples with B = 1.



Gibbs Example 2: The Ising Model



This is an undirected model with $x \in \{0, 1\}$.

$$p_{\theta}(x) = \frac{1}{Z} \exp\left(\sum_{s \in V} \theta_s x_s + \sum_{(s,t) \in E} \theta_{st} x_s x_t\right)$$

・ロト ・西ト ・ヨト ・ヨー うらぐ

Gibbs Example 2: The Ising Model



- ▶ The Markov blanket of x_s is A, B, C, D
- In general for undirected graphical models

$$p(x_s \mid x_{-s}) = p(x_s \mid x_{N(s)})$$

N(s) is the neighbors of s.

The Gibbs update is

$$p(x_s = 1 \mid x_{N(s)}) = \frac{1}{\exp(-(\theta_s + \sum_{t \in N(s)} \theta_{st} x_t)) + 1}$$

- 日本 - 1 日本 - 1 日本 - 1 日本

- A Markov chain is defined by a transition matrix $T(X' \mid X)$
- \blacktriangleright Certain Markov chains have a stationary distribution π such that $\pi=T\pi$
- Gibbs sampler is such a Markov chain with $T_i((X_{-i}, x'_i) | (X_{-i}, x_i)) = p(x'_i | X_{-i})$, and stationary distribution $p(x_Q | X_E)$
- But it takes time for the chain to reach stationary distribution (mix)
 - Can be difficult to assert mixing
 - In practice "burn in": discard $X^{(0)}, \ldots, X^{(T)}$
 - Use all of $X^{(T+1)}, \ldots$ for inference (they are correlated)

Do not thin

Collapsed Gibbs Sampling

- ▶ In general, $\mathbb{E}_p[f(X)] \approx \frac{1}{m} \sum_{i=1}^m f(X^{(i)})$ if $X^{(i)} \sim p$
- Sometimes X = (Y, Z) where Z has closed-form operations

If so,

$$\mathbb{E}_p[f(X)] = \mathbb{E}_{p(Y)} \mathbb{E}_{p(Z|Y)}[f(Y,Z)]$$
$$\approx \frac{1}{m} \sum_{i=1}^m \mathbb{E}_{p(Z|Y^{(i)})}[f(Y^{(i)},Z)]$$

 $\text{ if }Y^{(i)}\sim p(Y)$

- No need to sample Z: it is collapsed
- ► Collapsed Gibbs sampler $T_i((Y_{-i}, y'_i) | (Y_{-i}, y_i)) = p(y'_i | Y_{-i})$
- Note $p(y'_i \mid Y_{-i}) = \int p(y'_i, Z \mid Y_{-i}) dZ$

Example: Collapsed Gibbs Sampling for LDA



Collapse θ, ϕ , Gibbs update:

$$P(z_{i} = j \mid \mathbf{z}_{-i}, \mathbf{w}) \propto \frac{n_{-i,j}^{(w_{i})} + \beta n_{-i,j}^{(d_{i})} + \alpha}{n_{-i,j}^{(\cdot)} + W\beta n_{-i,\cdot}^{(d_{i})} + T\alpha}$$

- ▶ n^(w_i)_{-i,j}: number of times word w_i has been assigned to topic j, excluding the current position
- n^(d_i)_{-i,j}: number of times a word from document d_i has been assigned to topic j, excluding the current position
- n^(·)_{-i,j}: number of times any word has been assigned to topic j, excluding the current position
- ▶ $n_{-i,:}^{(d_i)}$: length of document d_i , excluding the current position

- Gibbs sampling
- Not covered: block Gibbs, Metropolis-Hastings

Unbiased (after burn-in), but can have high variance

To learn more, come to Prof. Prasad Tetali's tutorial "Markov Chain Mixing with Applications" 2pm Monday.

Outline

Life without Graphical Models

Representation

Directed Graphical Models (Bayesian Networks) Undirected Graphical Models (Markov Random Fields)

Inference

Exact Inference Markov Chain Monte Carlo Variational Inference Loopy Belief Propagation Mean Field Algorithm Exponential Family Maximizing Problems

Parameter Learning

Structure Learning

Outline

Life without Graphical Models

Representation

Directed Graphical Models (Bayesian Networks) Undirected Graphical Models (Markov Random Fields)

Inference

Exact Inference Markov Chain Monte Carlo Variational Inference Loopy Belief Propagation Mean Field Algorithm Exponential Family Maximizing Problems

Parameter Learning

Structure Learning

- Also known as belief propagation (BP)
- Exact if the graph is a tree; otherwise known as "loopy BP", approximate
- ▶ The algorithm involves passing *messages* on the factor graph

(日)、(型)、(E)、(E)、(E)、(Q)

Alternative view: variational approximation (more later)

Example: A Simple HMM

The Hidden Markov Model template (not a graphical model)



 $\pi_1 = \pi_2 = 1/2$

• Observing $x_1 = R, x_2 = G$, the directed graphical model



 $P(z_l)P(x_l/z_l)$

 $P(z_2|z_1)P(x_2|z_2)$

A message is a vector of length K, where K is the number of values x takes.

There are two types of messages:

- 1. $\mu_{f \to x}$: message from a factor node f to a variable node x $\mu_{f \to x}(i)$ is the ith element, $i = 1 \dots K$.
- 2. $\mu_{x \to f}$: message from a variable node x to a factor node f

(日)、(型)、(E)、(E)、(E)、(Q)

Leaf Messages

- Assume tree factor graph. Pick an arbitrary root, say z₂
- Start messages at leaves.
- ▶ If a leaf is a factor node f, $\mu_{f \to x}(x) = f(x)$
- ▶ If a leaf is a variable node x, $\mu_{x \to f}(x) = 1$



 $\mu_{f_1 \to z_1}(z_1 = 1) = P(z_1 = 1)P(R|z_1 = 1) = 1/2 \cdot 1/2 = 1/4$ $\mu_{f_1 \to z_1}(z_1 = 2) = P(z_1 = 2)P(R|z_1 = 2) = 1/2 \cdot 1/4 = 1/8$



 $\pi_1 = \pi_2 = 1/2$

▲□▶ ▲圖▶ ▲臣▶ ▲臣▶ ―臣 … のへで

Message from Variable to Factor

- A node (factor or variable) can send out a message if all other incoming messages have arrived
- Let x be in factor f_s .

$$\mu_{x \to f_s}(x) = \prod_{f \in ne(x) \setminus f_s} \mu_{f \to x}(x)$$

• $ne(x) \setminus f_s$ are factors connected to x excluding f_s .



