

LECTURE 1
INTRODUCTION TO GM

Types of Graphical models

Directed models

Bayesian, markov models

Undirected models

Markov random fields, factor graphs

If graph is tree, 'message-passing' algos

Graphical models useful when dataset is limited & domain knowledge useful

For high dimensions \rightarrow too many parameters
Huge computational time

LECTURE 2:
CONDITIONAL INDEPENDENCE & DIRECTED MODELS

Conditional Independence

Definition:

$$X \perp Y | Z \Leftrightarrow P(X, Y | Z) = P(X | Z) \cdot P(Y | Z)$$

Directed acyclic graph

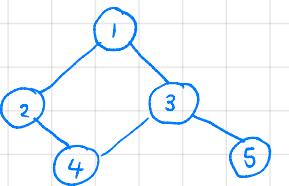
Set of nodes, no directed

Cycles

In a DAG, can ensure

parent nodes have lower values

(topological ordering)



Chain rule of probability

$$P(x_1, x_2, x_3, \dots) = P(x_1) \cdot P(x_2 | x_1) \cdot P(x_3 | x_1, x_2) \dots$$

$$= \prod_i P(x_i | X_{\text{Pa}(i)}) \quad \text{Pa} \rightarrow \text{Parent set}$$

Directed model

Assume we have DAG, parents have lower number

$$x_i \perp x_1 \dots x_{i-1} \mid x_{\text{Pa}(i)}$$

Proof:

$$\begin{aligned} P(x_1, \dots, x_n) &= \prod_{i=1}^n P(x_i | x_1, \dots, x_{i-1}) \\ &= \prod_{i=1}^n P(x_i | X_{\text{Pa}(i)}) \end{aligned}$$

Upshot:
Reduces no. of
Free parameters
Factorized
distribution of CI

A DAG encodes all assumptions

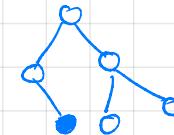
Example:

$$\begin{array}{c} 1 \quad 2 \quad 2 \\ P(x_1) \quad P(x_2 | x_1) \quad P(x_3 | x_1) \\ \text{Refer Fig (1)} \quad P(x_4 | x_2, x_3) \quad P(x_5 | x_3) \end{array} \quad \begin{array}{l} 11 \text{ free params} \\ (\text{original } 5-1) \end{array}$$

Topological ordering allows to drop CI.

Assume - A DAG, each var is CI of its non-descendants given parents

Starting CI's imply others (later)



Conditional Independence motivation

The world has true distribution.

& True factorized form maybe

Model may not correspond to True dist.

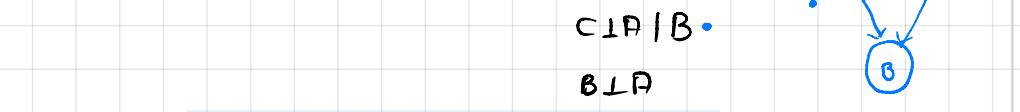
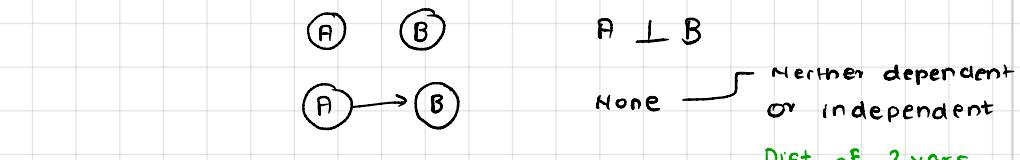
Why assumptions work?

- Domain-knowledge
- Convenience
- Intuition about causality
- Careful search / Structure learning

Causality

- Causal reasoning implies CI.
- However, CI does not imply causality

CI and Graphs



$A \perp C$ does not imply $A \perp C | B$

Example ① Boiler \rightarrow Alarm



Noise dependent on Alarm
i.e. hearing noise 'influences' estimate of Alarm

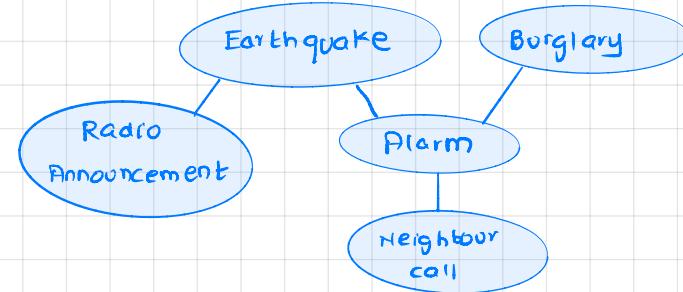
Example ② Boiler \rightarrow Alarm



Noise \perp Alarm

Since we know boiler is on, knowing that noise is present doesn't influence estimate

EG



$R \perp A E$	T
$E \perp B$	T
$H \perp E A$	T
$R \perp B A$	H
$E \perp B R, A$	H
* $E \perp B R, N$	N

LECTURE 3
BAYES BALL

D-separation & Bayes Ball

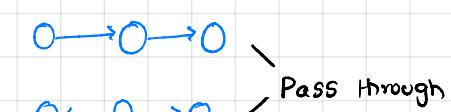
Given a directed model, is $x_i \perp x_j | X_A$?

If yes, they are d-separated

"observed"

① Take graph, color nodes in f.
Start a ball at x_i .

② Bounce ball with these rules:



Color of edge nodes doesn't matter



only center color matters



Note:
O -> O -<- O
Can be same node



③ If you hit x_j , not CI

④ If impossible to hit x_j , CI

Examples

genetics BP Age



$E \perp B \mid R$ Yes

$R \perp B \mid A$

genetics BP Age



Suppose given dataset

$$\theta_0^S = \frac{1}{2} \quad \theta_1^S = \frac{1}{2}$$

$$\theta_{0,0}^{CIS} = \frac{3}{4} \quad \theta_{1,0}^{CIS} = \frac{1}{4}$$

$$\theta_{0,1}^{CIS} = \frac{1}{2} \quad \theta_{1,1}^{CIS} = \frac{1}{2}$$

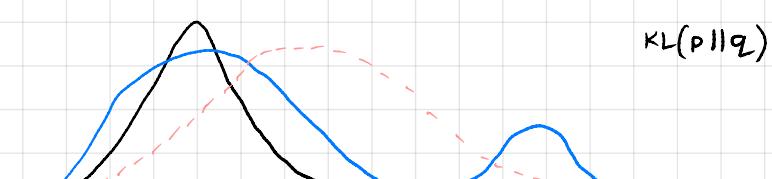
Smoker	Cancer
1	0
0	1
1	1
0	0
0	0
0	0
1	0
1	1

Example

best p
q



$KL(q \parallel p)$



$KL(p \parallel q)$

p doesn't put p over $q(x) = 0$ because of (2)

SUMMARY

DAGs \rightarrow CIS \rightarrow Factorized distribution

DAG + Bayes ball \rightarrow D separation

To discuss

- Given a dist, how to query it?
- Estimate dist from data

LECTURE 4

KL Divergence & Maximum Likelihood

Inference vs Learning

Let $\theta \rightarrow$ all params of P_θ

Inference: Query $P_\theta(A, B, \dots | C, D, \dots)$

Learning: Given data and P_θ , Fixed graph immutable DAG

Bayesian doesn't involve "learning"

MLE in Fully observed directed models

no missing data

Example

Smoker \rightarrow Cancer

$\Theta \rightarrow$ vector $\in \mathbb{R}^6$

$$\Theta = [\theta_0^S, \theta_1^S, \theta_{0,0}^{CIS}, \theta_{0,1}^{CIS}, \theta_{1,0}^{CIS}, \theta_{1,1}^{CIS}]$$

$$P_\theta(S) = \begin{cases} \theta_0^S, & S=0 \\ \theta_1^S, & S=1 \end{cases}$$

$$P_\theta(C|S) = \begin{cases} \theta_{0,0}^{CIS} \\ \theta_{0,1}^{CIS} \\ \theta_{1,0}^{CIS} \\ \theta_{1,1}^{CIS} \end{cases}$$

For fixed q.

(1) For $KL(q \parallel p)$ to be small, p must put prob. at all x s.t. $q(x) > 0$

(2) For $KL(p \parallel q)$ to be small, q can only put prob. at where $p(x) > 0$

$$\sum p \log \frac{p}{q}$$

All distributions
Distributions we can represent P_θ
Possible true distribution
Another true dist

Want to find best approximation to true distribution

- Lower is better
- Not a distance metric
- Measures the information "lost" when using encoding of one dist to another

Triangle property does not hold

Intuition

Maximum Likelihood

Given $x^{(1)}, \dots, x^{(m)}$, minimize KL divergence wrt P^*

$$(1) \quad KL(P^* \parallel P_\theta) = \sum_x P^*(x) \log \frac{P^*(x)}{P_\theta(x)}$$

$$(2) \quad \approx \frac{1}{m} \sum_{j=1}^m \log \frac{P^*(x^{(j)})}{P_\theta(x^{(j)})}$$

(3) From Monte Carlo

$$= \frac{1}{m} \sum_{j=1}^m \log P^*(x^{(j)}) - \underbrace{\frac{1}{m} \sum_{j=1}^m \log P_\theta(x^{(j)})}_{\text{constant}}$$

$$(4) \quad \underset{\theta}{\operatorname{argmin}} \quad KL(P^* \parallel P_\theta) = \underset{\theta}{\operatorname{argmax}} \sum_{j=1}^m \log P_\theta(x^{(j)})$$

$$\underset{\theta}{\operatorname{argmax}} \prod_{j=1}^m P_\theta(x^{(j)}) \rightarrow \text{maximum likelihood}$$

Mimimizing the Monte Carlo approximation of KL-Div of model & true dist is the max likelihood estimate θ

Learning Bernoulli

\sim why counting is enough

$$\text{① Data } x^{(1)} \dots x^{(m)} \text{ IID}, P_\theta(x) = \begin{cases} \theta_0^x & x=0 \\ \theta_1^x & x=1 \end{cases}$$

$$\text{② Log likelihood: } \sum_{i=1}^m \log P_\theta(x^{(i)}) = \#_0 \log \theta_0^x + \#_1 \log \theta_1^x$$

want to maximize under $\theta_0^x + \theta_1^x = 1$ $\theta_0^x \geq 0$ $\theta_1^x \geq 0$

use Lagrangian multipliers

$$\mathcal{L}(\theta, \lambda) = \#(0) \log \theta_0^x + \#(1) \log \theta_1^x + \lambda(1 - \theta_0^x - \theta_1^x)$$

$$\frac{\partial \mathcal{L}}{\partial \theta_0^x} = \frac{\#(0)}{\theta_0^x} - \lambda \quad \frac{\partial \mathcal{L}}{\partial \theta_1^x} = \frac{\#(1)}{\theta_1^x} - \lambda$$

LECTURE 5 Extended MLE & Undirected models

$$\theta_{\text{sc}}^x = \frac{\#(x)}{m} \rightarrow \text{After lagrangian}$$

Learning fully-observed discrete directed model via MLE

$$\text{① Data } x^{(1)} \dots x^{(m)} \quad x_i \in \mathbb{R}^n \quad \xrightarrow{\text{R-V}} \text{discrete} \quad \xrightarrow{\text{labels}}$$

$$\text{② } P_\theta(x) = \prod_{i=1}^n P_\theta(x_i | X_{\text{pa}(i)}) = \prod_{i=1}^n \theta_{x_i | X_{\text{pa}(i)}} \quad \xrightarrow{\text{values}}$$

Log Likelihood

$$\sum_{j=1}^m \log P_\theta(x^{(j)}) = \sum_{j=1}^m \sum_{i=1}^n \log P_\theta(x_i^{(j)} | X_{\text{pa}(i)})$$

$$= \sum_{j=1}^m \sum_{i=1}^n \log \theta_{x_i^{(j)} | X_{\text{pa}(i)}}^{(j)}$$

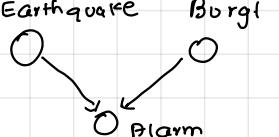
$$\star = \sum_{i=1}^n \sum_{x_{\text{pa}(i)}} \#(x_i, X_{\text{pa}(i)}) \log \theta_{x_i | X_{\text{pa}(i)}}^{(j)}$$

$\xrightarrow{\text{count of configs}}$

configurations occur many times in data.

