# IERG 5130 Probabilistic Models and Inference Algorithms for Machine Larning
## Term 1, 2018-19

Instructor: Prof. Dahua Lin
Scribe: Zihao FU

Distribution.

- Random Variable,

$$x \in \{0, 1\}$$

Variable
↓
Value

两个 随机变量.

$$\underset{\{0,1\}}{X} \qquad \underset{\{0,1,2\}}{Y}$$

| X\Y | 0 | 1 | 2 |
|-----|-----|-----|-----|
| 0 | 0.2 | 0.2 | 0 |
| 1 | 0 | 0.1 | 0.5 |

Joint  Distribution   $P(X, Y)$

若已知 $x = 0$

$$P(Y \mid x=0) = \frac{P(0, Y)}{\sum_{Y=0}^{2} (0, Y)}$$

↑

conditional  distribution

如果有 $n$ 个变量，每个有 $2$ 个值，
一共有 $2^n$ 个搭配，指数级增长。

stats vs Machine Learning

efficient

complexity          Direct calcute

两个模型: Bayes net work & Markov

Independence between Variables.

Model 和 Distribution 的区别

Normal Distribution

$$P(x \quad ) = \frac{1}{\sqrt{2\pi\sigma^2}} \exp\left(-\frac{(x-\mu)^2}{2\sigma^2}\right)$$

$x$: example
$\mu, \sigma$: parameters

改变, $\mu$ 和 $\sigma$ 得到 parametric family

Model 指的是这整套 family.

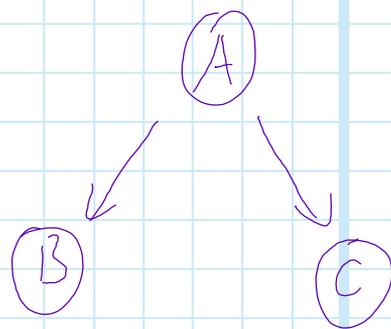A Graphical model Λ is use graphi to represent the family

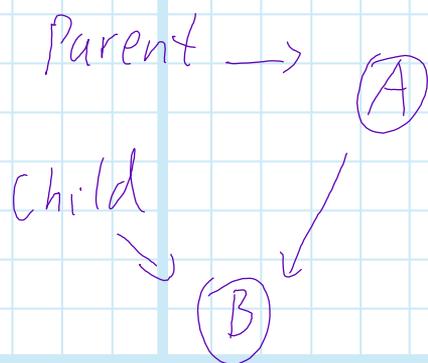How can we use graph to represent the family ?

Graph

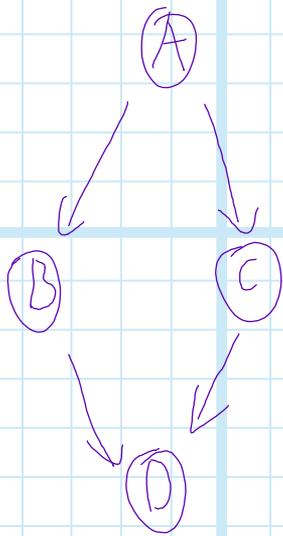$$p(A, B, C) = \underbrace{P(A) P(B|A) p(C|A)}_{factorization}$$

not always true here, $B \perp C | A$.



This is a Directed Acydic Graph (DAG)

Parent $\longrightarrow$ (A)

Child

(B)

也叫做 Bagesian Network.

A C D B
A B C D

Topological ordering:
  parents always before children
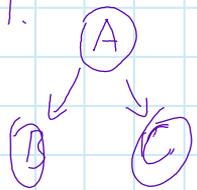
Inference always follow Topolgical order

如何从 BN来写 factorization

$$G = (V, E)$$

有几个变量就有几个乘项.
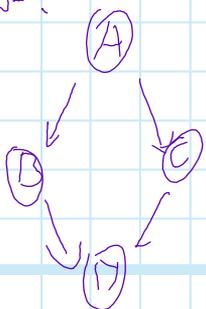
$$P(X_V) = \prod_{s \in V} P(X_s | X_{\pi(s)})$$

例1.

factorization formula

$$P(X_A | X_{\pi(A)}) P(X_B | X_{\pi(B)}) P(X_C | X_{\pi(C)})$$
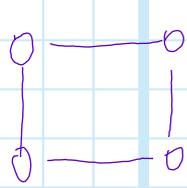
$$= P(X_A) P(X_B | X_A) P(X_C | X_A)$$

例2.

$$P(X_A | X_{\pi(A)}) P(X_B | X_{\pi(B)}) P(X_C | X_{\pi(C)}) P(X_D | X_{\pi(D)})$$

$$= P(X_A) P(X_B | X_A) P(X_C | X_A) P(X_D | X_B, X_C)$$
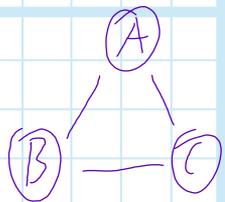
有时候并不清楚谁是 cause，此时就退无向图。