Message from Factor to Variable

• Let x be in factor f_s . Let the other variables in f_s be $x_{1:M}$.

$$\mu_{f_s \to x}(x) = \sum_{x_1} \dots \sum_{x_M} f_s(x, x_1, \dots, x_M) \prod_{m=1}^M \mu_{x_m \to f_s}(x_m)$$

$$\boxed{f_1 \qquad f_2 \qquad f_2} \qquad f_2$$

$$P(z_l)P(x_l/z_l) \qquad P(z_2/z_l)P(x_2/z_2)$$

$$\mu_{f_2 \to z_2}(s) = \sum_{s'=1}^2 \mu_{z_1 \to f_2}(s')f_2(z_1 = s', z_2 = s)$$

$$= 1/4P(z_2 = s|z_1 = 1)P(x_2 = G|z_2 = s)$$

$$+1/8P(z_2 = s|z_1 = 2)P(x_2 = G|z_2 = s)$$

▲□▶ ▲□▶ ▲□▶ ▲□▶ ▲□▶ ▲□▶ ▲□▶

We get $\begin{aligned} &\mu_{f_2 \to z_2}(z_2 = 1) = 1/32 \\ &\mu_{f_2 \to z_2}(z_2 = 2) = 1/8 \end{aligned}$

The message has reached the root, pass it back down



 $\pi_1 = \pi_2 = 1/2$

▲ロト ▲帰ト ▲ヨト ▲ヨト - ヨ - の々ぐ
Keep Passing Down





 $\pi_1 = \pi_2 = 1/2$

Once a variable receives all incoming messages, we compute its marginal as

$$p(x) \propto \prod_{f \in ne(x)} \mu_{f \to x}(x)$$

In this example

$$\begin{split} P(z_1|x_1, x_2) &\propto \mu_{f_1 \to z_1} \cdot \mu_{f_2 \to z_1} = \binom{1/4}{1/8} \cdot \binom{7/16}{3/8} = \binom{7/64}{3/64} \Rightarrow \binom{0.7}{0.3} \\ P(z_2|x_1, x_2) &\propto \mu_{f_2 \to z_2} = \binom{1/32}{1/8} \Rightarrow \binom{0.2}{0.8} \\ \end{split}$$
 One can also compute the marginal of the set of variables x_s

involved in a factor f_s

$$p(x_s) \propto f_s(x_s) \prod_{x \in ne(f)} \mu_{x \to f}(x)$$

Handling Evidence

Observing x = v,

- we can absorb it in the factor (as we did); or
- set messages $\mu_{x \to f}(x) = 0$ for all $x \neq v$

Observing X_E ,

► multiplying the incoming messages to x ∉ X_E gives the joint (not p(x|X_E)):

$$p(x, X_E) \propto \prod_{f \in ne(x)} \mu_{f \to x}(x)$$

The conditional is easily obtained by normalization

$$p(x|X_E) = \frac{p(x, X_E)}{\sum_{x'} p(x', X_E)}$$

- So far, we assumed a tree graph
- When the factor graph contains loops, pass messages indefinitely until convergence
- But convergence may not happen
- But in many cases loopy BP still works well, empirically

Outline

Life without Graphical Models

Representation

Directed Graphical Models (Bayesian Networks) Undirected Graphical Models (Markov Random Fields)

Inference

Exact Inference Markov Chain Monte Carlo Variational Inference Loopy Belief Propagation Mean Field Algorithm Exponential Family Maximizing Problems

Parameter Learning

Structure Learning

Example: The Ising Model



The random variables x take values in $\{0, 1\}$.

$$p_{\theta}(x) = \frac{1}{Z} \exp\left(\sum_{s \in V} \theta_s x_s + \sum_{(s,t) \in E} \theta_{st} x_s x_t\right)$$

◆□▶ ◆□▶ ◆目▶ ◆目▶ ● ● ● ●

The Conditional



• Markovian: the conditional distribution for x_s is

$$p(x_s \mid x_{-s}) = p(x_s \mid x_{N(s)})$$

N(s) is the neighbors of s.

This reduces to

$$p(x_s = 1 \mid x_{N(s)}) = \frac{1}{\exp(-(\theta_s + \sum_{t \in N(s)} \theta_{st} x_t)) + 1}$$

• Gibbs sampling would draw x_s like this.

The Mean Field Algorithm for Ising Model

$$p(x_s = 1 \mid x_{N(s)}) = \frac{1}{\exp(-(\theta_s + \sum_{t \in N(s)} \theta_{st} x_t)) + 1}$$

 \blacktriangleright Instead of Gibbs sampling, let μ_s be the estimated marginal $p(x_s=1)$

$$\boldsymbol{\mu_s} \leftarrow \frac{1}{\exp(-(\theta_s + \sum_{t \in N(s)} \theta_{st} \boldsymbol{\mu_t})) + 1}$$

・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・

- The µ's are updated iteratively
- The Mean Field algorithm is coordinate ascent and guaranteed to converge to a local optimal (more later).