$\xrightarrow{\text{Summing over all configs}}$

Eg



$$\begin{array}{c|cc|c} A & E & B \\ \hline 0 & 0 & 0 \\ 1 & 0 & 0 \\ 0 & 0 & 0 \end{array} \quad \xrightarrow{\sum_{j=1}^3 \log \theta_{x_j | e,b}} \quad \log \theta_{F \mid E,B} + \log \theta_{F \mid F,F} + \log \theta_{F \mid E,B}$$

Constraints: $\forall i, \forall_{\text{pa}(i)}, \sum_{x_i} \theta_{x_i | X_{\text{pa}(i)}} = 1$

Solving using Lagrangian

$$\begin{aligned} \mathcal{L} &= \sum_i \sum_{x_i} \sum_{X_{\text{pa}(i)}} \#(x_i, X_{\text{pa}(i)}) \log \frac{\theta_{x_i | X_{\text{pa}(i)}}}{\sum_{x_i} \theta_{x_i | X_{\text{pa}(i)}}} \\ &\quad + \sum_i \sum_{X_{\text{pa}(i)}} \lambda_i(X_{\text{pa}(i)}) \left(1 - \sum_{x_i} \theta_{x_i | X_{\text{pa}(i)}} \right) \\ &= \frac{\#(x_i, X_{\text{pa}(i)})}{\theta_{x_i | X_{\text{pa}(i)}}} - \frac{x_i | X_{\text{pa}(i)}}{\sum_{x_i} \theta_{x_i | X_{\text{pa}(i)}}} \\ \Rightarrow \theta_{x_i | X_{\text{pa}(i)}} &= \frac{\#(x_i, X_{\text{pa}(i)})}{\lambda_i(X_{\text{pa}(i)})} \\ &\downarrow \\ &= \frac{\#(x_i, X_{\text{pa}(i)})}{\sum_i \#(x_i, X_{\text{pa}(i)})} = \frac{\#(x_i, X_{\text{pa}(i)})}{\#(X_{\text{pa}(i)})} \end{aligned}$$

Note: $\theta_{x_i | X_{\text{pa}(i)}}$ can be computed by counting

Summary

- ① maximizing MLE \equiv minimize Monte Carlo approx. of KL divergence
- ② minimizing KL is reasonable even with modelling error

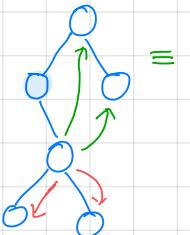
- ③ To maximize likelihood with fully observed discrete data

$$P_\theta(x_i | X_{\text{pa}(i)}) = \theta_{x_i | X_{\text{pa}(i)}} = \frac{\#(x_i, X_{\text{pa}(i)})}{\#(X_{\text{pa}(i)})}$$

Undirected models

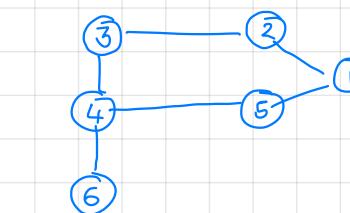
Asymmetry in directed models

- $x_i \perp x_{\text{nd}(i)} | X_{\text{pa}(i)}$
- Do not assert $x_i | X_{-i} | X_{\text{pa}(i)}$ except i



$x_i \rightarrow$ Production of regn i
How to model?

An undirected model asserts $x_A \perp x_B | x_C$ if C separates A & B



$x_1 \perp x_3 x_2$	X
$x_1 \perp x_3 x_2, x_4$	✓
$x_{6,1} \perp x_{3,5} x_2$	✗

when $P(x) > 0$ on undirected models asserts

$$P(x_i | x_{-i}) = P(x_i | X_{\text{nbr}(i)})$$

Formula for $P(x)$? $\equiv x_i \perp X_{-i} | X_{\text{nb}(i)}$

Hammersley - Clifford theorem

A positive dist. $P(x) > 0$ satisfies CI if and only if it can be written as

$$P(x_1, \dots, x_n) = \prod_{C \in G} \Phi_C(x_C)$$

↑ set of cliques

LECTURE 6 Undirected models

Undirected Graphs assert

- ① $x_A \perp x_B | x_C$ if C separates A & B
- ② $P(x_i | x_{-i}) = P(x_i | X_{\text{nbr}(i)})$

HC Theorem

A positive dist. ($P(x) > 0$) satisfies CI of UG iff P can be written as

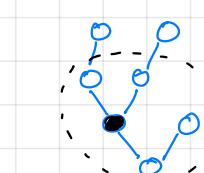
$$P(x) = \frac{1}{Z} \prod_{C \in C} \Phi_C(x_C)$$

↑ set of max factor

We know $x_1 \perp x_4 | x_2, x_3$
 $\therefore P(x) = \frac{1}{Z} \Phi_{12}(x_{1,2}) \cdot \Phi_{13}(x_{1,3}) \cdot \Phi_{24}(x_{2,4}) \cdot \Phi_{34}(x_{3,4})$

Markov blanket

Set of conditions that make x_i CI from everyone else
MB of node i is the only knowledge needed to predict the behaviour of that node & its children



HC Proof

Define $x^* = (0, 0 \dots)$ $\Phi(x) = \ln \left(\frac{P(x)}{P(x^*)} \right)$

Step 1 We can uniquely write Φ as

$$\begin{aligned}\Phi(x) &= \sum_i x_i G_i(x_i) + \sum_{i,j} x_i x_j G_{ij}(x_i, x_j) \\ &\quad + \sum_{i,j,k} x_i x_j x_k G_{ijk}(x_i, x_j, x_k) - \text{triples} \\ &\quad \dots \\ &\quad + \sum_{i_1, i_2, \dots, i_n} G_{i_1, i_2, \dots, i_n}(x_1, x_2, \dots, x_n) \quad (\mathbf{x} \in \mathbb{R}^n)\end{aligned}$$

↑ think tables

Finding G_i

All other terms disappear

Set $\mathbf{x} = (1, 0, 0 \dots) \rightarrow G_1(1)$

$\mathbf{x} = (2, 0, 0 \dots) \rightarrow G_1(2)$

$\mathbf{x} = (0, 1, 0 \dots) \rightarrow G_2(1)$

Finding G_{ij}

We know G_i , substitute

Set $\mathbf{x} = (1, 1, 0 \dots) \rightarrow G_{12}(x_1, x_2)$

Upshot → Easy to reason about Φ , given G_*

Step 2

Define $x^i \rightarrow i^{\text{th}} \text{ component} \leftarrow 0$ ($x_1, x_2 \dots 0, x_{i+1} \dots$)

Then $\exp(\Phi(\mathbf{x}) - \Phi(x^i)) \equiv \frac{P(\mathbf{x})}{P(x^i)} = \frac{P(x_i | \mathbf{x}_{-i})}{P(x_i^0 | \mathbf{x}_{-i})}$ (1)
 ↓ only depends on nbr(i)
 $\frac{P(x_i | \mathbf{x}_{-i})}{P(x_i^0 | \mathbf{x}_{-i})} = \frac{P(x_i | \mathbf{x}_{\text{nbr}(i)})}{P(x_i^0 | \mathbf{x}_{\text{nbr}(i)})} = \frac{P(x_i | \mathbf{x}_{\text{nbr}(i)})}{P(x_i^0 | \mathbf{x}_{\text{nbr}(i)})}$

Step 3

Pick node 1 without loss of generality (WLOG)

$$\Phi(x) - \Phi(x') = x_1 \left(G_1(x_1) + \sum_{i \neq j} x_j G_{ij}(x_{ij}) + \sum_{i \neq j \neq k} \dots \right)$$

will only have terms involving x_1

Step 4

Suppose j is not a neighbour of 1. ^{UG assertion}
 Then all terms involving x_j must be zero

→ Terms containing x_j in decomposition are 0

⇒ The only surviving terms are from cliques why?
 Because

$\Phi(x) - \Phi(x')$ doesn't depend on x_j

Because $P(x_i | \mathbf{x}_{-i})$ doesn't depend on j (1)
 j is not nbr, & UG assertion

Set $x_k = 0 \forall k \notin \{i, j\}$. Set everything except $i, j \rightarrow 0$

We get, $x_1 (G_1(x_1) + x_j G_{ij}(x_i, x_j))$

If $G_{ij} \neq 0$ Contradiction { .. $\Phi(x) - \Phi(x')$ doesn't depend on x_j

⇒ If j not a neighbour of 1 then $G_{ij} = 0$

For x_1, x_j, x_k we can argue similarly

Example



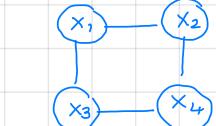
$$x_1 \perp x_3 | x_2$$



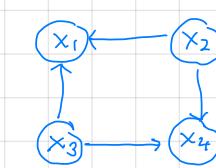
$$x_3 \perp x_1 | x_2$$



$$x_1 \perp x_3 | x_2$$



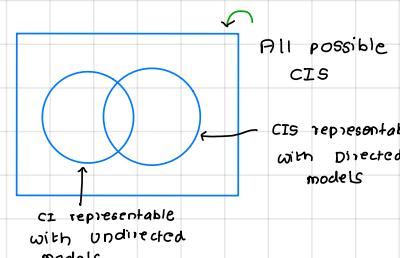
$$x_1 \perp x_4 | x_{2,3}$$



$$x_2 \perp x_3 | x_{1,4}$$

Not possible
to get CIS as directed

LECTURE 7 Markov Random Fields & CRFs



Bayesian nets → Represent distributions

Markov Random Field

$$P(x_1, x_2, \dots, x_n) = \frac{1}{Z} \prod_{c \in \mathcal{G}} \phi_c(x_c)$$

ϕ_c - arbitrary positive functions

↑ normalization factor

$$Z = \sum_{x_c} \prod_{c \in \mathcal{G}} \phi_c(x_c)$$

If we want to emphasize params

$$\Phi_c(x) = \phi_c(x | \theta) \quad \text{can depend on } \theta$$



$$\phi_c(x | \theta)$$

cliques of size 2

Ising model

$$x_i \in \{-1, 1\} \quad P(x) = \frac{1}{Z} \prod_{i,j \in \text{pairs}} \exp(\theta_{ij} x_i x_j)$$

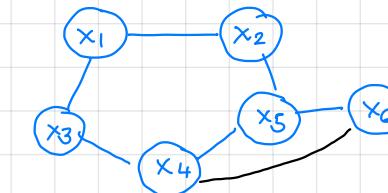
θ_{ij} is how much x_i on x_j "want" to have same value

Ising model is Markov Random Field (MRP)

$$\phi_{\{i,j\}}(x_i, x_j) = \begin{cases} e^{\theta_{ij}} & x_i = x_j \\ e^{-\theta_{ij}} & x_i \neq x_j \end{cases}$$

Conditional Random Field

Consider



$$P(x) = \frac{1}{Z} \Phi_{1,2}(x_1, x_2) \Phi_{1,3}(x_1, x_3) \Phi_{3,4}(x_3, x_4)$$

$$\Phi_{2,5}(x_2, x_5) \Phi_{4,5}(x_4, x_5) \Phi_{5,6}(x_5, x_6)$$

what is $P(x_1, x_2, x_3) x_{4,5,6}$ Fix $x_{4,5,6}$ & renormalize
 $\Phi_{4,5}, \Phi_{5,6}$ dropped → observed

$$= \frac{1}{Z(x_4, x_5, x_6)} \Phi_{1,2}(x_1, x_2) \Phi_{1,3}(x_1, x_3) \Phi_{3,4}(x_3, x_4) \Phi_{2,5}(x_2, x_5)$$

$$Z(x_4, x_5, x_6) = \sum_{x_1} \sum_{x_2} \sum_{x_3} \Phi_{1,2} \Phi_{1,3} \Phi_{3,4} \Phi_{2,5}$$

$x_{4,5}$ terms can

Now if we add edge 4,6, conditional distribution remains same

⇒ some modelling choice doesn't matter if we only care about conditional distribution

$$P(Y|x) = \frac{1}{Z(x)} \prod_{c \in \mathcal{G}} \phi_c(x_c, y_c)$$

$$Z(x) = \sum_y \prod_{c \in \mathcal{G}} \phi_c(x_c, y_c)$$

What you get from MRF & conditioning

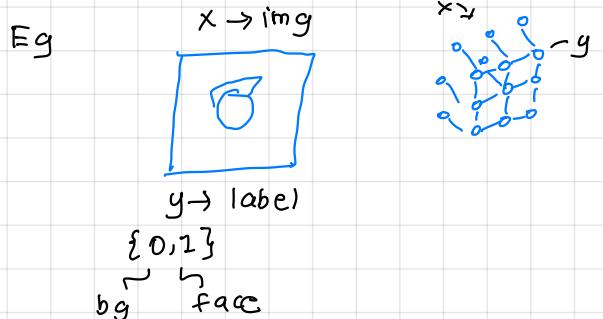
No need to define ϕ_c for cliques that only involve $x \rightarrow$ Only y also influences dist

No need to worry about curse of dimensionality (COD) of X , only y ?

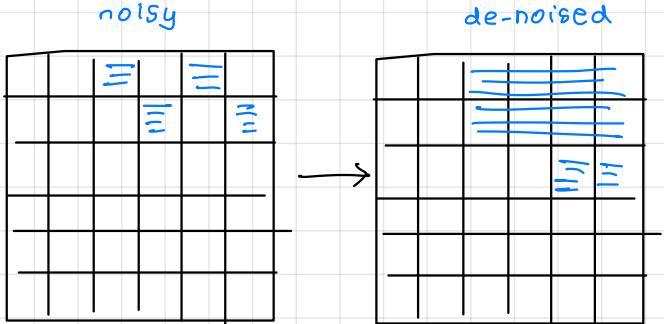
$$P_w(y|x) = \frac{1}{Z(x,w)} \prod_{c \in C} \phi_c(y_c, x; \omega)$$

$$Z(x,w) = \sum_y \prod_{c \in C} \phi_c(y_c, x; \omega)$$

ϕ_c parameterized in application dependent way
can depend on full input



Example \rightarrow Noisy Ising CRF



$$x_i \in \{-1, 1\}$$

$$y_i \in \{-1, 1\}$$

controls how much neighbouring pixels have same value

wants output pixel to match input pixel

$$P_w(y|x) = \frac{1}{Z(x,w)} \prod_{(i,j) \text{ pairs}} \exp(\omega_1 y_i y_j) \prod_i \exp(\omega_2 y_i x_i)$$

ω_1 controls how much neighbouring pixels have same value

wants output pixel to match input pixel

Why use a CRF?