　　　　如何表示这种关系的 likelihood.
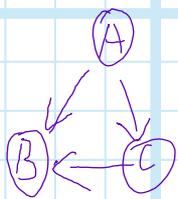
# Markov Random Field (MRF)
## or Markov Network.
　　用 undirected Graph 表示。

　　如何表示？

法1：加上依赖关系：

　　$p(A|B) \, p(A|C) \, p(B|C)$

　　　　　　但是强加了依赖关系。

法2、
$$P(A,B,C) = \frac{1}{z} \phi(A,B) \, \phi(B,C) \, \phi(A,C)$$

能否写成一个通式？

clique: A fully connected subset
　　　　　　　　$\{a, b\}$　$\{a, b, c\}$



　　　　不是团：$\{a, b, c, d\}$ × ∵ a,d 未连接。

Maximal Clique: add one more variable is not clique, it's a maximal clique.

$\{a,b\}$ 不是最大团, ∵ 添加 $c$ 了之后区是团。

$\{b, c, d, e\}$ 是最大团.

$\{b, c, d\}$ 是团, 但不是最大团.

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

↖ compatibility function.

$$\frac{1}{Z} \phi(a,b,c) \phi(b,c,d,e)$$

$$\phi(a,b,c) = \phi(a,b) \phi(a,c) \phi(b,c)$$

如何定义 compatibility function,
can be arbitry non-negative function

例: Ⓐ — Ⓑ

| A\B | 0 | 1 |
|-----|---|---|
| 0   | 2 | 1 |
| 1   | 1 | 2 |

# Independence : Most fundamental in probabilistic theory.

can be used to simplify the computation.

Independence.

$$P(A, B) = P(A) P(B) \quad A \perp B$$

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

下面来看期望.

$$E[f(x)] = \sum_{x \in X} P(x) f(x)$$

$$A \perp B \Rightarrow E[f(A) \cdot f(B)] = E[f(A)] E[f(B)]$$

用处: 减小计算复杂度

$$E[f(A) \cdot f(B)] : O(m^2)$$

$$E[f(A)] \cdot E[f(B)] : O(m).$$

Proof:

$$E(f(A) f(B)] = \sum_{A} \sum_{B} P(A, B) f(A) f(B)$$

$$= \sum_{A} \sum_{B} P(A) P(B) f(A) f(B)$$

$$= \sum_A P(A) f(A) \sum_B P(B) f(B)$$

$$= E[f(A)] E[f(B)]$$

推广: $E\left[\prod_{i=1}^{m} f_i(x_i)\right] = \prod_{i=1}^{m} E(f_i(x_i))$

if $x_i$ independent from each other.
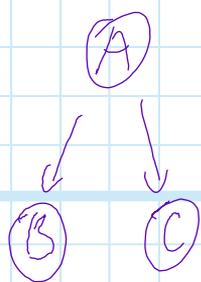
Conditional Independence.

$$P(A, B \mid C) = P(A \mid C) P(B \mid C)$$

↑

Conditional Independence.

It's very hard to verify.

But if we have graph, It's easy to see.



$\Rightarrow P(B, C \mid A) = P(B \mid A) P(C \mid A)$

$$B \perp C \mid A$$

a statement of conditional independent.

可见 Graphic Model 的作用:

1. give the factorization form.

2. encode conditional independent

# I-map.

model $\longrightarrow$ a set of conditional independent.

$$P \longrightarrow I(P)$$

family distribution

$$G \longrightarrow I(G)$$

I-map: $I(G) \subseteq I(P)$

We call $G$ an I-map of $P$

$G$ contains subset of the conditional independence

of $P$.



$$P(X_v) = \frac{1}{Z} \prod_S \phi_S(X_S)$$

$$\Downarrow$$

$I(G) \quad \underset{\text{what's the}}{\overset{?}{\underline{\phantom{xxx}}}} \quad I(P)$$

relation.

① Local Independence.

$$X_S \perp X_{V \setminus (S \cup N_G(S))} \mid X_{N_G(S)}$$

$X_s$ is conditinal indepent of the
rest of the world given its neighbour.

$$G \qquad P$$

① $\Downarrow$

$$I(G) \quad \overset{②}{\approx} \quad I(P)$$
$$I_g(G)$$

## 1. Local Independence

$$I_l(G) = \{ X_S \perp X_{V \setminus (S \cup N(S))} \mid X_{N(S)} \}$$



②: $N(S)$

## 2. Pairwise Independence.    是 Local Indep-nde·的拓展

$$I_p(G) = \{ X_A \perp X_B \mid X_{V \setminus (A \cup B)} : A \text{ and } B \text{ disjoint } \& \text{ no direct edge between } A \text{ and } B \}$$



← 这个不行, ∵有直接的边.

有 $I_l(G) \subseteq I_p(G)$

## 3. Global Independence.

$$I_g(G) = \{ X_A \perp X_B \mid X_C : C \text{ seperate } A \text{ and } B \}$$

Seperate: 从A到B无论怎么走都会经过C.

$$I_l(G) \subseteq I_p(G) \subset I_g(G)$$

← 能从图中得到的所有 independence.