Outline

Life without Graphical Models

Representation

Directed Graphical Models (Bayesian Networks) Undirected Graphical Models (Markov Random Fields)

Inference

Exact Inference Markov Chain Monte Carlo Variational Inference Loopy Belief Propagation Mean Field Algorithm Exponential Family Maximizing Problems

Parameter Learning

Structure Learning

- ▶ Let $\phi(X) = (\phi_1(X), \dots, \phi_d(X))^\top$ be d sufficient statistics, where $\phi_i : \mathcal{X} \mapsto \mathbb{R}$
- ▶ Note X is all the nodes in a Graphical model
- $\phi_i(X)$ sometimes called a feature function
- Let $\theta = (\theta_1, \dots, \theta_d)^\top \in \mathbb{R}^d$ be canonical parameters.
- The exponential family is a family of probability densities:

$$p_{\theta}(\mathbf{x}) = \exp\left(\theta^{\top}\phi(\mathbf{x}) - A(\theta)\right)$$

(日) (同) (三) (三) (三) (○) (○)

$$p_{\theta}(\mathbf{x}) = \exp\left(\theta^{\top}\phi(\mathbf{x}) - A(\theta)\right)$$

- The key is the inner product between parameters θ and sufficient statistics φ.
- A is the log partition function,

$$A(\theta) = \log \int \exp\left(\theta^{\top} \phi(\mathbf{x})\right) \nu(d\mathbf{x})$$

◆□▶ ◆□▶ ◆臣▶ ◆臣▶ 臣 の�?

 $\blacktriangleright A = \log Z$

Parameters for which the density is normalizable:

$$\Omega = \{\theta \in \mathbb{R}^d \mid A(\theta) < \infty\}$$

- ► A minimal exponential family is where the φ's are linearly independent.
- An overcomplete exponential family is where the φ's are linearly dependent:

$$\exists \alpha \in \mathbb{R}^d, \ \alpha^{\top} \phi(\mathbf{x}) = \text{constant } \forall \mathbf{x}$$

Both minimal and overcomplete representations are useful.

Exponential Family Example 1: Bernoulli

$$p(x) = \beta^x (1 - \beta)^{1-x}$$
 for $x \in \{0, 1\}$ and $\beta \in (0, 1)$.

- Does not look like an exponential family!
- Can be rewritten as

$$p(x) = \exp\left(x\log\beta + (1-x)\log(1-\beta)\right)$$

・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・

- Now in exponential family form with $\phi_1(x) = x, \phi_2(x) = 1 x, \ \theta_1 = \log \beta, \theta_2 = \log(1 \beta)$, and $A(\theta) = 0$.
- ▶ Overcomplete: $\alpha_1 = \alpha_2 = 1$ makes $\alpha^{\top} \phi(x) = 1$ for all x

Exponential Family Example 1: Bernoulli

$$p(x) = \exp\left(x\log\beta + (1-x)\log(1-\beta)\right)$$

Can be further rewritten as

$$p(x) = \exp\left(x\theta - \log(1 + \exp(\theta))\right)$$

(日) (同) (三) (三) (三) (○) (○)

 Minimal exponential family with φ(x) = x, θ = log β/(1-β), A(θ) = log(1 + exp(θ)).
Many distributions (e.g., Gaussian, exponential, Poisson, Beta) are in the exponential family, but not all (e.g., the Laplace

distribution).

Exponential Family Example 2: Ising Model



• Binary random variable $x_s \in \{0, 1\}$

► d = |V| + |E| sufficient statistics: $\phi(\mathbf{x}) = (\dots x_s \dots x_{st} \dots)^\top$

• This is a regular $(\Omega = \mathbb{R}^d)$, minimal exponential family.

Exponential Family Example 3: Potts Model



Similar to Ising model but generalizing $x_s \in \{0, \ldots, r-1\}$.

▶ Indicator functions $f_{sj}(\mathbf{x}) = 1$ if $x_s = j$ and 0 otherwise, and $f_{stjk}(\mathbf{x}) = 1$ if $x_s = j \land x_t = k$, and 0 otherwise.

$$p_{\theta}(\mathbf{x}) = \exp\left(\sum_{sj} \theta_{sj} f_{sj}(\mathbf{x}) + \sum_{stjk} \theta_{stjk} f_{stjk}(\mathbf{x}) - A(\theta)\right)$$

- $\blacktriangleright \ d = r|V| + r^2|E|$
- ▶ Regular but overcomplete, because $\sum_{j=0}^{r-1} \theta_{sj}(\mathbf{x}) = 1$ for any $s \in V$ and all \mathbf{x} .
- ► The Potts model is a special case where the parameters are tied: $\theta_{stkk} = \alpha$, and $\theta_{stjk} = \beta$ for $j \neq k$.

For sufficient statistics defined by indicator functions

- ▶ e.g., $\phi_{sj}(\mathbf{x}) = f_{sj}(\mathbf{x}) = 1$ if $x_s = j$ and 0 otherwise
- The marginal can be obtained via the mean

$$\mathbb{E}_{\theta}[\phi_{sj}(\mathbf{x})] = P(x_s = j)$$

Since inference is about computing the marginal, in this case it is equivalent to computing the mean.

Mean Parameters

- ▶ Let *p* be any density (not necessarily in exponential family).
- Given sufficient statistics ϕ , the mean parameters $\mu = (\mu_1, \dots, \mu_d)^\top$ is

$$\mu_i = \mathbb{E}_p[\phi_i(\mathbf{x})] = \int \phi_i(\mathbf{x}) p(\mathbf{x}) d\mathbf{x}$$

The set of mean parameters

$$\mathcal{M} = \{\mu \in \mathbb{R}^d \mid \exists p \text{ s.t. } \mathbb{E}_p[\phi(\mathbf{x})] = \mu\}$$

- If $\mu^{(1)}, \mu^{(2)} \in \mathcal{M}$, there must exist $p^{(1)}, p^{(2)}$
- \blacktriangleright The convex combinations of $p^{(1)}, p^{(2)}$ leads to another mean parameter in $\mathcal M$
- Therefore \mathcal{M} is convex

• Let
$$\phi_1(x) = x, \phi_2(x) = x^2$$

- ► For any p (not necessarily Gaussian) on x, the mean parameters $\mu = (\mu_1, \mu_2) = (\mathbb{E}(x), \mathbb{E}(x^2))^{\top}$.
- ▶ Note $\mathbb{V}(x) = \mathbb{E}(x^2) \mathbb{E}^2(x) = \mu_2 \mu_1^2 \ge 0$ for any p

• \mathcal{M} is not \mathbb{R}^2 but rather the subset $\mu_1 \in \mathbb{R}, \mu_2 \ge \mu_1^2$.