- we can fit on MRF $P_\theta(x, y)$ OR

- Fit on CRF $P_\theta(y|x)$

Suppose $P_\theta(x)$ is really complex
but $P_\theta(y|x)$ is simple \rightarrow CRF avoids modelling x

Any reason to use MRF?

Yes \rightarrow can have superior sample complexity

Summary

- ① $x_A \perp x_B | x_C$ if C separates A & B \rightarrow NO direct probabilistic interpretation
- ② HC Hm $\rightarrow P(x) = \frac{1}{Z} \prod_{c \in C} \phi_c(x_c)$
- ③ MRFs $\rightarrow P(x) = \frac{1}{Z} \prod_{c \in C} \phi_c(x_c)$
- ④ CRFs $\rightarrow P(y|x) = \frac{1}{Z(x,w)} \prod_c \phi_c(y_c, x; \omega)$

Lecture 8 Ising model

Generative model $\rightarrow P_\theta(x) = \frac{1}{Z_1} \prod_i \exp(\theta_{ij} x_i x_j)$

$$P_\theta(y|x) = \frac{1}{Z(x,w)} \prod_{(i,j)} \exp(w_i y_i y_j) + \exp(\omega_2 y_i x_i)$$

clear input pixels \downarrow pairs \downarrow noisy ising

Inference with model

Why do we need inf.

- Patient - nausea, headache
 $P(\text{ate food at Joe} | \text{nausea, headache})$

$P(j, n, h) = \sum_v \sum_f P(v, f, p, h, \dots)$

\downarrow Does not scale

virus, flu, ate food at Joe, food poisoning, nausea

Message passing on a chain

We can use normalizing constant Z ! How?



Works with directed models too

$$Z_1 = \sum_{x_1} \sum_{x_2} \sum_{x_3} \sum_{x_4} \Phi_1(x_1) \Phi_2(x_2) \Phi_3(x_3) \Phi_4(x_4)$$

$$\Phi_i(x_i) = \begin{cases} 1 & x_i = 0 \\ 2 & x_i = 1 \end{cases} \rightarrow \text{single potential}$$

$$\Phi_{ij}(x_i, x_j) = \begin{cases} 1 & x_i \neq x_j \\ 2 & x_i = x_j \end{cases} \rightarrow \text{pairwise potential}$$

Process:

Push sums to right

$$Z_1 = \sum_{x_1} \Phi_1(x_1) \sum_{x_2} \Phi_2(x_2) \Phi_{12}(x_1, x_2)$$

$$\sum_{x_3} \Phi_3(x_3) \Phi_{23}(x_2, x_3) \sum_{x_4} \Phi_4(x_4) \Phi_{34}(x_3, x_4)$$

need x_3
Define $m_{4 \rightarrow 3}(x_3)$

$$Z_1 = \sum_{x_1} \Phi_1(x_1) \sum_{x_2} \Phi_2(x_2) \Phi_{12}(x_1, x_2)$$

$$\sum_{x_3} \Phi_3(x_3) \Phi_{23}(x_2, x_3) m_{4 \rightarrow 3}(x_3)$$

need to know x_2
Define as $m_{3 \rightarrow 2}(x_2)$

$$Z_1 = \sum_{x_1} \Phi_1(x_1) \sum_{x_2} \Phi_2(x_2) \Phi_{12}(x_1, x_2) m_{3 \rightarrow 2}(x_2)$$

$m_{2 \rightarrow 1}(x_1)$

$$Z_1 = \sum_{x_1} \Phi_1(x_1) m_{2 \rightarrow 1}(x_1) \rightarrow \text{Sum to 1}$$

$$m_{4 \rightarrow 3}(x_3=0) = \Phi_4(0) \Phi_{3,4}(0,0) + \Phi_4(1) \Phi_{3,4}(0,1)$$

1 - 2 2 - 1 = 4

$$m_{4 \rightarrow 3}(x_3=1) = \Phi_4(0) \Phi_{3,4}(1,0) + \Phi_4(1) \Phi_{3,4}(1,1)$$

1 - 1 2 - 2 = 5

$$m_{4 \rightarrow 3} = [4 \quad 5]$$

m wants $x_3 = 1$

$$m_{4 \rightarrow 3}(x_3=0) \quad m_{4 \rightarrow 3}(x_3=1)$$

$$m_{3 \rightarrow 2}(x_2=0) = \Phi_3(0) \Phi_{2,3}(0,0) + \Phi_3(1) \Phi_{2,3}(0,1)$$

1 - 2 2 - 1 = 18

$$m_{3 \rightarrow 2}(x_2=1) = \Phi_3(0) \Phi_{2,3}(0,1) + \Phi_3(1) \Phi_{2,3}(1,1) \cdot m_{4 \rightarrow 3}(x_3=0)$$

$m_{4 \rightarrow 3}(x_3=0)$
 $= 24$

$$1 \times 1 \times 4 + 2 \times 2 \times 5 = 24$$

Finally we get $m_{2 \rightarrow 1}$

$$\sum_{x_i} \phi_i(x_i) m_{2 \rightarrow 1}(x_i) = \phi_1(x_1=0) \cdot m_{2 \rightarrow 1}(x_1=0) + \phi_1(x_1=1) m_{2 \rightarrow 1}(x_1=1)$$

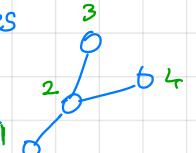
$\frac{1 \times 84}{2 \times 114} = 812$

We can compute forward chain ($L \rightarrow R$)

Define

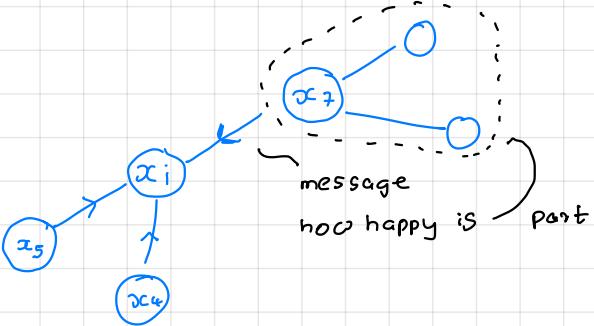
$$m_{i \rightarrow j}(x_j) = \sum_{x_i} \phi(x_i) \phi(x_i, x_j) \prod_{l \in \text{nb}(i) \setminus j} m_{l \rightarrow i}(x_i)$$

works for acyclic undirected graphs \rightarrow trees



General formula for Z

$$Z = \sum_{x_i} \phi(x_i) \prod_{j \in \text{nb}(i)} m_{j \rightarrow i}(x_i)$$



We can move sums around

$m_{7 \rightarrow i}(0)$ = How happy is subgraph starting with x_7 when x_i has value 0

What is $p(x_i)$

$$\frac{\sum_{x_2} \sum_{x_3} \sum_{x_4} (\phi_1(x_1) \phi_2(x_2) \dots \phi_4(x_4) - \phi_{12}(x_1, x_2) \dots)}{Z}$$

$$= \frac{\phi_1(x_1) m_{2 \rightarrow 1}(x_1)}{Z}$$

\Rightarrow From empirical eval, prefers to be 1

Gen formula for univariate marginals

$$p(x_i) = \frac{1}{Z} \phi(x_i) \prod_{j \in \text{nb}(i)} m_{j \rightarrow i}(x_i)$$

NOTES

factor

$$\tilde{p}(A, B, C, D) = \frac{\phi(A, B) \phi(B, C)}{\phi(C, D) \phi(D, A)}$$

score
(any func)

$$\phi(x, y) = \begin{cases} 0 & x=y=1 \\ 5 & x=y=0 \\ 1 & \text{o.w.} \end{cases}$$

probability dist

Indicate a level of coupling b/w variables

MRF

$$p(x_1, \dots, x_n) = \frac{1}{Z} \prod_{c \in C} \phi_c(x_c)$$

Advantages:

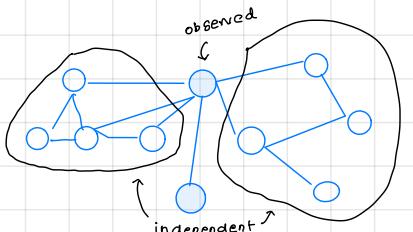
- Applied to many applications \rightarrow no directionality

Drawbacks

Z calculation expensive

Difficult to interpret

Independencies:

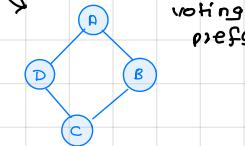


CRF

$$P(y|x) = \frac{1}{Z(x)} \prod_{c \in C} \phi_c(x_c, y_c)$$

partition func [depends on x]

friends \rightarrow tend to have similar voting prefs

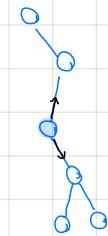


Summary

Take pair-wise MRF $p(x) = \frac{1}{Z} \prod_{i \in I} \phi_i(x_i) \prod_{(i,j) \in E} \phi_{ij}(x_i, x_j)$

Pass messages

$$m_{i \rightarrow j}(x_j) = \sum_{x_i} \phi(x_i) \phi(x_i, x_j) \prod_{l \in \text{nb}(i) \setminus j} m_{l \rightarrow i}(x_l)$$



We can find the dist. for x_i & pairs

$$p(x_i) = \frac{1}{Z} \phi(x_i) \prod_{l \in \text{nb}(i)} m_{l \rightarrow i}(x_l)$$

$$p(x_i, x_j) = \frac{1}{Z} \phi(x_i) \phi(x_j) \phi(x_i, x_j) \prod_{l \in \text{nb}(i) \setminus j} m_{l \rightarrow i}(x_l) \prod_{l \in \text{nb}(j) \setminus i} m_{l \rightarrow j}(x_l)$$

Discussion:

- WORKS FOR ANY TREE \rightarrow pass messages in right order

- Messages can get very large

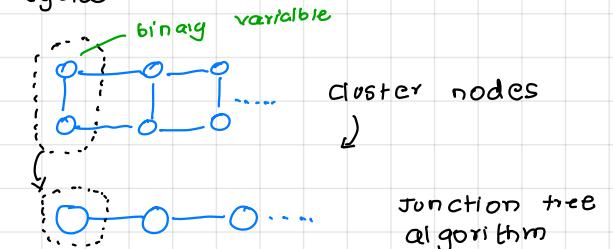
$\begin{cases} \text{use log-space} \\ \text{normalize messages so that } \sum_{x_j} m_{i \rightarrow j}(x_j) = 1 \end{cases}$

- we can use CRFs \rightarrow no change, fix values

- Can extend to graphs with larger cliques

- Can use same basic alg for directed models

For cycles



Tree width
 \hookrightarrow no. of nodes to cluster to form tree

You can use message-passing as a heuristic

Loopy Belief Propagation

Iterate graph till graph converges



Exponential Families

- VG are cases of exp fam

$$P_\theta(x) = h(x) \exp(\theta^T T(x) - A(\theta))$$

makes dist sum to 1
↳ feature extractor

θ - natural parameters

$T(x)$ - sufficient statistics

$A(\theta)$ - log partition function

$h(x)$ - scaling constant (usually 1)

↳ does not depend on θ , constant

$$A(\theta) = \log \int h(x) \exp(\theta^T T(x)) dx$$

$$= \log \sum_x h(x) \exp(\theta^T T(x)) \quad *$$

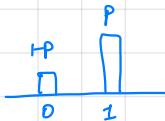
Why?

Graphical models are factorized distributions
↳ Product of simple functions

$$\exp(\theta^T T(x)) = \exp(\sum_i \theta_i T_i(x)) = \prod_i \exp(\theta_i T_i(x))$$

Bernoulli example

$$x \in \{0, 1\} \quad Ber(x) = \mu^x (1-\mu)^{1-x}$$



$$T(x) = (\mathbb{I}[x=0], \mathbb{I}[x=1])$$

$$h(x) = 1$$

$$\bar{\theta} = (\theta_1, \theta_2) = (\log \mu, \log(1-\mu))$$

$$P_\theta(x) = \exp(\theta^T T(x) - A(\theta))$$

$$A(\theta) = \log \sum_x \exp(\theta^T T(x)) = \log(\theta_0 T(0) + \theta_1 T(1)) = \log(e^{\theta_1} + e^{\theta_2})$$

$$= \exp(\theta_1 T_1(x) + \theta_2 T_2(x) - A(\theta))$$

$$= \exp(\quad)$$

Option B

$$T(x) = x \quad \theta = \log\left(\frac{\mu}{1-\mu}\right)$$

$\log \sum \exp(\theta^T T(x))$
softmax ↗
 $= \log(1 + e^\theta)$

$$P_\theta(x) = \exp(\theta^T T(x) - A(\theta))$$

$$= \exp\left(\log\left(\frac{\mu}{1-\mu}\right) - A(\theta)\right)$$

Isling model



$$P(x) = \frac{1}{Z} \prod_{(i,j)} \exp(w_i x_i x_j)$$

$$= \frac{1}{Z} \exp(w_1 x_1 x_2) \exp(\dots)$$

$$T(x) = (x_1 x_2 + x_2 x_3 + \dots) \quad \theta = w$$

Isling

$$P(x) = \frac{1}{Z} \prod_{(i,j)} \exp(w_{ij} x_i x_j)$$

$$T(x) = [x_i x_j \quad \forall (i,j) \in \text{pairs}]$$

$$\theta = [w_{ij} \quad \forall (i,j) \in \text{pairs}]$$

→ exponential family

Properties of exponential families

Critical property #1

Gradients of $A(\theta)$ are expected values

$$\frac{d}{d\theta} A(\theta) = \mathbb{E}_{x \sim P_\theta(x)} [T(x)]$$

$$p_\theta(x) = h(x) \exp(\theta^T T(x) - A(\theta))$$

$T(x)$: Features of x (sufficient statistics)

$\theta^T T(x)$: Scoring of x

$\exp(\theta^T T(x))$: Makes the scoring/linear combination of x always come out to be positive

$Z(\theta) = \exp(A(\theta))$: Normalizing constant to make this a distribution

$h(x)$: A term we add to make this formulation more general but the above intuition is the real takeaway

Sufficient statistic - A statistic that summarizes all the information in a sample about chosen parameter

Eg → Consider data $(1, 2, 3, 4, 5)$. It has sample mean 3.