I-map:

G is an I-map of P if $I(G) \subseteq I(P)$

① Soundness

P factorize according to G

$G = (V, E)$ $\quad P(x) = \frac{1}{Z} \prod_{c \in C(CG)} \phi_c(x_c)$

G is an I-map of P, i.e. $I(G) \subseteq I(P)$

证明:

把图分为三部分

$Ⓐ — Ⓒ — Ⓑ$ $\qquad a,b,c$: specific assignment of A,B,C

$p(a, b, c) = \frac{1}{Z} \psi_A(a) \psi_B(b) \psi_C(c) \phi_{AC}(a,c) \phi_{BC}(b,c)$

$\phi'_{AC}(a,c) = \psi_A(a) \psi_C(c) \phi_{AC}(a,c)$

$\phi'_{BC}(b,c) = \psi_B(b) \phi_{BC}(b,c)$

$\Rightarrow = \frac{1}{Z} \phi'_{AC}(a,c) \phi'_{BC}(b,c)$

问题转化为证明: $A \perp B \mid C$ $\qquad p(a,b,c) = \frac{1}{Z} \phi_{AC}(a,c) \phi_{BC}(b,c)$

$q_{A|C}(a,c) = \dfrac{\phi_{AC}(a,c)}{\sum\limits_{a'} \phi_{AC}(a',c)}$ $\qquad \psi_1(c) = \sum\limits_{a'} \phi_{AC}(a',c)$

$\phi_{AC}(a,c) = q_{A|C}(a|c) \psi_1(c)$

$\phi_{BC}(b,c) = q_{B|C}(b|c) \psi_2(c)$

$\therefore p(a,b,c) = \frac{1}{Z} \phi_{AC}(a,c) \phi_{BC}(b,c) = \frac{1}{Z} \psi_1(c) \psi_2(c) q_{A|C}(a|c) q_{B|C}(b|c)$

上面证明了如果一个概率分布能分解为一个图,那

幺图生的 Independence 项 都能在 概率分布中扇色.

下面来看逆命题

Hammergley-Clifford

- P: positive distribution. ← density value is positive almost everywhere

- $G = (U, E)$ is an I-map of P.

$$I(G) \subseteq I(P)$$

$\Rightarrow P$ factorize according to $G$.

我们能否用相同的方法处理 Bayes Network? 可以.

无向: Ⓐ — Ⓑ — Ⓒ

有向:

① Ⓐ → Ⓑ → Ⓒ

Ⓐ ← Ⓑ ← Ⓒ

Ⓐ → Ⓑ ← Ⓒ

Ⓐ ← Ⓑ → Ⓒ

$A \perp C | B$ 在哪些里面成立.

下面研究 BNs 的 conditional independence.

Graph       Formulation

①   Ⓐ → Ⓒ → Ⓑ      $P(a)P(c|a)P(b|c)$    ✓

②   Ⓐ ← Ⓒ ← Ⓑ      $P(b)P(c|b)P(a|c)$    ✓

③   Ⓐ ← Ⓒ → Ⓑ      $P(c)P(a|c)P(b|c)$    ✓

④   Ⓐ → Ⓒ ← Ⓑ      $P(a)P(b)P(c|a,b)$    ✗

$A \perp B | C$ 是否成立?

即 $P(a,b|c) = P(a|c)P(b|c)$

$P(a)P(c|a)P(b|c) = P(a,c)P(b|c)$
$= P(c)P(a|c)P(b|c)$ ← when c is given
     a, and b are independent from each other

∴ 成立

同理成立

$P(c)P(a|c)P(b|c) = P(c)P(a,b|c)$
成立

不成立. 反例:

cause: A: work hard   B: sleep well

effect: C: pass exam

| A | B | P(c|A,B) |
|---|---|---|
| 0 | 0 | 0.3 |
| 0 | 1 | 0.8 |
| 1 | 0 | 0.8 |
| 1 | 1 | 0.9 |

Suppose c=1 observed

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

| A\B | 0 | 1 |
|---|---|---|
| 0 | $\frac{3}{20}$ | $\frac{8}{20}$ |
| 1 | $\frac{8}{20}$ | $\frac{9}{20}$ |

共 28 / 20   归一 28

$P(A=1|C=1) = \frac{17}{20}$

$P(B|C=1) = \frac{17}{20}$

$P(A=1, B=1|C=1) = \frac{9}{20}$

∴ $P(A=1|C=1)P(B=1|C=1) > P(A=1, B=1|C=1)$

Trick:

     写 formulation 的时候, 一般按拓扑排序来写.

Explain Away:



     一个结果有很多原因, 如果一个原因的概率高,
     那么其它的原因概率则低. 则不是条件独立了.

Active Trial



A    C    B

无向图中, observe 的变量 block 住该 path, block 住了就是 inactive, 如果所有 path inactive 则条件独立. 有向图中 block 和 active 的关系有四个特殊见课件 29 页.

$A \perp B | C$

A and B are blocked by C since all path from A to B is blocked/inactive by C

下面看有向图如何描述 block.

     ① ② ③   C block 而 ④   active when C is observed
                                 inactive when C is not observed.