- ▶ The marginal polytope is defined for discrete x_s
- Recall $\mathcal{M} = \{ \mu \in \mathbb{R}^d \mid \mu = \sum_{\mathbf{x}} \phi(\mathbf{x}) p(\mathbf{x}) \text{ for some } p \}$
- p can be a point mass function on a particular x.
- In fact any p is a convex combination of such point mass functions.
- M = conv{φ(x), ∀x} is a convex hull, called the marginal polytope.

Marginal Polytope Example

Tiny Ising model: two nodes $x_1, x_2 \in \{0, 1\}$ connected by an edge.

- minimal sufficient statistics $\phi(x_1, x_2) = (x_1, x_2, x_1 x_2)^{\top}$.
- only 4 different $\mathbf{x} = (x_1, x_2)$.
- the marginal polytope is $\mathcal{M} = conv\{(0,0,0), (0,1,0), (1,0,0), (1,1,1)\}$



- the convex hull is a polytope inside the unit cube.
- ► the three coordinates are node marginals µ₁ ≡ E_p[x₁ = 1], µ₂ ≡ E_p[x₂ = 1] and edge marginal µ₁₂ ≡ E_p[x₁ = x₂ = 1], hence the name.

For any regular exponential family, $A(\theta)$ is convex in θ .

- Strictly convex for minimal exponential family.
- Nice property: $\frac{\partial A(\theta)}{\partial \theta_i} = \mathbb{E}_{\theta}[\phi_i(\mathbf{x})]$
- Therefore, $\nabla A = \mu$, the mean parameters of p_{θ} .

The conjugate dual function A^* to A is defined as

$$A^*(\mu) = \sup_{\theta \in \Omega} \theta^\top \mu - A(\theta)$$

Such definition, where a quantity is expressed as the solution to an optimization problem, is called a *variational* definition.

- ► For any $\mu \in \mathcal{M}$'s interior, let $\theta(\mu)$ satisfy $\mathbb{E}_{\theta(\mu)}[\phi(\mathbf{x})] = \nabla A(\theta(\mu)) = \mu.$
- ► Then $A^*(\mu) = -H(p_{\theta(\mu)})$ the negative entropy.
- The dual of the dual gives back A: A(θ) = sup_{μ∈M} μ^Tθ − A^{*}(μ)
- ► For all $\theta \in \Omega$, the supremum is attained uniquely at the $\mu \in \mathcal{M}^0$ by the moment matching conditions $\mu = \mathbb{E}_{\theta}[\phi(\mathbf{x})]$.

Example: Conjugate Dual for Bernoulli

- ► Recall the minimal exponential family for Bernoulli with $\phi(x) = x, A(\theta) = \log(1 + \exp(\theta)), \Omega = \mathbb{R}.$
- By definition

$$A^*(\mu) = \sup_{\theta \in \mathbb{R}} \theta \mu - \log(1 + \exp(\theta))$$

Taking derivative and solve

$$A^*(\mu) = \mu \log \mu + (1 - \mu) \log(1 - \mu)$$

i.e., the negative entropy.

$$A(\theta) = \sup_{\mu \in \mathcal{M}} \mu^{\top} \theta - A^*(\mu) \text{ is attained by } \mu = \mathbb{E}_{\theta}[\phi(\mathbf{x})].$$

- Want to compute the marginals P(x_s = j)? They are the mean parameters μ_{sj} = E_θ[φ_{ij}(x)] under standard overcomplete representation.
- Want to compute the mean parameters μ_{sj}? They are the arg sup to the optimization problem above.
- This variational representation is exact, not approximate (will relax it next to derive loopy BP and mean field)

The Difficulties with Variational Representation

$$A(\theta) = \sup_{\mu \in \mathcal{M}} \mu^{\top} \theta - A^*(\mu)$$

- Difficult to solve even though it is a convex problem
- Two issues:
 - Although the marginal polytope *M* is convex, it can be quite complex (exponential number of vertices)
 - The dual function $A^*(\mu)$ usually does not admit an explicit form.
- Variational approximation modifies the optimization problem so that it is tractable, at the price of an approximate solution.
- Next, we cast mean field and sum-product algorithms as variational approximations.

- ► The mean field method replaces *M* with a simpler subset *M*(*F*) on which *A*^{*}(µ) has a closed form.
- \blacktriangleright Consider the fully disconnected subgraph $F=(V,\emptyset)$ of the original graph G=(V,E)
- Set all $\theta_i = 0$ if ϕ_i involves edges
- ► The densities in this sub-family are all fully factorized:

$$p_{\theta}(\mathbf{x}) = \prod_{s \in V} p(x_s; \theta_s)$$

・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・

The Geometry of $\mathcal{M}(F)$

- Let *M*(*F*) be the mean parameters of the fully factorized sub-family. In general, *M*(*F*) ⊂ *M*
- Recall \mathcal{M} is the convex hull of extreme points $\{\phi(\mathbf{x})\}$.
- It turns out the extreme points $\{\phi(\mathbf{x})\} \in \mathcal{M}(F)$.
- Example:
 - The tiny Ising model $x_1, x_2 \in \{0, 1\}$ with $\phi = (x_1, x_2, x_1 x_2)^\top$
 - ► The point mass distribution $p(\mathbf{x} = (0, 1)^{\top}) = 1$ is realized as a limit to the series $p(\mathbf{x}) = \exp(\theta_1 x_1 + \theta_2 x_2 A(\theta))$ where $\theta_1 \to -\infty$ and $\theta_2 \to \infty$.