Imagine we don't know data but know sample mean, finding out data doesn't improve estimate

Lecture 10 EXponential Family Contd.

Form:

$$P_\theta(x) = h(x) \exp(\theta^T T(x) - A(\theta))$$

↓ score
↓ to make +ve
Gradient property

Critical property #2

$$\frac{d^2 A(\theta)}{d\theta d\theta^T} = \mathbb{C}_{x \sim P_\theta(x)} [T(x)]$$

Hessian

$$\frac{d^2 A(\theta)}{d\theta d\theta^T} = \frac{d}{d\theta^T} \sum_x P_\theta(x) T(x)$$

$$= \frac{d}{d\theta^T} \sum_x h(x) \exp(\theta^T T(x) - A(\theta)) T(x)$$

$$= \sum_x h(x) \underbrace{\exp(\theta^T T(x) - A(\theta)) T(x)}_{P_\theta(x)} \frac{d}{d\theta^T} (\theta^T T(x) - A(\theta))$$

$$= \sum_x P_\theta(x) T(x) (T(x)^T - \mathbb{E}[T(x)]^T)$$

$$= [\mathbb{E}[T(x) T(x)^T] - \mathbb{E}[T(x)] \mathbb{E}[T(x)]^T]$$

Maximum Likelihood

$$\text{Fix } P_\theta(x) = h(x) \exp(\dots)$$

$$L(\theta) = \frac{1}{n} \sum_{i=1}^n \log P_\theta(x^{(i)})$$

$$= \frac{1}{n} \sum_{i=1}^n \left\{ \log(h(x^{(i)})) + (\theta^T T(x^{(i)}) - A(\theta)) \right\}$$

Data lives here

$$= \theta^T \left(\frac{1}{n} \sum_{i=1}^n T(x^{(i)}) \right) - A(\theta) + \frac{1}{n} \sum_{i=1}^n \log(h(x^{(i)}))$$

Expected value of data

↑
We only need for MLE

Hence $T(D) \rightarrow$ sufficient statistics

At the minimum $\frac{dL(\theta)}{d\theta} = 0$

$$\frac{dL}{d\theta} = \frac{1}{n} \sum_{i=1}^n (\tau(x^{(i)}) - \mathbb{E}[\tau(x)]) = 0$$

$$\Rightarrow \mathbb{E}[\tau(x)] = \frac{1}{n} \sum_{i=1}^n \tau(x^{(i)})$$

MLE finds θ s.t. Expectation under data = True expectation

Practical MLE computation

- Compute / Approximate $\mathbb{E}_{P_\theta(x)} [\tau(x)]$

- Find $\frac{dL}{d\theta} \rightarrow \theta^\top \leftarrow \theta^{\top-1} + \alpha \frac{dL}{d\theta}$

- Repeat

$$\frac{dA(\theta)}{d\theta} = [P(x_0=0, x_1=0), \dots]$$

So Far

- $P_\theta(x) = h(x) \exp(\theta^\top \tau(x) - A(\theta))$

- Bernoulli, Ising, Gaussian, Chi-squared, etc
in exp Family

- $\frac{dA(\theta)}{d\theta} = \frac{\mathbb{E}[\tau(x)]}{P_\theta(x)}$

$$\frac{d^2A}{d\theta d\theta^\top} = \frac{1}{P_\theta(x)} [\tau(x)]$$

- Likelihood $L(\theta) = \theta^\top \bar{\tau} - A(\theta) + C$

$$\bar{\tau} = \frac{1}{n} \sum_{i=1}^n \tau(x^{(i)}) \quad C = \frac{1}{n} \sum_i \log h(x^{(i)})$$

MRF as Exponential Family



binary

$$P(x) = \prod_{c \in C} \phi(x_c, x_{c'}) \quad \text{12 free params}$$

modeling $\tau(x)$ & θ

$$\tau(x) = [\mathbb{I}(x_1=0, x_2=0), \mathbb{I}(x_1=0, x_2=1), \mathbb{I}(x_1=1, x_2=0), \dots, \mathbb{I}(x_2=0, x_3=0), \dots] \quad \text{12 components}$$

$$\theta = [\theta(x_1=0, x_2=0), \dots, \theta(x_3=1, x_4=1)]$$

For a config of x , exactly 3 components in θ will be 1 corresponding to 3 cliques

$$P_\theta(x) = \exp(\theta^\top \tau(x) - A(\theta))$$

$$= \exp(\theta(x_1, x_2) + \theta(x_2, x_3) + \theta(x_3, x_4) - A(\theta))$$

$$= \frac{\exp(\theta(x_1, x_2)) \exp(\theta(x_2, x_3)) \dots}{\exp(A(\theta))} = \frac{\phi(x_1, x_2)}{Z(\theta)} \phi(x_2, x_3) \phi(x_3, x_4)$$

$$\frac{dA(\theta)}{d\theta} = \mathbb{E}[\tau(x)] = [P_\theta(x_1=0, x_2=0), P_\theta(\dots)]$$

= marginals

LECTURE 11 MLE in MRF

General MRFs as exponential families

$$P(x) = \frac{1}{Z} \prod_{c \in C} \phi_c(x_c)$$

Note: x_c clique can have many variables $x_c = \{x_1, x_2\}$

Exp family $\rightarrow \tau(x)$ has indicator funcs for every config. of every clique

$$\tau(x) = [\mathbb{I}[x_c = x_{c'}] \quad \forall c \in C \quad \forall x_{c'}] \quad *$$

$$\theta = [\theta(x_c) \quad \forall c \in C, \forall x_c]$$

$\phi_c(x_c)$ has no. of parameters = no. of configuration of x_c

$$\phi_c \begin{cases} \theta_1 & x_c = \{0, 1, 3\} \\ \theta_2 & x_c = \{0, 2, 3\} \end{cases}$$

$$\textcircled{R} \downarrow \quad \mathbb{E}_{P_\theta(x)} [\tau(x)] = \mathbb{E} [\mathbb{I}(x_c = x_{c'}) \quad \forall c \in C \quad \forall x_c]$$

$$= [P(x_c = x_{c'}) \quad \forall c \in C \quad \forall x_c]$$

model marginals

$$\frac{1}{n} \sum_{i=1}^n \tau(x^{(i)}) = \left[\frac{\#(x_c)}{n} \quad \forall c \in C \quad \forall x_c \right]$$

data marginals

Model marginal match data marginal at MLE

Conditional Exponential Family

$$P_\theta(x, y) = h(x, y) \exp(\theta^\top \tau(x, y) - A(\theta))$$

what is $P_\theta(y|x)$

$$P_\theta(y|x) = \frac{P_\theta(x, y)}{P_\theta(x)} = \frac{h(x, y) \exp(\theta^\top \tau(x, y) - A(\theta))}{\sum_{y'} h(x, y') \exp(\theta^\top \tau(x, y') - A(\theta))}$$

$A(\theta)$ cancels

$$= h(x, y) \exp(\theta^\top \tau(x, y) - \log \sum_{y'} h(x, y') \exp(\theta^\top \tau(x, y'))) \quad \text{Different exponential family form! } A'(\theta)$$

Definition

Conditional Exp. family

$$P_\theta(y|x) = h(x, y) \exp(\theta^\top \tau(x, y) - A(x, \theta))$$

L depends on x

$$A(x, \theta) = \log \sum_y h(x, y) \exp(\theta^\top \tau(x, y))$$

Note: IF $x=\{y\}$, becomes traditional exp family

$$\frac{dA(x; \theta)}{d\theta} = \sum_y P_\theta(y|x) \tau(x, y) = \frac{\mathbb{E}[\tau(x, y)]}{P(y|x)}$$

Learning in Conditional Exp Family

Given $(x^{(i)}, y^{(i)})$

$$\text{Likelihood is } L(\theta) = \frac{1}{n} \sum_{i=1}^n \log P(y^{(i)}|x^{(i)})$$

$$L(\theta) = \frac{1}{n} \left[\sum_{i=1}^n (\theta^\top \tau(x^{(i)}, y^{(i)}) - A(x^{(i)}, \theta)) + \log h(x^{(i)}, y^{(i)}) \right]$$

$$\frac{dL(\theta)}{d\theta} = \frac{1}{n} \left(\sum_{i=1}^n \tau(x^{(i)}, y^{(i)}) - \frac{d}{d\theta} A(x^{(i)}, \theta) \right)$$

At optimum $\frac{dL}{d\theta} = 0$

$$\Rightarrow \frac{1}{n} \sum_{i=1}^n \tau(x^{(i)}, y^{(i)}) = \frac{1}{n} \sum_{i=1}^n \frac{d}{d\theta} A(x^{(i)}, \theta)$$

$$= \frac{1}{n} \sum_{i=1}^n \frac{1}{P(y|x^{(i)})} \mathbb{E}[\tau(x^{(i)}, y)]$$

$$\hat{\mathbb{E}}[\tau(x, y)] = \hat{\mathbb{E}}[\mathbb{E}[\tau(x, y)] \mid x, P(y|x)] \rightarrow \text{No need to model } x \mid P(y|x)$$

we can use data expectation of x

Latent variables

$$P_{\theta}(x, u) = h(x, u) \exp(\theta^T \tau(x, u) - A(\theta))$$

unobserved
marginalize u

$$\log P_{\theta}(x) = \log \sum_u P_{\theta}(x, u) = \log \left(\sum_u h(x, u) \exp(\theta^T \tau(x, u) - A(\theta)) \right)$$

$$= \log \sum_u h(x, u) \exp(\theta^T \tau(x, u)) - A(\theta)$$

$$\log P_{\theta}(x) = A(x; \theta) - A(\theta)$$

$$P_{\theta}(x) = \exp(A(x; \theta) - A(\theta))$$

$$\frac{dA}{d\theta}(x, \theta) = \hat{\mathbb{E}}_{P_{\theta}(u|x)} [\tau(u, x)]$$

$$\frac{dA}{d\theta} = \hat{\mathbb{E}}_{P_{\theta}(x, u)} [\tau(x, u)]$$

Given $x^{(1)}, \dots, x^{(n)}$. The marginal log likelihood is

$$L(\theta) = \frac{1}{n} \sum_{i=1}^n \log P_{\theta}(x^{(i)}) = \frac{1}{n} \sum_{i=1}^n A(x^{(i)}; \theta) - A(\theta)$$

IF θ is MLE, $\frac{dL}{d\theta} = 0$

$$\frac{1}{n} \sum_{i=1}^n \hat{\mathbb{E}}_{P_{\theta}(u|x)} \tau(x^{(i)}, u) = \hat{\mathbb{E}}_{P_{\theta}(x, u)} \tau(x, u)$$

$$\hat{\mathbb{E}}_x \hat{\mathbb{E}}_{P_{\theta}(u|x)} \tau(x, u) = \hat{\mathbb{E}}_{P_{\theta}(x, u)} \tau(x, u)$$

Using the model to
fill in data
Different than prev

LECTURE 12 LEARNING THEORY

$$\hat{\mathbb{E}}_x [f(x)] = \frac{1}{n} \sum_{i=1}^n f(x^{(i)})$$

$$\mathbb{E}_{P_{\theta}(x)} [f(x)] = \sum_x P_{\theta}(x) f(x)$$

Why not instead of $\max_{\theta} \frac{1}{n} \sum_{i=1}^n \log P_{\theta}(x^{(i)})$

$$\text{use } \min_{\theta} \left\| \hat{\mathbb{E}}_x [\tau(x)] - \mathbb{E}_{P_{\theta}(x)} [\tau(x)] \right\|_2^2$$

Fill as much model on left & as much data on right?

Max Likelihood vs Conditional Max likelihood

Data $\{x^{(1)}, y^{(1)}, \dots, x^{(n)}, y^{(n)}\}$

Create model $P_{\theta}(x, y)$

Only care about $P_{\theta}(y|x)$ accuracy Choosing is subtle!