     用 d-separation 来描述.

# D-seperation

observed.

$$Ⓐ \rightarrow Ⓑ \rightarrow Ⓒ \leftarrow Ⓓ \rightarrow Ⓔ$$

$A \perp D \mid C$ ? 找到 A 到 D 的所有 path. 如果 path 很长, 就 sub sesceion by sub session 地分析

而 A subsession
$$\begin{cases} ① A \rightarrow B \rightarrow Ⓒ & active \\ ② B \rightarrow Ⓒ \leftarrow D & active \end{cases}$$

∴ $B - B - C - D$ is active

∴ $A \perp D \mid C$ 不成立.

observe 的 变量子还是 blocked, 不能 通过. 而 active 则 根据 block 的情况由上面四种来看.

所有 BNs 都能写成 MRF

$$Ⓐ \qquad Ⓑ$$
$$\searrow \quad \swarrow$$
$$Ⓒ$$
$$\downarrow$$
$$Ⓓ$$

$$p(a) p(b) p(c \mid a, b) p(d \mid c)$$

⇓ 可以直接把概率定义为 compatibility function

$$\frac{1}{z} \phi_a(a) \phi_b(b) \phi_{c,a,b}(c, a, b) \phi_{d,c}(d, c)$$

⇓ 画成 MRF

$$Ⓐ \rightarrow Ⓑ$$
$$\searrow \quad \swarrow$$
$$Ⓒ$$
$$\mid$$
$$Ⓓ$$

下面研究如何从有向图直接画无向图.

# Moralization

结是 effect, 把所有的 parent 连起来的过程.

有向 G ⟹ 无向 M

问题: $I(G) = I(M)$ ?

答案: 不相等.

实质上 $I(M) \subseteq I(G)$ 例如图中 $A \perp B \in I(G)$
$A \perp B \notin I(M)$

转成无向图, Conditional Independence 信息会丢失.

为了不丢失信息, 需要研究 factor graph.

factor graph 是二部图.

Factor Graph:



也可以画成:



解决问题流程.

1. Formulate the model

2. estimate the model based on data

3. apply model

下面讲如何 formulate a graphical model.

如何建模.

理解问题, 用如下描述:

- Variables

- Relations

- Constraints & Assumptions.

Gaussian Mixture Model.

sample 都有 vector space representation



自然生成模型和 GAN
那种生成模型是两回事

cluster     prior           tradeoff: complexity
$\downarrow$         $\downarrow$                      expressivity
components   components

$\downarrow$

points

GMM

Assumptions:

1. Each component is a Gaussian Distribution $(\mu, \Sigma)$

$$P(x|\mu, \Sigma) = \frac{1}{\sqrt{(2\pi)^d |\Sigma|}} exp\left(-\frac{1}{2}(x-\mu)^T \Sigma^{-1}(x-\mu)\right)$$

2、 $k$: index of the component $\quad k = 1:K$

$$\{(\mu_k, \Sigma_k)\} \leftarrow \text{parameters.}$$

Prior for component choice. $\pi \quad \pi = (\pi_1, \cdots, \pi_k)$

Given the Model, the generate Process:

1. Choose a component:
$$z_i \in \{1, \cdots, K\} \sim \pi$$

2、 $x_i \sim N(\mu_k, \Sigma_k)$ where $k = z_i$

或者写作 $x_i \sim N(\mu_{z_i}, \Sigma_{z_i})$

$X = (x_1, \cdots x_n) \quad Z = (z_1, \cdots, z_n)$

$$P(X, Z \mid \theta) = \prod_{i=1}^{n} P(x_i, z_i \mid \theta) = \prod_{i=1}^{n} P(z_i \mid \pi) P(x_i \mid \mu_{z_i}, \Sigma_{z_i})$$

$(\theta, \mu_{z_i}, \Sigma_{z_i})$ 是 parameters. parameter 是从数据里

学得, 学完了就固定的.

看如何画图



不能表示 各样本



Factor Graph

Question:

$G_1$



$G_2$

# GMM

$\pi$

$\pi, \mu_k, \Sigma_k$ 是 parameters

$z_i$

$\mu_k$
$\Sigma_k$
$x_i$

$G_1$

$G_2$

如何控制 $G_1$, $G_2$ 这种不同类型.

## Group_wise GMM

① 

m: index of group

②

$z_i$

$\mu_k$
$\Sigma_k$
$x_i$
$N_m$

$M$

$z_i$

$\mu_k$
$\Sigma_k$
$x_i$
$N_m$

$M$

①: All component 共用一套参数

②: 每个 component 有自己的参数

选① 还是② 靠观察 dataset 来决定.

parameter:

$$\{\pi_1, \pi_2, \cdots, \pi_M\}$$

$$\{(\mu_1, \Sigma_1), \cdots, (\mu_k, \Sigma_{(c)})\}$$



问题: train: $G_1, \cdots, G_M$

如果有一个新 group $G'$ 我们没有参数 $\pi'$

解决方案: 把参数 $\{\pi_1, \pi_2, \cdots \pi_M\}$

变为 random variables.