- This series is in F because $\theta_{12} = 0$.
- Hence the extreme point $\phi(\mathbf{x}) = (0, 1, 0)$ is in $\mathcal{M}(F)$.



- ▶ Because the extreme points of *M* are in *M*(*F*), if *M*(*F*) were convex, we would have *M* = *M*(*F*).
- But in general $\mathcal{M}(F)$ is a true subset of \mathcal{M}
- Therefore, $\mathcal{M}(F)$ is a nonconvex inner set of \mathcal{M}



The Mean Field Method as Variational Approximation

Recall the exact variational problem

$$A(\theta) = \sup_{\mu \in \mathcal{M}} \mu^{\top} \theta - A^*(\mu)$$

attained by solution to inference problem $\mu = \mathbb{E}_{\theta}[\phi(\mathbf{x})].$

• The mean field method simply replaces \mathcal{M} with $\mathcal{M}(F)$

$$\mathcal{L}(\theta) = \sup_{\mu \in \mathcal{M}(F)} \mu^{\top} \theta - A^{*}(\mu)$$

- Obvious $\mathcal{L}(\theta) \leq A(\theta)$.
- The original solution μ^* may not be in $\mathcal{M}(F)$
- Even if $\mu^* \in \mathcal{M}(F)$, may hit local maximum and not find it
- Why both? Because A^{*}(µ) = −H(p_θ(µ)) has a very simple form for M(F)

Example: Mean Field for Ising Model

- ► The mean parameters for the Ising model are the node and edge marginals: µ_s = p(x_x = 1), µ_{st} = p(x_s = 1, x_t = 1)
- ▶ Fully factorized $\mathcal{M}(F)$ means no edge. $\mu_{st} = \mu_s \mu_t$
- ▶ For $\mathcal{M}(F)$, the dual function $A^*(\mu)$ has the simple form

$$A^*(\mu) = \sum_{s \in V} -H(\mu_s) = \sum_{s \in V} \mu_s \log \mu_s + (1 - \mu_s) \log(1 - \mu_s)$$

Thus the mean field problem is

$$\mathcal{L}(\theta) = \sup_{\mu \in \mathcal{M}(F)} \mu^{\top} \theta - \sum_{s \in V} (\mu_s \log \mu_s + (1 - \mu_s) \log(1 - \mu_s))$$
$$= \max_{(\mu_1 \dots \mu_m) \in [0,1]^m} \left(\sum_{s \in V} \theta_s \mu_s + \sum_{(s,t) \in E} \theta_{st} \mu_s \mu_t + \sum_{s \in V} H(\mu_s) \right)$$

◆□▶ ◆□▶ ◆三▶ ◆三▶ 三三 のへぐ

Example: Mean Field for Ising Model

$$\mathcal{L}(\theta) = \max_{(\mu_1 \dots \mu_m) \in [0,1]^m} \left(\sum_{s \in V} \theta_s \mu_s + \sum_{(s,t) \in E} \theta_{st} \mu_s \mu_t + \sum_{s \in V} H(\mu_s) \right)$$

- Bilinear in µ, not jointly concave
- But concave in a single dimension μ_s , fixing others.
- ▶ Iterative coordinate-wise maximization: fixing μ_t for $t \neq s$ and optimizing μ_s .
- Setting the partial derivative w.r.t. μ_s to 0 yields:

$$\mu_s = \frac{1}{1 + \exp\left(-(\theta_s + \sum_{(s,t) \in E} \theta_{st} \mu_t)\right)}$$

as we've seen before.

► Caution: mean field converges to a local maximum depending on the initialization of µ₁...µ_m.

$$A(\theta) = \sup_{\mu \in \mathcal{M}} \mu^{\top} \theta - A^*(\mu)$$

The sum-product algorithm makes two approximations:

- ▶ it relaxes *M* to an *outer* set *L*
- it replaces the dual A^* with an approximation.

$$A(\theta) = \sup_{\mu \in L} \mu^{\top} \theta - \tilde{A}^*(\mu)$$

The Outer Relaxation

- ► For overcomplete exponential families on discrete nodes, the mean parameters are node and edge marginals µ_{sj} = p(x_s = j), µ_{stjk} = p(x_s = j, x_t = k).
- The marginal polytope is $\mathcal{M} = \{\mu \mid \exists p \text{ with marginals } \mu\}.$
- Now consider τ ∈ ℝ^d₊ satisfying "node normalization" and "edge-node marginal consistency" conditions:

$$\sum_{j=0}^{r-1} \tau_{sj} = 1 \qquad \forall s \in V$$
$$\sum_{k=0}^{r-1} \tau_{stjk} = \tau_{sj} \qquad \forall s, t \in V, j = 0 \dots r - 1$$
$$\sum_{j=0}^{r-1} \tau_{stjk} = \tau_{tk} \qquad \forall s, t \in V, k = 0 \dots r - 1$$

• Define $L = \{\tau \text{ satisfying the above conditions}\}$.

The Outer Relaxation

- If the graph is a tree then $\mathcal{M} = L$
- If the graph has cycles then $\mathcal{M} \subset L$
 - L is too lax to satisfy other constraints that true marginals need to satisfy
- Nice property: L is still a polytope, but much simpler than \mathcal{M} .