Could maximize

$$\textcircled{1} \quad \frac{1}{n} \sum_{i=1}^n \log P_{\theta}(x^{(i)}, y^{(i)}) \quad \text{— Full likelihood}$$

$$\textcircled{2} \quad \frac{1}{n} \sum_{i=1}^n \log P_{\theta}(y^{(i)} | x^{(i)}) \quad \text{— cond. likelihood}$$

We know:

$$\frac{1}{n} \sum_{i=1}^n \log P_{\theta}(x^{(i)}, y^{(i)}) = \frac{1}{n} \sum_{i=1}^n \log P_{\theta}(y^{(i)} | x^{(i)}) + \frac{1}{n} \sum_{i=1}^n \log P_{\theta}(x^{(i)})$$

Which should we use?

IF $P_{\theta}(y|x)$ and $P_{\theta}(x)$ use separate parts of θ , doesn't matter

/ no error, for some θ

IF model is exact, both find true $P^*(y|x)$. But full likelihood has less variance.

IF model has error — full likelihood with infinite data finds

$$\arg \min_{\theta} \sum_{x, y} P^*(x, y) \log \frac{P^*(x, y)}{P_{\theta}(x, y)} = \arg \min_{\theta} KL(P^*, P_{\theta})$$

Conditional likelihood finds

$$\arg \min_{\theta} \sum_{x, y} P^*(x, y) \log \frac{P^*(y|x)}{P_{\theta}(y|x)} = \arg \min_{\theta} KL(P^*, P_{\theta})$$

— conditional divergence

Gradient ascent

$$\nabla L(\theta) = \hat{\mathbb{E}}_{x, y} [\tau(x, y)] - \mathbb{E}_{P_{\theta}(x, y)} [\tau(x, y)]$$

Joint

Inference over x & y together

For conditional likelihood

$$\nabla L(\theta) = \hat{\mathbb{E}}_x \tau(x, y) - \mathbb{E}_{P_{\theta}(y|x)} [\tau(x, y)]$$

Repeated n times

Inference over y only with x fixed

Summary

Next

- True dist P_{θ_0}
- Data $x^{(1)}, \dots, x^{(n)} \sim P_{\theta_0}$ Don't know θ_0
- maximize likelihood to find $\hat{\theta}_n$
- $\hat{\theta}_n$ will be different from θ_0

How different?

LECTURE 13 FISCHER INFORMATION

We have true dist $P_{\theta}(x)$. Data $x^{(1)}, \dots, x^{(n)}$ from $P_{\theta}(x)$
don't know θ

MLE $\rightarrow \hat{\theta}_n$

How different are θ_0 & $\hat{\theta}_n$

Central limit theorem

If x_1, x_2, \dots are IID variables with mean μ and covariance Σ then $\bar{x}_n = \frac{1}{n} \sum_{i=1}^n x_i$

$$\bar{x}_n \stackrel{D}{\sim} \mathcal{N}(\mu, \frac{1}{n} \Sigma)$$

$$\sqrt{n} (\bar{x}_n - \mu) \xrightarrow{\text{"rate" of convergence}} \mathcal{N}(0, \Sigma)$$

L converges in distribution

Fisher Information

$$x \sim P_{\theta_0}(x)$$

Claim: $\mathbb{E}_{P_{\theta_0}(x)} [\nabla_{\theta} \log P_{\theta_0}(x)] = 0$

Proof:

$$\sum_x P_{\theta}(x) = 1$$

$$\nabla_{\theta} \sum_x P_{\theta}(x) = 0 \quad \downarrow \text{log derivative trick}$$

$$\sum_x P_{\theta}(x) \nabla_{\theta} \log P_{\theta_0}(x) = 0$$

$$\mathbb{E}_{P_{\theta}} [\nabla_{\theta} \log P_{\theta_0}(x)] = 0$$

(1)

$$I(\theta) = -\mathbb{E}_{P_{\theta}(x)} [\nabla^2 \log P_{\theta}(x)]$$

Claim: $\mathbb{V}_{P_{\theta_0}(x)} [\nabla_{\theta} \log P_{\theta_0}(x)] = -\mathbb{E}_{P_{\theta_0}(x)} [\nabla^2 \log P_{\theta_0}(x)]$

Proof: $\nabla^2 \sum_x P_{\theta}(x) = 0$ — matrix Fisher Information

$$= \nabla^T \nabla \sum_x P_{\theta}(x)$$

$$= \nabla^T \sum_x P_{\theta_0}(x) \nabla \log P_{\theta}(x)$$

$$= \sum_x \nabla \log P_{\theta}(x) (\nabla P_{\theta}(x))^T + \sum_x P_{\theta}(x) \nabla^2 \log P_{\theta}(x)$$

$$= \sum_x P_{\theta}(x) \nabla \log P_{\theta}(x) (\nabla \log P_{\theta}(x))^T \quad E[\nabla^2 \log P_{\theta}(x)]$$

$E[\log P_{\theta}(x)^2] \equiv \text{Var}(\dots)$
 $\approx E[x] = 0$

Exponential Family

$$P_{\theta}(x) = h(x) \exp(\theta^T T(x) - A(\theta))$$

$$\Rightarrow \nabla^2 \log P_{\theta}(x) = -\nabla^2 A(\theta) \Rightarrow I(\theta) = \nabla^2 A(\theta)$$

Asymptotic dist of likelihood grad at θ_0 of θ

Get θ_0 , find maximum likelihood, find gradient if grad ≈ 0 , we are at me of data

$$\nabla L_i(\theta_0) = \nabla \log P_{\theta}(x) \text{ is R.V}$$

$$\nabla L_n(\theta_0) = \frac{1}{n} \sum_{i=1}^n \nabla \log P_{\theta}(x^{(i)}) \text{ is a R.V}$$

What is dist of $\nabla L_n(\theta_0)$

$$\sqrt{n} \nabla L_n(\theta_0) \xrightarrow{D} N(0, I(\theta)) \quad *$$

Using prior bounds $\nabla L_n(\theta_0)$ is 0-mean & var $I(\theta)$

Asymptotic dist of max likelihood

$$\text{Define } \hat{\theta}_n = \arg \max_{\theta} L_n(\theta)$$

$$\text{we know } \nabla L_n(\hat{\theta}_n) = 0$$

Assume model is consistent. For large n $\hat{\theta}_n \approx \theta_0$

$$0 = \nabla L_n(\hat{\theta}_n) \approx \nabla L_n(\theta_0) + \nabla^2 L_n(\theta_0)(\hat{\theta}_n - \theta_0)$$

Taylor Exp?!

$$\sqrt{n}(\hat{\theta}_n - \theta_0) \approx (-\nabla^2 L_n(\theta_0))^{-1} \sqrt{n} L_n(\theta_0)$$

$\xrightarrow{D} N(0, I(\theta_0))$

Fixed value distribution Slutsky's thm

$$\sqrt{n}(\hat{\theta}_n - \theta_0) \xrightarrow{D} I^{-1}(\theta_0) N(0, I(\theta_0))$$

$$\sqrt{n}(\hat{\theta}_n - \theta_0) \xrightarrow{D} N(0, I(\theta_0)^{-1}) \xrightarrow{C[x]} A C[x] A^T$$

very important result in statistics

Bernoulli example

$$I(\theta) = -\mathbb{E} \left[\frac{d^2 \log P_{\theta}(x)}{d\theta^2} \right]$$

$$= -\theta \frac{d^2 \log \theta}{d\theta^2} - (1-\theta) \frac{d^2 \log (1-\theta)}{d\theta^2}$$

$$= \frac{1}{\theta} + \frac{1}{1-\theta} = \frac{1}{\theta(1-\theta)}$$

$$\text{Maximum likelihood } \hat{\theta}_n = \hat{E}[x]$$

$$\theta \rightarrow \text{high var} \rightarrow 0.5 \quad ?$$

$$\theta \rightarrow \text{low var} \rightarrow (0 \rightarrow 1)$$

If we consider tree-structured graph, we can do this quickly

LECTURE 14 CRAMER-RAO & CHOW-LIU

Recall:

$$\sqrt{n}(\hat{\theta}_n - \theta^*) \xrightarrow{D} N(0, I(\theta^*)^{-1})$$

$\hat{\theta}_n$ ML params θ^* true params

Cramer-Rao bound

1) Take $x^1, x^2, \dots, x^n \sim P_{\theta^*}$ i.i.d

2) Let $\hat{\theta}$ be any unbiased estimator of θ^*

$$\hat{\theta}(x^1, \dots, x^n) - \text{Avg of } \hat{\theta} = \theta^* - E[\hat{\theta}(x^1, \dots, x^n)] = \theta^*$$

3) Then, under some regularity assumptions

$$\mathbb{V}[\hat{\theta}] \geq \frac{1}{n} I(\theta)^{-1} \quad A \geq B \Leftrightarrow A \text{ is P.S.D}$$

Implies \rightarrow Maximum likelihood is as well as we can do

Caveats

that the true distribution has some "structure"

① Need to assume well-specified model

② max likelihood is not unbiased, but asymptotically unbiased

IF n is large, error is coming from variance, than bias

Chow-Liu Structure learning

What if we don't know the graph?

Take dataset $x^{(1)}, \dots, x^{(n)}$. How to pick graph G ?

Naive:

- Loop over all graphs — Computationally expensive

- Find max likelihood params

- Return graph with highest likelihood

It will return fully connected graph \rightarrow 1 clique

Solution:

① Limit G with multiple edges

Choosing the "best" tree

Max Likelihood on trees

Take a tree-structured MRF with "full tables"

$$P_\theta(x) = \prod_i \phi_i(x_i) \prod_{(i,j) \in T} \phi_{ij}(x_i, x_j)$$

Given T and data $x^{(1)}, \dots, x^{(n)}$ want θ

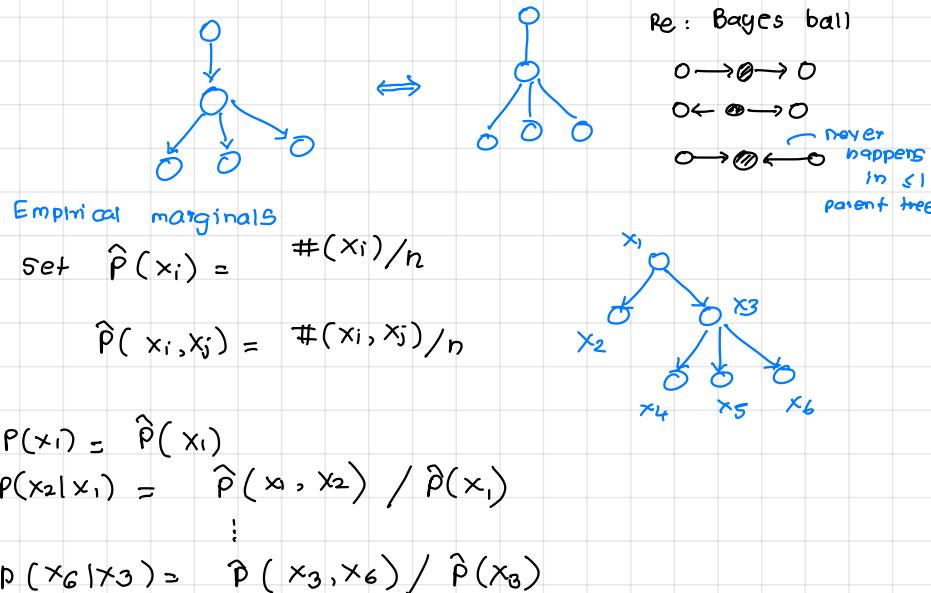
Intuition

Tree structured MRF same structured as Directed models

Can learn tree-structured models by counting

Claim

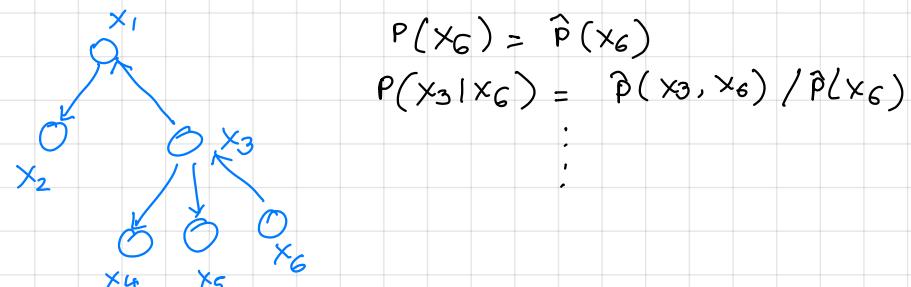
If each variable in a directed model has ≤ 1 parent then it is equivalent to model with same edges



Joint distribution

$$P(x_1, x_2, \dots) = \frac{\hat{p}(x_1) \hat{p}(x_2 | x_1) \hat{p}(x_3 | x_1, x_2) \hat{p}(x_4 | x_3, x_1, x_2) \hat{p}(x_5 | x_3, x_4, x_1, x_2) \hat{p}(x_6 | x_3, x_5, x_4, x_1, x_2)}{\hat{p}(x_3, x_5) \hat{p}(x_3, x_6)}$$

What if we had a different tree? Now,



General Formula:

$$P(x) = \prod_i \hat{p}(x_i)^{1 - |\text{nb}(i)|} \prod_{(i,j) \in T} \hat{p}(x_i, x_j)$$

Equivalence

$$P(x) = \prod_i \hat{p}(x_i) \prod_{(i,j) \in T} \frac{\hat{p}(x_i, x_j)}{\hat{p}(x_i) \hat{p}(x_j)}$$

Lecture 15 CHOW-LIU + MARKOV CHAIN MONTE CARLO

Previously:

Want tree T with max likelihood for a dataset

Reasonable approach even if true dist. not tree-structured

Learning tree-structured MRF

General Form

$$P(x) = \prod_i \hat{p}(x_i)^{1 - |\text{nb}(i)|} \prod_{(i,j) \in T} \hat{p}(x_i, x_j)$$

Equivalently

$$P(x) = \prod_i \hat{p}(x_i) \prod_{(i,j) \in T} \frac{\hat{p}(x_i, x_j)}{\hat{p}(x_i) \hat{p}(x_j)}$$

We want to find the tree

To learn a tree-structured MRF, set —

$$\phi_i(x_i) = \hat{p}(x_i)^{1 - |\text{nb}(i)|} \quad \phi_{ij}(x_i, x_j) = \hat{p}(x_i, x_j)$$

Or,

$$\phi_i(x_i) = \hat{p}(x_i) \quad \phi_{ij}(x_i, x_j) = \frac{\hat{p}(x_i, x_j)}{\hat{p}(x_i) \hat{p}(x_j)}$$

Same speed as Fully-observed directed model

What is final likelihood?

$$\begin{aligned} L(\theta) &= \sum_{i=1}^n \log P_\theta(x^{(i)}) = \hat{E} \log P_\theta(x) \\ &= \hat{E} \left[\sum_i \log \phi_i(x_i) + \sum_{i,j \in T} \log \phi_{i,j}(x_i, x_j) \right] \end{aligned}$$

The normalizing constant = 1

$$\begin{aligned} &= \hat{E} \left[\sum_i \log \hat{p}(x_i) + \sum_{i,j \in T} \underbrace{\log \frac{\hat{p}(x_i, x_j)}{\hat{p}(x_i) \hat{p}(x_j)}}_{-\text{ve entropy}} \underbrace{\log \frac{\hat{p}(x_i, x_j)}{\hat{p}(x_i) \hat{p}(x_j)}}_{\text{mutual information}} \right] \\ &= -\sum_i H_{\hat{p}}(x_i) + \sum_{i,j \in T} I_{\hat{p}}(x_i, x_j) \end{aligned}$$

does not depend on T

$$H_{\hat{p}}(x_i) = -\sum_{x_i} \hat{p}(x_i) \log \hat{p}(x_i)$$

$$I_{\hat{p}}(x_i, x_j) = \sum_{x_i, x_j} \hat{p}(x_i, x_j) \log \frac{\hat{p}(x_i, x_j)}{\hat{p}(x_i) \hat{p}(x_j)}$$

So Far:

Given dataset D and tree T , we can compute the likelihood for best MRF on tree T in terms of H and I

Define L_T = likelihood of ML params for tree T

we know

$$L_T = -\sum_i H_i + \sum_{(i,j) \in T} I_{i,j}$$

$$H_{\hat{p}}(x_i) \quad I_{\hat{p}}(x_i, x_j)$$

can be pre-calculated

To get best T ? —— max spanning tree

Only true for trees

Chow-Liu algorithm

- ① Input $x^{(1)}, \dots, x^{(n)}$
- ② Define $\hat{p}(x_i, x_j) = \frac{\#(x_i, x_j)}{n}$
- ③ For all pairs (i,j)

$$I_{i,j} = \sum_{x_i} \sum_{x_j} \hat{p}(x_i, x_j) \log \frac{\hat{p}(x_i, x_j)}{\hat{p}(x_i) \hat{p}(x_j)}$$

$$H_i = -\sum_{x_i} \hat{p}(x_i) \log \hat{p}(x_i)$$
- ④ Find max. spanning tree T with weights $I_{i,j}$
- ⑤ Return $P(x) = \prod_i \hat{p}(x_i) \prod_{(i,j) \in T} \frac{\hat{p}(x_i, x_j)}{\hat{p}(x_i) \hat{p}(x_j)}$
- ⑥ Optional: Find likelihood $-\sum_i H_i + \sum_{i,j \in T} I_{i,j}$

LECTURE 17 Mc MC contd

Idea:

Input P dist. → want to estimate $\mathbb{E}_{P(x)} f(x)$

Create a markov chain with stationary dist p.
Run markov chain for long time - $x_1, \dots, x_{t_{\max}}$

Use approximation $\mathbb{E}_{P(x)} f(x) \approx \frac{1}{t_{\max}} \sum_{t=1}^{t_{\max}} f(x_t)$

Note: $x_1, \dots, x_{t_{\max}}$ not i.i.d. → they are correlated
 okay as long as not too correlated decays
 not sampled from P, but from close to P if chain run for long time

Questions:

- ① How to create markov chain with stationary dist P
- ② How to be sure that p is the only stationary p
- ③ How long to run the chain?

Note for a T there can be multiple stationary distribution. T = Identity every p is ↗

Stationary distribution & Linear Algebra

π is stationary on T if $\pi(y) = \sum_{x \in E} \pi(x) T(y|x)$
 $\pi \equiv \pi T$

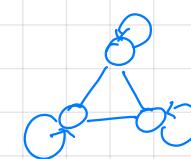
Principle,

- ① T has unique stationary dist iff it has 1 eigenvector with eigenvalue 1
- ② Speed of convergence is determined by spectral gap ↗ diff b/w 1st & 2nd largest eigen value

Can only do this when state space E is small → in which case we can sample exactly with SVD
 → our applications require large E



Imagine



D-B does not hold
stationary dist possible

Regularity

A markov chain is regular if there exists a finite value t means
 $\forall i,j \quad (T^t)_{ij} > 0$
 - we can start from any state and can go to any state with non-zero prob in t iterations

Theorem

A regular markov chain has a unique stationary dist

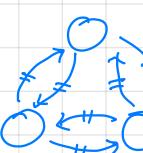
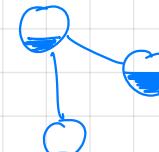
Detailed Balance

A markov chain T satisfies detailed balance w.r.t distribution π if $\forall x, x'$

$$\pi(x) T(x'|x) = \pi(x') T(x|x')$$

Intuition

- Put $\pi(x)$ "water" at each node
- For each node x send a fraction $T(x'|x)$ of its water to x'
- π is stationary if total amount of water at each node is unchanged
- For detailed balance, flow along each pair must be equal
- If T satisfies D.B. w.r.t π , then π is a stationary dist of T



Proof:

Let $\pi' = \pi T$ be result of running chain once

$$\begin{aligned} \text{Then } \pi'(x') &= \sum_x \pi(x) T(x'|x) \quad (\text{definition}) \\ &\equiv \sum_{x'} \pi(x') T(x|x') \quad (\text{def. of DB}) \end{aligned}$$

Gibbs Sampling

Init $x^{(0)}$
 For $t = 1, 2, \dots, t_{\max}$

Randomly pick $i \in \{1, \dots, D\}$ $D = \text{len}(x)$
 Sample $r \sim P(x_i | x_{-i})$ ← full conditional
 \ sample one dimen

Set $x_i \leftarrow r$

Record $x^{(t)} \leftarrow x$

Return $x^{(1)}, \dots, x^{(t_{\max})}$

T is huge, exp in # vars

Questions

- ① Is p stationary? → Detailed balance
- ② Is stationary dist unique? → Regularity

LECTURE 18 GIBBS, METROPOLIS HASTING

Linear Algebra SVD decomp not useful for large state spaces E

Regularity for Gibbs

want to show $\exists t \text{ s.t. } (T^t)_{ij} > 0 \ \forall i, j$

Assume $P(X_i=x_i | X_{-i}=x_{-i}) > 0$ always
 It means its possible to reach any point from any point in finite no. of steps
 Thus, regularity holds

Gibbs sampling detailed balance

$T_i(x'|x) \rightarrow$ Transition prob assuming i chosen

$$\equiv \mathbb{I}[x'_i = x_{-i}] P(x'_i | x_{-i})$$

$\tau(x'|x) =$ Full transition probabilities

$$= \frac{1}{D} \sum_{i=1}^D T_i(x'|x)$$

$$P(x) \tau_i(x'|x) = P(x) P(x'_i | x_{-i}) \mathbb{I}[x'_i = x_{-i}]$$

$$= P(x_{-i}) P(x_i | x_{-i}) P(x'_i | x_{-i}) \\ \mathbb{I}[x'_i = x_{-i}]$$

$$P(x') \tau_i(x|x') = P(x') P(x_i | x'_{-i}) \mathbb{I}[x_i = x'_{-i}]$$

$$= P(x'_{-i}) P(x'_i | x'_{-i}) P(x_i | x'_{-i}) \\ \mathbb{I}[x_i = x'_{-i}]$$

$$= P(x_{-i}) P(x'_i | x_{-i}) P(x_i | x_{-i}) \mathbb{I}[x_i = x'_{-i}]$$

$$= P(x) \tau_i(x'|x)$$

$$P(x) \tau(x'|x) = P(x) \frac{1}{D} \sum_{i=1}^D \tau_i(x'|x) \quad \text{Defn of } \tau \text{ i chosen uniform}$$

$$= P(x') \frac{1}{D} \sum_{i=1}^D \tau_i(x|x') \quad \text{Defn of } \tau \text{ i chosen uniform}$$

$$= P(x') \tau(x|x')$$

thus detailed balance holds for Gibbs sampling

Summary

Has unique stationary dist because regular

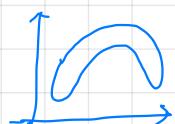
Has p as stationary dist because of detailed balance

Advantages

- Efficient per iteration
- Seems to require little training

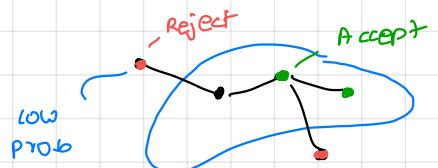
Disadvantages

- Need to sample from $P(x_i | x_{-i})$
- mixing slow when strong correlations



To overcome drawbacks

- ① In each iteration, try some kind of random move
- ② "Accept" or "Reject" carefully to guarantee regularity & detailed balance



Proposal / Acceptance MCMC

Input $p(x)$

Init x

For $t = 1, 2, \dots, t_{\max}$

- Sample $x' \sim Q(x'|x)$
- Look at x and x' and calc probability $\alpha(x', x)$ of keeping x' depend on p
- Choose $r \in [0, 1]$ uniformly
- If $r < \alpha(x', x)$, $x \leftarrow x'$
else nothing
- $x^{(t)} \leftarrow x$

Return $x^{(1)}, x^{(2)}, \dots, x^{(t_{\max})}$

How do we choose α ?

want to choose α s.t. p is stationary dist

use detailed balance

what are transition prob?

$$\text{If } x' \neq x \quad \begin{matrix} \text{Propose} \\ \downarrow \\ T(x'|x) = Q(x'|x) \alpha(x', x) \end{matrix} \quad \begin{matrix} \text{Accept} \\ \downarrow \\ \end{matrix}$$

If $x' = x$
don't care

$$P(x) \tau(x'|x) = P(x') \tau(x|x')$$

Lecture 19 METROPOLIS HASTINGS

Want an acceptance probability α that guarantees detailed balance

If $x' \neq x \Rightarrow \tau(x'|x) = Q(x'|x) \cdot \alpha(x', x)$
else \Rightarrow don't care \hookrightarrow DB always satisfied

Proposal 1 - TODO

Proposal 2

$$\alpha(x', x) = \min \left(1, \frac{P(x') Q(x|x')}{P(x) Q(x'|x)} \right)$$

Claim: DB holds

Proof:

If $P(x) Q(x'|x) \geq P(x') Q(x|x')$ then

$$P(x) \tau(x'|x) = P(x) Q(x'|x) \alpha(x', x)$$

$$= P(x) Q(x'|x) \frac{P(x') Q(x|x')}{P(x) Q(x'|x)}$$

$$= P(x') Q(x|x') \quad \text{can add because } \alpha(x, x') = 1$$

$$= P(x') Q(x|x') \alpha(x, x')$$

$$= P(x') \tau(x|x')$$

If $P(x) Q(x'|x) \leq P(x') Q(x|x')$, switch primes

\Rightarrow same results

{ one of $\alpha(x, x')$, $\alpha(x', x)$ will be 1

$$\theta_0 \sim \pi(\theta)$$

for each iteration: \hookrightarrow need to specify

$$\textcircled{1} \quad \theta_t^* \sim \mathcal{N}(\theta_{t-1}, \sigma)$$

proposed

proposal dist (should satisfy) $P(\theta_a \rightarrow \theta_b) = P(\theta_b \rightarrow \theta_a)$

$$\textcircled{2} \quad r = \frac{P(x|\theta_t^*) P(\theta_t^*)}{P(x|\theta_{t-1}) P(\theta_{t-1})}$$

$$\textcircled{3} \quad \text{if } r > u \sim U(0, 1)$$

$$\theta_t = \theta_t^*$$

else: ignore

Metropolis Hastings Algo

```

Init x
For t=1,2 ... tmax
    x' ~ Q(x'|x)
    Choose r ∈ [0,1] randomly
    IF r ≤  $\frac{P(x')Q(x|x')}{P(x)Q(x'|x)}$ , then x ← x'
        ↳ don't need min(1,...)
    xt ← x
Return x1, x2, ..., xtmax

```

$$P(x) = \frac{1}{Z} \bar{P}(x), \text{ we don't need } Z \quad \frac{P(x)}{P(x')} = \frac{\bar{P}(x)}{\bar{P}(x')}$$