那么问题变成了如何生成 r.v. $\pi_m$?

$\pi_m$ is a probability vector

i.e. $\sum_k \pi_m(k) = 1$

$\pi_m(k) \geq 0 \quad k = 1, \cdots, k$.

Dirichlet Distribution.

Group-wise GMM v2

设计原则:

parameter 不要放在框里, 因为这样新的
group 来了没法弄。 解法是把 parameter
变成 variable, 再引入先验参数

下面看 temprol 的建模方式:
Temporal Dynamic GMM



① $\pi_t$

② $z_i$

③ $x_i$

dynamics: how things change
over time.

Dynamic 可以在 $x_i$, $z_i$, $\pi_t$ 三个 level, 下面分别对应

① Dynamic on $x_i$

dynamic model: 已知当前点, 下一个点的出现方式。

如何描述点的变化?

物理中 Motion:

$$\frac{d\overset{location}{x_t}}{dt} = \mu_t$$

velocity

$$dx_t = \mu_t \cdot dt$$

$$dx_t = \mu(x_t, t)dt$$ 速度是位置和时间的函数

$\hookrightarrow$ DE: differential Equation (deterministic)

用微分方程经过初始状态，状态是确定

计算。下面引入 uncertainty

把 DE 变成 stochastic Differential Equation

$$dx_t = \mu(x_t, t) \cdot dt + \sigma(x_t, t) \underline{dB_t}$$

$$\downarrow$$

Brownian
Motion

$$dx_t = \sigma \cdot dB_t$$

叫做 perfect Browian Motion

$$dx_t = x_{t+dt} - x_t$$

虽然 $x_t$ 无法明确的计算，但是有分布：

$$dx_t \sim N(0, \sigma \cdot dt)$$

越往后方差越大



$$dx_t = \mu_t \cdot dt$$

那么怎么写成 conditional distribution 的形式？

Discretize：只观注 整点的 time step.

$$x_t, \quad x_{t+1} \ldots$$

$p(x_{t+1}|x_t)$ 为了建模，做如下简化：

—Stationary： 去掉和时间相关

$$dx_t = \mu(x_t)dt + \sigma \cdot dB_t$$

- Linear:

$$\mu(x_t) = Ax_t + b$$

determinstic part! 决定运动

布朗运动
决定运动

$$\therefore dx_t = x_{t+1} - x_t \sim N(\mu(x_t), \sigma^2 I)$$

$$\sim N(Ax_t + b, \sigma^2 I)$$

$$\therefore p(x_{t+1} | x_t) \sim N(x_t + Ax_t + b, \sigma^2 I)$$

$$= N(A' x_t + b, \sigma^2 I)$$

这样完成了布朗运动的建模(条件分布)

卡尔曼滤波也是用它建模.

回到之前的 GMM, 可以应用.



A, b, ∧ 就是布朗运动里的那个参数.

详细说明李伪布朗运动

$$dx_t = \mu(x_t) \cdot dt + \sigma \cdot dB_t$$

$$x_{t+1} - x_t$$

$$= \int_t^{t+1} \mu(x_t)dt + \int_t^{t+1} \sigma \cdot dB_t$$

黎曼积分     伊藤积分

确定的积分     积出来是 r.v.

$$\int_t^{t+1} dB_t \sim N(0, I)$$

$$\mu(x_{t+\delta t}) = \mu(x_t) = \mu(x_t) \int_t^{t+1} dt$$

$$\delta t \in (0, 1) \qquad = \mu(x_t) + \sigma \cdot y$$

之前的 dynamic 是在 $x_i$ 上，上面有
dynamic 在 $z_i$ 上，时离散的 HMM

② Dynamics on $z_i$

$$\boxed{z_0} \longrightarrow \boxed{z_1} \longrightarrow \boxed{z_2} \to \cdots \longrightarrow \boxed{z_T}$$

$$\downarrow \qquad\qquad \downarrow \qquad\qquad \downarrow \qquad\qquad \downarrow$$

$$\boxed{x_0} \qquad \boxed{x_1} \qquad \boxed{x_2} \qquad\qquad \boxed{x_T}$$

这里 dynamic 不直接在 $x_i$ 上，而是在
$z_i$ 上。

$$p(z_{t+1} | z_t)$$

step:
$$z_0 \sim \pi$$
$$x_0 \sim N(\mu_{z_0}, \Sigma_{z_0})$$
$$z_{t+1} | z_t \sim$$

$$
\begin{array}{c|ccc}
z_t \backslash z_{t+1} & 1 & \cdots & K \\
\hline
\vdots & & & \\
\vdots & & \text{conditional} & \\
\vdots & & \text{table.} & \\
K & & & \\
\end{array}
$$

$$x_t \sim N(\mu_{z_t}, \Sigma_{z_t})$$

$\therefore$ HMM 是 GMM dynamic 在 $z_i$ 上的扩展

下面看能否把 dynamic 再上升一个 level

③ Dynamics on $\pi_t$.

$$\boxed{\pi_0} \longrightarrow \boxed{\pi_1} \longrightarrow \cdots \longrightarrow \boxed{\pi_T}$$

我们希望 $\pi_{t+1}$ 与 $\pi_t$ 接近

可以用高斯模型建模，但是无法
保证新的 $\pi_{t+1}$ 是概率分布

即 $P(\pi_{t+1}|\pi_t) \sim N(\pi_t, \sigma)$

不保证 $\sum_k \pi_{t+1}(k) = 1$

可以用 Transformed (Warped) Model.



$$\pi(i) = \frac{e^{n(i)}}{\sum_i e^{n(i)}}$$ softmax transform.