- Recall μ are node and edge marginals
- If the graph is a tree, one can exactly reconstruct the joint probability

$$p_{\mu}(\mathbf{x}) = \prod_{s \in V} \mu_{sx_s} \prod_{(s,t) \in E} \frac{\mu_{stx_sx_t}}{\mu_{sx_s}\mu_{tx_t}}$$

If the graph is a tree, the entropy of the joint distribution is

$$H(p_{\mu}) = -\sum_{s \in V} \sum_{j=0}^{r-1} \mu_{sj} \log \mu_{sj} - \sum_{(s,t) \in E} \sum_{j,k} \mu_{stjk} \log \frac{\mu_{stjk}}{\mu_{sj}\mu_{tk}}$$

Neither holds for graph with cycles.
Define the *Bethe entropy* for $\tau \in L$ on loopy graphs in the same way:

$$H_{Bethe}(p_{\tau}) = -\sum_{s \in V} \sum_{j=0}^{r-1} \tau_{sj} \log \tau_{sj} - \sum_{(s,t) \in E} \sum_{j,k} \tau_{stjk} \log \frac{\tau_{stjk}}{\tau_{sj}\tau_{tk}}$$

Note H_{Bethe} is not a true entropy. The second approximation in sum-product is to replace $A^*(\tau)$ with $-H_{Bethe}(p_{\tau})$.

With these two approximations, we arrive at the variational problem

$$\tilde{A}(\theta) = \sup_{\tau \in L} \tau^{\top} \theta + H_{Bethe}(p_{\tau})$$

- Optimality conditions require the gradients vanish w.r.t. both τ and the Lagrangian multipliers on constraints $\tau \in L$.
- The sum-product algorithm can be derived as an iterative fixed point procedure to satisfy optimality conditions.
- At the solution, A
 (θ) is not guaranteed to be either an upper or a lower bound of A(θ)

 $\blacktriangleright \tau$ may not correspond to a true marginal distribution

- The sum-product algorithm (loopy belief propagation)
- The mean field method
- Not covered: Expectation Propagation

Efficient computation. But often *unknown* bias in solution.

・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・

Outline

Life without Graphical Models

Representation

Directed Graphical Models (Bayesian Networks) Undirected Graphical Models (Markov Random Fields)

Inference

Exact Inference Markov Chain Monte Carlo Variational Inference Loopy Belief Propagation Mean Field Algorithm Exponential Family

Maximizing Problems

Parameter Learning

Structure Learning

Maximizing Problems

Recall the HMM example



$$\pi_1 = \pi_2 = 1/2$$

There are two senses of "best states" $z_{1:N}$ given $x_{1:N}$:

- 1. So far we computed the marginal $p(z_n|x_{1:N})$
 - We can define "best" as $z_n^* = \arg \max_k p(z_n = k | x_{1:N})$
 - However $z_{1:N}^*$ as a whole may not be the best
 - In fact $z_{1:N}^*$ can even have zero probability!
- 2. An alternative is to find

$$z_{1:N}^* = \arg\max_{z_{1:N}} p(z_{1:N}|x_{1:N})$$

- finds the most likely state configuration as a whole
- The max-sum algorithm solves this
- Generalizes the Viterbi algorithm for HMMs

Simple modification to the sum-product algorithm: replace \sum with \max in the factor-to-variable messages.

$$\begin{split} \mu_{f_s \to x}(x) &= \max_{x_1} \dots \max_{x_M} f_s(x, x_1, \dots, x_M) \prod_{m=1}^M \mu_{x_m \to f_s}(x_m) \\ \mu_{x_m \to f_s}(x_m) &= \prod_{f \in ne(x_m) \setminus f_s} \mu_{f \to x_m}(x_m) \\ \mu_{x_{\mathsf{leaf}} \to f}(x) &= 1 \\ \mu_{f_{\mathsf{leaf}} \to x}(x) &= f(x) \end{split}$$

◆□▶ ◆□▶ ◆三▶ ◆三▶ 三三 のへぐ

- As in sum-product, pick an arbitrary variable node x as the root
- Pass messages up from leaves until they reach the root
- Unlike sum-product, do not pass messages back from root to leaves
- At the root, multiply incoming messages

$$p^{\max} = \max_{x} \left(\prod_{f \in ne(x)} \mu_{f \to x}(x) \right)$$

This is the probability of the most likely state configuration

- ► To identify the configuration itself, keep *back pointers*:
- When creating the message

$$\mu_{f_s \to x}(x) = \max_{x_1} \dots \max_{x_M} f_s(x, x_1, \dots, x_M) \prod_{m=1}^M \mu_{x_m \to f_s}(x_m)$$

for each x value, we separately create M pointers back to the values of x_1, \ldots, x_M that achieve the maximum.

◆□▶ ◆□▶ ◆三▶ ◆三▶ 三三 のへぐ

At the root, backtrack the pointers.



• Message from leaf f_1 $\mu_{f_1 \to z_1}(z_1 = 1) = P(z_1 = 1)P(R|z_1 = 1) = 1/2 \cdot 1/2 = 1/4$ $\mu_{f_1 \to z_1}(z_1 = 2) = P(z_1 = 2)P(R|z_1 = 2) = 1/2 \cdot 1/4 = 1/8$

The second message

$$\mu_{z_1 \to f_2}(z_1 = 1) = 1/4$$

$$\mu_{z_1 \to f_2}(z_1 = 2) = 1/8$$



 $\pi_1 = \pi_2 = 1/2$

$$\mu_{f_2 \to z_2}(z_2 = 1)$$

$$= \max_{z_1} f_2(z_1, z_2) \mu_{z_1 \to f_2}(z_1)$$

$$= \max_{z_1} P(z_2 = 1 \mid z_1) P(x_2 = G \mid z_2 = 1) \mu_{z_1 \to f_2}(z_1)$$

$$= \max(1/4 \cdot 1/4 \cdot 1/4, 1/2 \cdot 1/4 \cdot 1/8) = 1/64$$

Back pointer for $z_2 = 1$: either $z_1 = 1$ or $z_1 = 2$

◆□▶ ◆□▶ ◆臣▶ ◆臣▶ ○臣 - の々ぐ



 $\pi_1 = \pi_2 = 1/2$

The other element of the same message:

$$\mu_{f_2 \to z_2}(z_2 = 2)$$

$$= \max_{z_1} f_2(z_1, z_2) \mu_{z_1 \to f_2}(z_1)$$

$$= \max_{z_1} P(z_2 = 2 \mid z_1) P(x_2 = G \mid z_2 = 2) \mu_{z_1 \to f_2}(z_1)$$

$$= \max(3/4 \cdot 1/2 \cdot 1/4, 1/2 \cdot 1/2 \cdot 1/8) = 3/32$$

Back pointer for $z_2 = 2$: $z_1 = 1$

◆□▶ ◆□▶ ◆臣▶ ◆臣▶ 三臣 - のへで

At root z_2 ,



$$\max_{s=1,2} \mu_{f_2 \to z_2}(s) = 3/32$$
$$z_2 = 2 \to z_1 = 1$$

$$z_{1:2}^* = \arg\max_{z_{1:2}} p(z_{1:2}|x_{1:2}) = (1,2)$$

In this example, sum-product and max-product produce the same best sequence; In general they differ.