Discussion:

Performance depends critically on choosing proposal distribution

How "wide" should Q(x'|x) be? ↳ want to jump far ↳ want to jump often

Paper - Weak convergence & optimal scaling { Gelman et al , 97
Acceptance rate - .234 works well in practice}

One iteration of MH touches all variables ⇒ more expensive
↳ Gibbs only 1 dim

Like all MCMC → consider mixing time

Choosing Q

Ideal Q(x'|x) → P(x) ↳ can't sample
this would "mix completely" with single accepted sample

Hamiltonian Monte Carlo (HMC)

Idea: Use gradients to hit a "far away" point that has high P(x)

HMC is Metropolis with a particular proposal

uses augmented space with "momentum"

LECTURE 20 BAYESIAN INFERENCE

Example

Gender	Hand	Tall
F	L	110
M	R	113
F	R	130

$$\omega \equiv \text{weights} \quad \omega_t = (\mu_{FL}, \mu_{FR}, \mu_{ML}, \mu_{MR})$$

$$P(g, h, t | \omega) = P(g|\omega_g) P(h|\omega_h) P(t|g, h, \omega_t)$$

$$\textcircled{1} \quad P(g|\omega_g) = \begin{cases} \omega_g & g=F \\ 1-\omega_g & g=M \end{cases} \quad \textcircled{2} \quad P(h|\omega_h) = \begin{cases} \omega_h & h=R \\ 1-\omega_h & h=L \end{cases}$$

$$\textcircled{3} \quad P(t|g, h, \omega) = \mathcal{N}(t | \mu_t, \sigma^2) \quad \hookrightarrow \text{param for diff. configs}$$

Max Likelihood Solution ↳ Frequentist

$$\hat{\omega} = \arg \max_{\omega} \sum_{i=1}^N \log P(g^{(i)}, h^{(i)}, t^{(i)} | \omega)$$

Bayesian Solution

want probability of ω given data → P(ω | Data)

$$P(\omega | \text{Data}) = \frac{P(\omega) P(\text{Data} | \omega)}{P(\text{Data})} \quad \begin{matrix} \text{Posterior} & \text{Prior} & \text{Likelihood} \\ \downarrow & \downarrow & \downarrow \\ \text{Evidence} & & \end{matrix}$$

Can calculate

$$\hookrightarrow P(\text{Data} | \omega) = \prod_{i=1}^N P(x^{(i)} | \omega)$$

P(ω) → Prior? Need to specify (Using domain knowledge)

$$P(\text{Data}) = \int_{\omega} P(\omega) P(\text{Data} | \omega) d\omega \quad \hookrightarrow \text{don't need}$$

Possible Prior

$$\textcircled{1} \quad P(\omega_g) = \begin{cases} 1 & 0 \leq \omega_g \leq 1 \\ 0 & \text{o.w.} \end{cases} = \text{Uniform}(0,1)$$

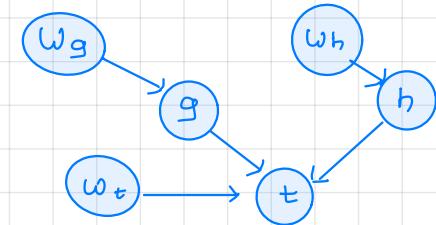
Similarly

$$\textcircled{2} \quad P(\omega_h) = \text{Uniform}(0,1)$$

$$\textcircled{3} \quad P(\omega_t) = \prod_g \prod_h \mathcal{N}(\mu_{gh} | 120, 30^2) \quad \hookrightarrow \text{Empirical?}$$

Treat $\mu_{FL}, \mu_{FR}, \mu_{ML}, \mu_{MR}$ is independent

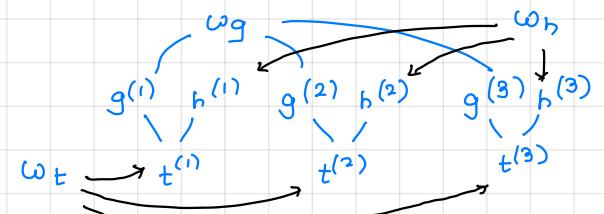
we can specify graphical model for $P(\omega, x)$



Guideline for checking model

- Specify $P(\omega, x)$
- Sample $\omega, x \sim P$
- Think of x as "synthetic dataset"
- Ask if it looks reasonable
- Repeat

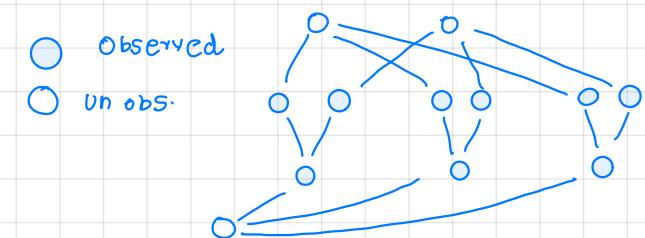
model for $P(\omega, x^{(1)}, x^{(2)}, x^{(3)})$



For $N=3$

	g ⁽ⁱ⁾	h ⁽ⁱ⁾	t ⁽ⁱ⁾
i=1	F	L	110
i=2	M	R	113
i=3	F	R	130

$P(\omega | \text{Data})$ is just the conditional dist of directed model with certain variables observed



Can draw samples using MCMC

using direct samples hard because children observed

How Bayesian Works

Design $P(\text{Data}, \omega) \approx P(\omega) \prod_i P(x_i | \omega)$
↳ Specifies a graphical model

Conditional dist $P(\omega | \text{Data})$ is posterior

use $P(\omega | \text{Data})$ in inference

↳ draw samples using MCMC

No separation of "learning" & "inference"

Game:

80 A coins .6 Head

20 B → .4 →

Pick random coin, flip 10 times $\text{reg} \approx (\omega_{\text{pred}} - \omega_{\text{true}})^2$
 $\text{Util}_{\text{ig}} \rightarrow U(\omega_{\text{pred}}, \omega_{\text{true}}) = \text{reward}$

How happy to make pred

What strategy maximizes expected reward?

Decision rule $\alpha(x^{(1)} \dots x^{(10)})$

↳ what ω_{pred} to choose for every possible dataset

Expected Utility

$$\mathbb{E} \mathbb{E} U(\alpha(x^{(1)} \dots x^{(10)}), \omega)$$

Prior $\sim P(\omega) P(x^{(1)} \dots x^{(10)} | \omega)$
 likelihood
 all datasets that can be generated using "true" coin

expectation over possible "true" coin bias

$$P(\omega = a) = .8 \quad P(x|\omega) \rightarrow \begin{array}{c} \omega \\ \hline a & b \\ \hline .4 & .6 \\ \hline 1 & .6 & .4 \end{array}$$

Best strategy claim

$$\alpha(x^{(1)}, \dots x^{(10)}) = \underset{\omega^*}{\operatorname{argmax}} \mathbb{E}_{\omega} U(\omega^*, \omega)$$

LECTURE 21 BAYESIAN INF continued

Goal: Choose decision rule (α) to make Expected Utility as high as possible

$$\text{Expected Utility } \mathbb{E}_{\omega} \mathbb{E}_{P(x| \omega)} U(\alpha(x^1 \dots x^{10}), \omega)$$

↳ how good is decision making

Best Strategy:

$$\alpha(x^1 \dots x^{10}) = \underset{\omega^*}{\operatorname{argmax}} \mathbb{E}_{\omega} U(\omega^*, \omega)$$

Why: Expr Ut is equivalently

$$\mathbb{E}_{P(x^1 \dots x^{10})} \mathbb{E}_{P(\omega | x^1 \dots x^{10})} U(\alpha(x^1 \dots x^{10}), \omega)$$

No other strategy can have higher utility on average

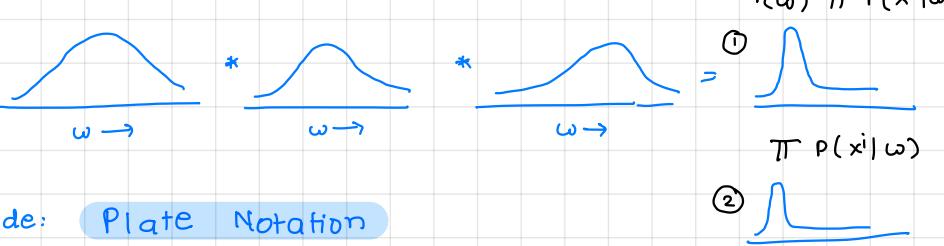
Assumptions about Bayesian

- Need to know prior $P(\omega)$

- Need to know model $P(x|\omega)$ ↗ data generation process

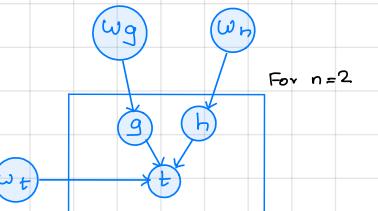
does not influence posterior when data large

$$P(\omega | \text{Data}) \propto P(\omega) \prod_{i=1}^n P(x^i | \omega)$$

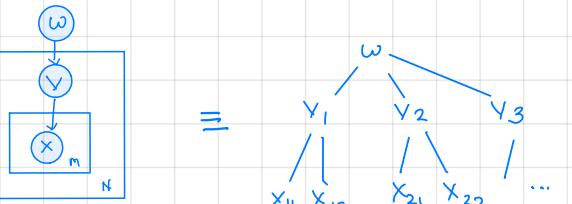


Aside: Plate Notation

Draw a box around repetitions



Plates can be nested.



Emulate MAP with Bayesian.

$$U(\omega_{\text{pred}}, \omega) = \mathbb{I}[\omega_{\text{pred}} = \omega] \quad \text{if discrete} \\ = \delta(\omega_{\text{pred}} - \omega) \quad \text{if continuous}$$

Want $P(g=F, h=R, t=110)$

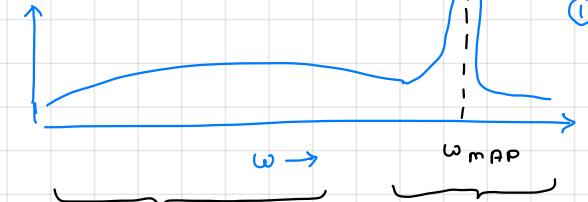
Given a specific ω , we can compute $P(g=F \dots | \omega)$

Bayesians want $\int P(\omega | \text{Data}) P(g=F, h=R, t=110 | \omega)$

Sample $\omega^1 \dots \omega^{t_{\text{max}}}$ and approximate

$$\approx \frac{1}{t_{\text{max}}} \sum_{t=1}^{t_{\text{max}}} P(g=F, h=R, t=110 | \omega^t)$$

$$P(\omega | \text{Data})$$



These ω think many right handed Females

These ω think less R, F

① Bayes will give high prob. because many samples from left area

② MAP will give low prob

Advantages

- Use domain knowledge
- Some way "optimal"

- Hard inference
- Computationally expensive

LECTURE 22 VARIATIONAL INFERENCE

General Strategy

- ① Input $P(z, x)$ and observed X
- ② Let $q_{\omega}(z)$ be a simple dist with params ω
- ③ Minimize $\text{KL}(q_{\omega}(z) || P(z|x))$

ELBO decomposition

Evidence lower bound optimization

$$\log P(x) = \sum_z q_{\omega}(z) \log \frac{P(z, x)}{q_{\omega}(z)} + \sum_z q_{\omega}(z) \log \frac{q_{\omega}(z)}{P(z|x)}$$

↓
ELBO

Proof:

$$= \sum_z q_{\omega}(z) \log P(z, x) - q_{\omega}(z) \log q_{\omega}(z) \\ + \sum_z q_{\omega}(z) \log q_{\omega}(z) - \sum_z q_{\omega}(z) \log P(z|x)$$

$q_{\omega}(z) \rightarrow$ has to be simple to sample from
 eg: Normal

Note:

- ① Can't compute KL term → don't know $P(z|x)$
- ② Can approx ELBO → $z \sim q_{\omega}$
 Compute $\log \frac{P(z, x)}{q_{\omega}}$
- ③ $\text{KL} \geq 0$
- ④ $\Rightarrow \text{ELBO} \leq \log P(x)$ (lower bound)
- ⑤ Maximizing ELBO \equiv minimize KL :: $\log P(z, x) = \text{constant}$

Uses of VI

- ① Want lower bound on $\log P(x) \rightarrow$ by learning params
- ② Want $q(z) \approx$ posterior for inference

Why not minimize $KL(p \parallel q)$

$$KL(q \parallel p) = \sum_z q(z) \log \frac{q(z)}{p(z)} \quad \text{Expectation wrt } q$$

$$KL(p \parallel q) = \sum_z p(z) \log \frac{p(z)}{q(z)} \quad \text{Exp wrt } p$$

Posterior is unknown, can't use q

VI Algorithm

Input $P(x, z)$ and fixed x

We know $ELBO(q_w \parallel p_{\theta}(x, z))$

Maximize $ELBO$ wrt w somehow

Use $q_w(z)$ as proxy for $p(x, z)$

Expectation - Maximization

Given X , want to maximize $P_{\theta}(x)$ wrt θ

$$P_{\theta}(x) = \sum_z P_{\theta}(x, z)$$

Usual derivation

$$\begin{aligned} \log P_{\theta}(x) &= \log \sum_z P_{\theta}(x, z) \\ &= \log \sum_z q(z) \frac{P_{\theta}(x, z)}{q(z)} \\ &\geq \sum_z q(z) \log \frac{P_{\theta}(x, z)}{q(z)} \quad \text{— Jensen's inequality} \end{aligned}$$

Algorithm

- Set $q(z) = P_{\theta}(z|x)$
- maximize $\sum_z q(z) \log \frac{P_{\theta}(x, z)}{q(z)}$ wrt θ
- Repeat

Produces local maximum at $P_{\theta}(x)$

Variational EM

learning & inference simultaneously

Goal: maximize $P_{\theta}(x)$

Strategy

- Define $ELBO(q_w(z) \parallel P_{\theta}(z, x))$

- we know $\forall w, \theta \quad ELBO(q_w(z) \parallel P_{\theta}(z, x)) \leq \log P_{\theta}(x)$

- maximize $ELBO$ over both w, θ simult.