Softmax 可以把任意的向量变成
probabilty vector.

这个模型叫 Warped Gaussian
Copula Process

上面讲了 Bayes Network，下面讲如何
对 Markov Random Field 建模

Image as grid



$j \in N_i$ neighbour.

Smoothness Assumption.

在 (可用) compatibility function 表示 Smoothness Assumption?

$$\frac{1}{Z} \prod_{(i,j) \in P_N} \phi_{ij}(x_i, x_j)$$

$\downarrow$

$$\frac{1}{Z} \exp\left( \sum_{(i,j) \in P_N} e_{ij}(x_i, x_j) \right)$$

Gibbs
$\swarrow$ Distribution

← Neighboud 集合

basic principle of energy function:

$e_{ij}$ ⟶ large value (undesirable)
       ⟶ small value (desirable)

$$e_{ij}(x_i, x_j) = \frac{1}{2}(x_i - x_j)^2$$

为什么要选这个? 为了 mathematical convinent.

Limitation : Smoothness Assumption

在图像中的世界处并不成立.

如何克服: 采用 Gated Markov Random Field.

Gated MRF:

$z_{ij}$ : $i$ $j$ 间是否有 boundary.

$z_{ij}$ 是一个 Gate.

$$\frac{\alpha(z_{ij})}{2}(x_i - x_j)^L$$

如何 develop model:

先有个 simple model, 然后观察 Assumption 是否不完全符合, 如果不是所有的都符合, 就引入 cutten variable 来 indicate 每个变量的类型.

Field of Experts.

Textures:



这种用 MRF smoothness 不好建模.

Field of Experts



Kernals:



pattern

$k_1$    $k_2$    $k_3$

4    2    0

Single Expert:

$$exp(k_i^T I)$$

Product of Experts

$$k_1, \cdots, k_m$$

$$exp\left(\sum_{j=1}^{n} \alpha_j k_j^T I\right)$$

但是这种只能用 Kernel 和啮图一样大的
情况，如何打脱？

Field of Experts

$$\frac{1}{Z} exp\left(\sum_c \sum_{j=1}^{m} \alpha_j \cdot (K_j^T \cdot I_c)\right)$$

Long Tail Distribution 问题



$$K_j^T I_c$$

用到波包法把

$$f(u) = exp(u) \Rightarrow f(u) = \frac{1}{1 + \frac{1}{2}u^2}$$

$$exp\left(log \frac{1}{1 + \frac{u^2}{2}}\right)$$

Exponential Family

$$P_\theta(x) = \frac{h(x)}{Z(\theta)} exp\left(\eta^T(\theta) \phi(x)\right)$$

normalizing
constant.
(partition function)

function of
parameter

function of
data.

$$= h(x) \cdot \exp\left(\eta(\theta)^T \phi(x) - \log Z(\theta)\right)$$

$$= h(x) \cdot \exp\left(\eta(\theta)^T \phi(x) - \underline{A(\theta)}\right) \longrightarrow \text{log partition function}$$

$$Z(\theta) = \int \exp\left(\eta(\theta)^T \phi(x)\right) h(x) \nu(dx)$$

$$A(\theta) = \int \exp\left(\eta(\theta)^T \phi(x)\right) h(x) \underline{\nu(dx)}$$

$\searrow$ base measure.

※ $\eta^T(\theta)\phi(x)$ 可选的例子:

$$f(\theta, x) = \theta \cdot x^2 \quad \checkmark$$

$$f(\theta, x) = \frac{1}{1 + \theta \cdot x} \quad \times \quad \text{不能分解或嵌套的形式}$$

※ base measure: 统一离散和连续的表达形式.

在概率中有两种分布 离散和连续

Discrete Distribution

$$p_1 \cdots p_n$$

$$f: X \to \mathbb{R}$$

$$E_p[f] = \sum_{i=1}^{n} p_i f_i \quad \text{如何写成统一的积分形式}$$

$$= \int f(x) p(x) \underline{\mu(dx)}$$

$\uparrow$ counting measure $\to$ |—|—|—|—|—|—|$\longrightarrow$

连续情况下:

$$E_p[f] = \int f(x) p(x) \, dx$$

$$= \int f(x) p(x) \underline{\nu(dx)} \quad \text{用勒贝格测度}$$

$\nwarrow$ 转化为标准零积分.