From Max-Product to Max-Sum

The *max-sum algorithm* is equivalent to the max-product algorithm, but work in log space to avoid underflow.

$$\mu_{f_s \to x}(x) = \max_{x_1 \dots x_M} \log f_s(x, x_1, \dots, x_M) + \sum_{m=1}^M \mu_{x_m \to f_s}(x_m)$$
$$\mu_{x_m \to f_s}(x_m) = \sum_{f \in ne(x_m) \setminus f_s} \mu_{f \to x_m}(x_m)$$
$$\mu_{x_{\mathsf{leaf}} \to f}(x) = 0$$
$$\mu_{f_{\mathsf{leaf}} \to x}(x) = \log f(x)$$

When at the root,

$$\log p^{\max} = \max_{x} \left(\sum_{f \in ne(x)} \mu_{f \to x}(x) \right)$$

The back pointers are the same.

Outline

Life without Graphical Models

Representation

Directed Graphical Models (Bayesian Networks) Undirected Graphical Models (Markov Random Fields)

▲ロト ▲帰ト ▲ヨト ▲ヨト - ヨ - の々ぐ

Inference

Exact Inference Markov Chain Monte Carlo Variational Inference Loopy Belief Propagation Mean Field Algorithm Exponential Family Maximizing Problems

Parameter Learning

Structure Learning

- Assume the graph structure is given
- Learning in exponential family: estimate θ from *iid* data x₁...x_n.
- Principle: maximum likelihood
- Distinguish two cases:
 - fully observed data: all dimensions of x are observed
 - ▶ *partially observed data*: some dimensions of x are unobserved.

(日)、(型)、(E)、(E)、(E)、(Q)

Fully Observed Data

$$p_{\theta}(\mathbf{x}) = \exp\left(\theta^{\top}\phi(\mathbf{x}) - A(\theta)\right)$$

• Given iid data $\mathbf{x}_1 \dots \mathbf{x}_n$, the log likelihood is

$$\ell(\theta) = \frac{1}{n} \sum_{i=1}^{n} \log p_{\theta}(\mathbf{x}_i) = \theta^{\top} \left(\frac{1}{n} \sum_{i=1}^{n} \phi(\mathbf{x}_i) \right) - A(\theta) = \theta^{\top} \hat{\mu} - A(\theta)$$

- $\hat{\mu} \equiv \frac{1}{n} \sum_{i=1}^{n} \phi(\mathbf{x}_i)$ is the mean parameter of the empirical distribution on $\mathbf{x}_1 \dots \mathbf{x}_n$. Clearly $\hat{\mu} \in \mathcal{M}$.
- ► Maximum likelihood: $\theta^{ML} = \arg \sup_{\theta \in \Omega} \theta^{\top} \hat{\mu} A(\theta)$
- ► The solution is θ^{ML} = θ(µ̂), the exponential family density whose mean parameter matches µ̂.
- ▶ When $\hat{\mu} \in \mathcal{M}^0$ and ϕ minimal, there is a unique maximum likelihood solution θ^{ML} .

- \blacktriangleright Each item (\mathbf{x},\mathbf{z}) where \mathbf{x} observed, \mathbf{z} unobserved
- Full data $(\mathbf{x}_1, \mathbf{z}_1) \dots (\mathbf{x}_n, \mathbf{z}_n)$, but we only observe $\mathbf{x}_1 \dots \mathbf{x}_n$
- ► The incomplete likelihood $\ell(\theta) = \frac{1}{n} \sum_{i=1}^{n} \log p_{\theta}(\mathbf{x}_i)$ where $p_{\theta}(\mathbf{x}_i) = \int p_{\theta}(\mathbf{x}_i, \mathbf{z}) d\mathbf{z}$
- Can be written as $\ell(\theta) = \frac{1}{n} \sum_{i=1}^{n} A_{\mathbf{x}_i}(\theta) A(\theta)$
- ▶ New log partition function of $p_{\theta}(\mathbf{z} \mid \mathbf{x}_i)$, one per item:

$$A_{\mathbf{x}_i}(\theta) = \log \int \exp(\theta^\top \phi(\mathbf{x}_i, \mathbf{z}')) d\mathbf{z}'$$

Expectation-Maximization (EM) algorithm: lower bound A_{xi}

EM as Variational Lower Bound

Mean parameter realizable by any distribution on z while holding x_i fixed:
Mathematical for some m

 $\mathcal{M}_{\mathbf{x}_i} = \{ \mu \in \mathbb{R}^d \mid \mu = \mathbb{E}_p[\phi(\mathbf{x}_i, \mathbf{z})] \text{ for some } p \}$

- ► The variational definition $A_{\mathbf{x}_i}(\theta) = \sup_{\mu \in \mathcal{M}_{\mathbf{x}_i}} \theta^\top \mu A^*_{\mathbf{x}_i}(\mu)$
- ► Trivial variational lower bound: $A_{\mathbf{x}_i}(\theta) \ge \theta^\top \mu^i - A^*_{\mathbf{x}_i}(\mu^i), \forall \mu^i \in \mathcal{M}_{\mathbf{x}_i}$

Lower bound L on the incomplete log likelihood:

$$\ell(\theta) = \frac{1}{n} \sum_{i=1}^{n} A_{\mathbf{x}_{i}}(\theta) - A(\theta)$$

$$\geq \frac{1}{n} \sum_{i=1}^{n} \left(\theta^{\top} \mu^{i} - A_{\mathbf{x}_{i}}^{*}(\mu^{i}) \right) - A(\theta)$$

$$\equiv \mathcal{L}(\mu^{1}, \dots, \mu^{n}, \theta)$$

・ロト・日本・モート モー うへぐ

The EM algorithm is coordinate ascent on $\mathcal{L}(\mu^1, \ldots, \mu^n, \theta)$.