Since $ELBO$ is lowerbound \rightarrow approach reasonable

IF $q_w(z) = P_{\theta}(z|x) \Rightarrow ELBO = \log P_{\theta}(x)$

Examples

- Regular EM \rightarrow where $q_w(z)$ can represent $P(z|x)$ exactly

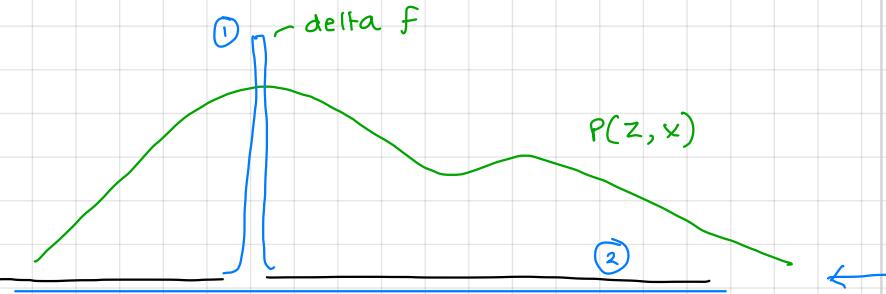
- Reg. VI \rightarrow where θ is fixed

- VAE $q_w(z)$ and P_{θ} represented with neural nets

Intuition

$$ELBO = \sum_z q_w(z) \log P(z, x) - \sum_z q_w(z) \log q_w(z)$$

↑ ↑
encourages $q_w(z)$ encourages $q_w(z)$
① to be high where ② to be random
 $P(z|x)$ is high (high entropy)



Stochastic Gradient Variational Inference

Black-box VI

Idea: Need $q_w(z)$ to sample from

- evaluate
- compute gradient

Need $P(z, x)$ - evaluate
- maybe compute gradient

- maybe compute gradient

Algorithm

- Start with some w
- For $t = 1, 2, \dots, t_{max}$

- Get unbiased estimate g of $\nabla_w ELBO$

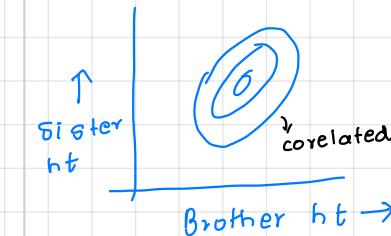
- Take gradient step - RMSprop, Adam, etc

- Return w

How to estimate $\nabla_w ELBO$?

should be cheap to compute
low variance

Lecture 23 SGVI



We have avg ht of sisters = 64 in $P(\text{Brother}, \text{Sister})$
we get to sample 1 from $P(B,S)$
want avg ht of brother

Can estimate brother
Better \rightarrow Brother + (μ - Sister)

Input $P(z, x)$. Fixed x
Want w s.t. $q_w(z) \approx P(z|x)$
posterior estimate

Idea:

For $t = 1, 2, \dots$

- Get unbiased estimate of $\nabla_w ELBO$
- Take a step
- $\approx -KL(q_w(z) \parallel p(z|x))$

$$ELBO = \int q_w(z) \log P(z, x) dz - \int q_w(z) \log q_w(z) dz$$

find high prob.
regions

q_w should have
high divergence

Often Entropy is known in closed form

For 1-d Gaussian $w = (\mu, \sigma) \rightarrow H = \frac{1}{2} \ln(2\pi e \sigma^2) - \text{ind. of } \mu$

n-d Gaussian $\rightarrow H = \frac{1}{2} \ln |2\pi e \Sigma|$
 $w = (\mu, \Sigma)$

If we have exact entropy \rightarrow can compute gradient

$$\text{Mean evidence} \rightarrow \int_z q_\omega(z) \log P(z, x) dz \\ = E_{q_\omega} [\log P(z, x)]$$

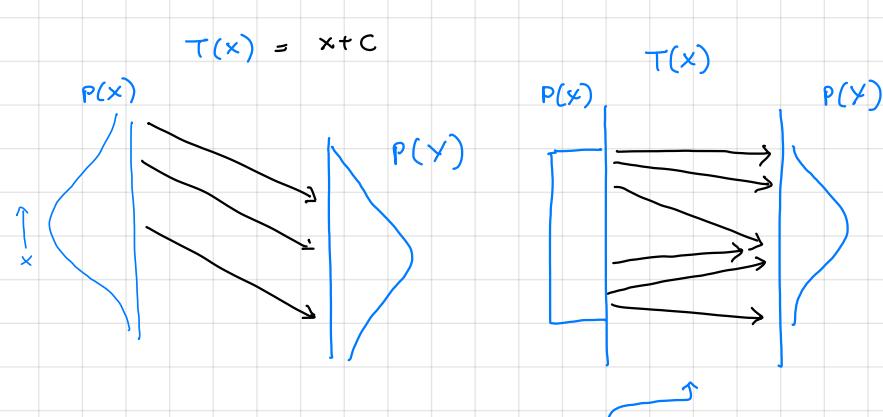
estimate by (1) drawing a single sample from $q_\omega(z)$
(2) Return $\log p(x, z)$

However, want estimate of gradient $\nabla_\omega \int_z q_\omega(z) \log P(z, x) dz$

Aside: Transformation of densities

if $Y = T(X)$ then
 \sim transformation

$$P(X=x) = P(Y=T(x)) |\det(\nabla T(x))|$$



We need det. to model non-linear

If $T(x)$ changes near x quickly

Re-parameterization trick estimator

Find - dist $q(\varepsilon)$ independent of ω
- Transformation $T_\omega(\varepsilon)$

such that $\forall \omega$

$$T_\omega(\varepsilon) \triangleq z, \quad \varepsilon \sim q(\varepsilon), \quad z \sim q_\omega(z)$$

$q(\varepsilon)$ is fixed, transform it non-linearly to get $q(\omega)$

$$\text{For any } f, \quad \underset{q_\omega(\varepsilon)}{\mathbb{E}} f(z) = \underset{q(\varepsilon)}{\mathbb{E}} f(T_\omega(\varepsilon))$$

$$\text{Thus: } \underset{q_\omega(z)}{\nabla_\omega} \underset{q_\varepsilon(z)}{\mathbb{E}} f(z) = \underset{q(\varepsilon)}{\mathbb{E}} \nabla_\omega f(T_\omega(\varepsilon))$$

Since, $q(\varepsilon)$ ind. of ω

Example: Gaussian re-parameterization

FACT

$$x \sim \mathcal{N}(0, I) \text{ then } Cx + \mu \sim \mathcal{N}(\mu, CC^T)$$

$$\omega = (\mu, C) \text{ with } \Sigma = CC^T$$

then,

$$q(\varepsilon) = \mathcal{N}(\varepsilon | 0, I)$$

$$T_\omega(\varepsilon) = C\varepsilon + \mu$$

Compute $\nabla_\omega f(T_\omega(\varepsilon))$ by auto-diff

Mean-Evidence

$$\mathbb{E}_{\varepsilon} \nabla_\omega \log P(T_\omega(\varepsilon), x) = \nabla_\omega \int_z q_\omega(z) \log P(z, x) dz$$

Given $x = (x^1 \dots x^n)$

$$P(z, x) = P_z(z) \prod_{i=1}^n P(x^i | z)$$

then,

$$\begin{aligned} \nabla_\omega \log P(T_\omega(\varepsilon), x) &= \nabla_\omega \log P_z(T_\omega(\varepsilon)) \\ &\quad + \sum_{i=1}^n \nabla_\omega \log P_{xz}(x^i | T_\omega(\varepsilon)) \end{aligned}$$

expensive for large n
(log calc.)

Reducing expense via sub-sampling

$$\log P(z, x) \approx \log P_z(z) + \frac{1}{m} \sum_{i \in \text{mini batch}} \log P_{xz}(x^i | z)$$

Doubly-stochastic estimator

$$\nabla_\omega \log P_z(T_\omega(\varepsilon)) + \frac{1}{m} \sum_{i \in \text{mini batch}} \log P_{xz}(x^i | T_\omega(\varepsilon))$$

- stochastic over ε & mini-batch

$$q_\omega(z) = \mathcal{N}(z | \omega, 0.5 \cdot I)$$

Variational distribution

$$\int q_\omega(z) \log P(z, \text{Data}) dz$$

Mean Evidence

LECTURE 24
V.I & REVIEW

Re-parameterization discussion

Pros - low-ish variance

Cons - Only works for continuous z
Need to find $q(\varepsilon)$ & $T_\omega(\varepsilon)$

Usually use auto-diff

Score function Estimator

We know,

$$\nabla_\omega \int_z q_\omega(z) \log P(z, x) dz$$

Expectation

$$= \int q_\omega(z) \log P(z, x) \nabla_\omega \log q_\omega(z) dz$$

$$\therefore \nabla_\omega \log q_\omega(z) = \frac{1}{q_\omega(z)} \nabla_\omega q_\omega(z)$$

Estimator

sample $z \sim q_\omega(z)$

Set $g = \log P(z, x) \nabla_\omega \log q_\omega(z)$

$$\text{Then} - \mathbb{E} g = \nabla_\omega \mathbb{E}_{q_\omega(z)} \log P(z, x)$$

What if we don't have closed-form entropy?

$$\begin{aligned} -\nabla_\omega \int_z q_\omega(z) \log q_\omega(z) &= -\int (\nabla_\omega q_\omega(z)) \log q_\omega(z) \\ &\quad \text{chain rule} \\ &= -\int q_\omega(z) (\nabla_\omega \log q_\omega(z)) \\ &\quad \text{Can estimate} \\ &\quad \text{like mean Evidence} \\ &\quad \downarrow \\ &\quad \nabla \log x = \perp \nabla x \\ &= \nabla \frac{q_\omega}{q_\omega} q_\omega = \nabla 1 \\ &= 0 \end{aligned}$$

$$\nabla_\omega \text{ELBO} = \left[\nabla_\omega \mathbb{E}_{q_\omega(z)} (\log P(z, x) - \log q_\omega(z)) \right]_{v=\omega}$$

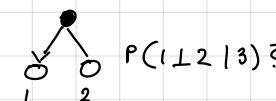
DISCUSSION OF Score function

Pros - simple, easy, works with discrete z

Cons - Typically high variance

REVIEW

① Conditional Ind. in directed model



$$P(x_1 \perp x_{\text{not}(1)} | x_{\text{pa}(1)})$$

② Cond. ind. in Undirected model

③ HC Theorem ↗ Something involving proof

Change assumption or ... work?

$$P(x) > 0$$

$$G_{ij} > 0$$

④ Message Passing efficiency (Time complexity)

⑤ Msg Pass inference on tiny model

⑥ Chow Liu In a tiny dataset what is ML tree

⑦ Exponential Family ↗ cheat sheet

given a simple dist & dataset
is θ an optimal of likelihood

⑧ Bayesian Inf. — Conceptual

8-10ish ↗

3-5 parts

Is decision rule optimal?
Compare to ML, MAP

⑨ MCMC — Conceptual ↗ Regularity

Detailed balance

Given Metropolis / Gibbs change some aspect of algo → check if work
Give argument?

⑩ VI — Conceptual → ELBO

⑪ SGVI — Given simple p, q derive gradient estimator

init ω

for $t_1 \dots t_{\max}$

$$\nabla_{\omega} = \nabla_{\omega} \text{ELBO}$$

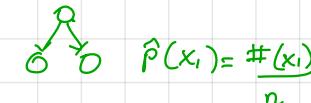
$$\omega \leftarrow \omega + \alpha \nabla \omega$$

$$\sum_z q_{\omega}(z) \log \frac{p(z, x)}{q_{\omega}(z)}$$

$p(z) \prod p(x^i | z)$

$$\int_{\mathcal{E}} q(\varepsilon) \log p(\varepsilon, \omega) \prod p(x^i | \varepsilon)$$

$$P(x_i \perp x_{\text{not}(i)} | x_{\text{nb}(i)})$$



$$\hat{P}(x_i) = \frac{\#(x_i)}{n}$$

Tree-structured MRFs \equiv DAG \Rightarrow Counting reasonable

$$p(x) = \prod_i \phi_i(x_i) \prod_j \phi_j(x_i, x_j)$$

$\hat{P}(x_i)$ $\hat{P}(x_i, x_j)$

$$\log \sum_i \hat{P}(x_i)$$

Ind. Free

$$\log \sum_i x_i x_j$$

mutual Inf.