※

$$\underset{\substack{\text{base} \\ \text{density}}}{h(x)} \exp\left(\underset{\substack{\text{canonical} \\ \text{parameter}}}{\eta(\theta)^T} \underset{\substack{\text{sufficient} \\ \text{statistics}}}{\phi(x)} - \underset{\text{log partition function}}{A(\theta)}\right)$$

Sufficient statistics: 只要知道期望 $E_p[\phi(x)]$
整个分布就确定了

canonical parameter: 就是说所有的 parameter
划成这种标准形式

为什么 exp-family 重要?

$\therefore$ 很多常见分布都是 exp family

下面看如何写成 exp family.

Bernolli Distribution: (最简单的分布)

$x \in \{0, 1\}$

$p(x) = \begin{cases} P_0 & (x=0) \\ P_1 & (x=1) \end{cases}$

$P_0 + P_1 = 1$

$P(x) = \begin{cases} \exp(\log(P_0)) & (x=0) \\ \exp(\log(P_1)) & (x=1) \end{cases}$

$= \exp\left( \mathbb{1}(x=0) \log(P_0) + \mathbb{1}(x=1) \log(P_1) \right)$

$= \exp\left( (1-x) \log(P_0) + x \cdot \log(P_1) \right)$

$= \exp\left( \begin{bmatrix} 1-x \\ x \end{bmatrix} \cdot \begin{bmatrix} \log P_0 \\ \log P_1 \end{bmatrix} \right)$

# Exponential Family

$$P_\theta(x) \quad \underline{h(x)} \exp\left( \underline{\eta(\theta)^T} \, \underline{\phi(x)} - A(\theta) \right)$$

base measure    canonical parameter    sufficient statistics    → log partition function

下面来看 Poisson Distribution

自然数的分布

$$P_\lambda(x) = \frac{\lambda^x}{x!} e^{-\lambda} \qquad x \in \{0, 1, \cdots\}$$

$$\lambda^x = \exp(x \cdot \log\lambda) \quad 代入原式:$$

$$P_\lambda(x) = \frac{1}{x!} \exp(x \cdot \log\lambda - \lambda)$$

$$\underbrace{\phantom{\frac{1}{x!}}}_{h(x)} \qquad \underbrace{\phantom{x}}_{\phi(x)} \underbrace{\phantom{\log\lambda}}_{\eta(\theta)} \qquad \underbrace{\phantom{\lambda}}_{A(\theta)}$$

# Exponential Distribution

$$P_\lambda(x) = \lambda e^{-\lambda x}$$

used to capture the time. is high related to Poisson Distr.



$t_1 \quad t_2 \quad t_3$

rate: in unit time, how many happens.

the rate is the $\lambda$ in Poisson,

$$P_\lambda(x) = \lambda e^{-\lambda x}$$

$$= \exp(-\lambda x + \log \lambda)$$

$$-\lambda : \eta(\theta) \qquad \text{or} : \lambda : \eta(\theta)$$
$$x : \phi(x) \qquad \qquad -x : \phi(x)$$

# Normal Distribution

$$P_{(\mu, \sigma^2)}(x) = \frac{1}{\sqrt{2\pi\sigma^2}} \exp\left(-\frac{(x-\mu)^2}{2\sigma^2}\right)$$

$$\frac{(x-\mu)^2}{2\sigma^2} = \frac{x^2 - 2\mu x + \mu^2}{2\sigma^2}$$

$$= \frac{1}{2\sigma^2}x^2 - \frac{\mu}{\sigma^2}x + \frac{\mu^2}{2\sigma^2}$$

$$= \begin{bmatrix} \frac{1}{\sigma^2} \\ \frac{\mu}{\sigma^2} \end{bmatrix} \cdot \begin{bmatrix} \frac{x^2}{2} \\ -x \end{bmatrix} + \frac{\mu^2}{2\sigma^2}$$

$\frac{1}{\sigma^2}$ : precision 记作 $a$

$\frac{\mu}{\sigma^2}$ : potential coefficient 记作 $b$

$\therefore P(x) = \exp\left(-\frac{a}{2}x^2 + bx + A(a,b)\right)$

$\therefore$ 只要 exp 里面 是个二次函数，
那么就是正态分布！

Uniform Distribution 不属于 Exp Family 的特例.

canonical parameter: 普通参数经过
canonical transform 变为标准参数
小变换后重新作的形式为:

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

$\Omega$: Domain of parameters.

$$Z(\theta) = \int_X \exp\left(\theta^T \phi(x)\right) h(dx)$$

我们希望 $Z(\theta)$ 是 finite value.

The set of valid parameters:

$$\Omega = \left\{\theta : \int_X \exp(\theta^T \phi(x)) h(dx) < +\infty\right\}$$

sample space $\uparrow$

sufficient statistics $\uparrow$

$\Omega$ 由 sample space 和 sufficient statistics 共同决定.

Regular Family

$\Omega$ is an open subset of $\mathbb{R}^d$

open subset 是个挺好的概念.

开集的好处: 领域都在 $\Omega$,

所以在任一地方都可以求导.

幸运的是, 几乎所有常见分布都

是 Regular Family

基础上, 很多 puper 引入各种条件

往往为了引入某种性质, 比如是

可导.