 \blacktriangleright In the E-step, maximizes each μ^i

$$\mu^i \leftarrow \arg \max_{\mu^i \in \mathcal{M}_{\mathbf{x}_i}} \mathcal{L}(\mu^1, \dots, \mu^n, \theta)$$

- Equivalently, $\operatorname{argmax}_{\mu^i \in \mathcal{M}_{\mathbf{x}_i}} \theta^\top \mu^i A^*_{\mathbf{x}_i}(\mu^i)$
- This is the variational representation of the mean parameter $\mu^i(\theta) = \mathbb{E}_{\theta}[\phi(\mathbf{x}_i, \mathbf{z})]$

The E-step is named after this E_θ[] under the current parameters θ

• In the M-step, maximize θ holding the μ 's fixed:

$$\theta \leftarrow \arg \max_{\theta \in \Omega} \mathcal{L}(\mu^1, \dots, \mu^n, \theta) = \arg \max_{\theta \in \Omega} \theta^\top \hat{\mu} - A(\theta)$$

$$\blacktriangleright \hat{\mu} = \frac{1}{n} \sum_{i=1}^{n} \mu^i$$

- The solution $\theta(\hat{\mu})$ satisfies $\mathbb{E}_{\theta(\hat{\mu})}[\phi(\mathbf{x})] = \hat{\mu}$
- Standard fully observed maximum likelihood problem, hence the name M-step

(日)、(型)、(E)、(E)、(E)、(Q)

For loopy graphs E-step often intractable.

► Can't maximize

$$\max_{\mu^i \in \mathcal{M}_{\mathbf{x}_i}} \theta^\top \mu^i - A^*_{\mathbf{x}_i}(\mu^i)$$

- Improve but not necessarily maximize: "generalized EM"
- The mean field method maximizes

$$\max_{\mu^i \in \mathcal{M}_{\mathbf{x}_i}(F)} \theta^\top \mu^i - A^*_{\mathbf{x}_i}(\mu^i)$$

- up to local maximum
- recall $\mathcal{M}_{\mathbf{x}_i}(F)$ is an inner approximation to $\mathcal{M}_{\mathbf{x}_i}$
- Mean field E-step leads to generalized EM
- The sum-product algorithm does not lead to generalized EM

Outline

Life without Graphical Models

Representation

Directed Graphical Models (Bayesian Networks) Undirected Graphical Models (Markov Random Fields)

▲ロト ▲帰ト ▲ヨト ▲ヨト - ヨ - の々ぐ

Inference

Exact Inference Markov Chain Monte Carlo Variational Inference Loopy Belief Propagation Mean Field Algorithm Exponential Family Maximizing Problems

Parameter Learning

Structure Learning

- Let *M* be all allowed candidate features
- Let $M \subseteq \mathcal{M}$ be a log-linear model structure

$$P(X \mid M, \theta) = \frac{1}{Z} \exp\left(\sum_{i \in M} \theta_i f_i(X)\right)$$

- ► A score for the model M can be $\max_{\theta} \ln P(\mathsf{Data} \mid M, \theta)$
- The score is always better for larger M needs regularization

(日)、(型)、(E)、(E)、(E)、(Q)

• M and θ treated separately

Structure Learning for Gaussian Random Fields

- \blacktriangleright Consider a p-dimensional multivariate Gaussian $N(\mu,\Sigma)$
- The graphical model has p nodes x_1, \ldots, x_p
- ▶ The edge between x_i, x_j is absent if and only if $\Omega_{ij} = 0$, where $\Omega = \Sigma^{-1}$
- Equivalently, x_i, x_j are conditionally independent given other variables





・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・

Structure Learning for Gaussian Random Fields

- Let data be $X^{(1)}, \ldots, X^{(n)} \sim N(\mu, \Sigma)$
- ► The log likelihood is $\frac{n}{2} \log |\Omega| \frac{1}{2} \sum_{i=1}^{n} (X^{(i)} \mu)^{\top} \Omega(X^{(i)} \mu)$
- ► The maximum likelihood estimate of ∑ is the sample covariance

$$S = \frac{1}{n} \sum_{i} (X^{(i)} - \bar{X})^{\top} (X^{(i)} - \bar{X})$$

・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・

where \bar{X} is the sample mean

• S^{-1} is not a good estimate of Ω when n is small

▶ For centered data, minimize a regularized problem instead:

$$-\log |\Omega| + \frac{1}{n} \sum_{i=1}^{n} X^{(i)^{\top}} \Omega X^{(i)} + \lambda \sum_{i \neq j} |\Omega_{ij}|$$

◆□▶ ◆□▶ ◆三▶ ◆三▶ 三三 のへぐ

Known as glasso

• Given $GM = joint distribution p(x_1, \ldots, x_n)$

- BN or MRF
- conditional independence

▶ Do inference =
$$p(X_Q | X_E)$$
, in general $X_Q \cup X_E \subset \{x_1 \dots x_n\}$

- exact, MCMC, variational
- ▶ If $p(x_1, ..., x_n)$ not given, estimate it from data

(日)、(型)、(E)、(E)、(E)、(Q)

parameter and structure learning

Much on-going research!