Idenifiable.

$$P(x) = \begin{cases} P & (x=1) \\ 1-P & (x=0) \end{cases} \qquad (\text{Bernulli Distr.})$$

$$P(x) = \exp\left( (1-x)\underbrace{\log(1-P)}_{a_0} + x\underbrace{\log P}_{a_1} \right)$$

$$= \exp\left( a_0(1-x) + a_1 x - A(a_0, a_1) \right)$$

当 $a_0 = a_1 = 1$ 或 $a_0 = a_1 = 2$ 时
都有 $P = 0.5$.

Identifiable problem is a fundamentable problem. Two sets of parameters are undisdiguishable.

形式上:

$$\theta_1 \neq \theta_2 \Rightarrow P_{\theta_1} \neq P_{\theta_2} : \text{Identifiable}.$$

$$\exists \theta_1 \neq \theta_2 \quad s.t. \quad P_{\theta_1} = P_{\theta_2} : \text{Unidentifiable}.$$

∴ 例中 椰个 Bernulli Distr. is unidentifiable

下面研究如何把 unidentifiable 转换为 identifiable.

$$P(x) \cdot \exp(a_0(1-x) + a_1 x - A(a_0, a_1)) \qquad \text{B1}$$

$$\downarrow$$

$$P(x) = \exp(\quad \theta x - A(\theta)) \qquad \text{B2}$$

下面的问题变为一个 Family 是否 identifiable 取决于如何设充分统计量.

Overcomplete representation

$$\exists a \neq 0 \quad a^T \phi(x) = b$$

Minimal representation

  otherwise

对于 Ⓑ :

$$\phi(x) = \begin{pmatrix} 1-x \\ x \end{pmatrix} \quad a = \begin{pmatrix} 1 \\ 1 \end{pmatrix}$$

$$a^T \phi(x) = (1-x) + x = 1 \quad \text{over complete}$$

$$(b=1)$$

对于 Ⓑ2

$$\phi(x) = x \quad \therefore \text{Minimal}.$$

那么 over complete 和 Identifvable 的关系?


Regular Family

  $\exists \ a \neq 0 : \quad a^T \phi(x) = b \quad a.e.$

  $\rightarrow$ unidentifiable.

  注意 有两个条件
  ① Regular Family
  ② $a^T \phi(x) = b$.

记住! 
  $\theta : \quad P_\theta(x) = h(x) \exp\left( \theta^T \phi(x) - A(\theta) \right)$

由此决定

This component determin
the distribution

$$\theta' = \theta + \lambda a \quad \text{where} \quad \lambda \neq 0 \quad \theta' \neq \theta.$$

$$P_{\theta'}(x) = h(x) \cdot \exp(\theta'^T \phi(x) - A(\theta'))$$

$$= h(x) \cdot \exp(\theta^T \phi(x) + \lambda a^T \phi(x) - A(\theta'))$$

$$= h(x) \exp(\theta^T \phi(x) + \lambda b - A(\theta'))$$

$$\Rightarrow P_\theta(x) = P_{\theta'}(x) \qquad Q.E.D$$

(这里用) Regular 性质, θ 往任一方向都收敛界.

$\{\theta + \lambda a, \forall \lambda \in R\}$ 都能得到相同分布.

那么 Unidentifiable 能否推出 Over complete?

证 unidentifiable 比 证 identifiable 容易,

∵ 只要举反例即可. 这个证明以后再证

# Categorical Distribution

$$x = \{1, \cdots, k\}$$

$$P(x) = \begin{cases} \ell_1 & (x=1) \\ \vdots & \\ \ell_k & (x=k) \end{cases} \qquad \sum_{k=1}^{k} P_k = 1$$

和 Bernulli 一样:

$$P(x) = \exp\left(\sum_{k=1}^{K} \theta_k \cdot \mathbb{I}(x=k) - A(\theta)\right)$$

这不是 Minimal Distr.

例如

$$\theta_1 = \theta_2 = \cdots \theta_K = 1$$

$$\theta_1 = \theta_2 = \cdots \theta_K = 2$$

是相同分布.

那么如何写成 Minimal Representation 体现.

下面看充分统计量为什么是充分, 以及如何决定整个分布.

$P_\theta$   $\phi(x)$     充分统计量的均值:

$$\mu = E_P[\phi(x)] = \int P_\theta(x)\,\phi(x)\,\nu(dx)$$

如果知道了 $\mu$, 那么整个分布即确定.

We call it mean parameter.

$\mu \in \mathbb{R}^d$

$$P_\theta \xrightarrow{\text{realize}} \mu$$

$M_\phi$ realizable mean

$\mu$ 可以由一套参数来 realize.

$$= \{\mu : \exists \theta \in \Omega, \text{ s.t. } E_{P_\theta}[\phi(x)] = \mu\}$$

世合 $\langle \phi: \{\mu: \exists p \in P(x) \ E_p[\phi(x)] = \mu\}$

$M_\phi$ is convex set

Convex:
is a subset of real vector space



Convex          non-Convex

$S$ is a convex set when
$\forall x_1, x_2 \in S$

$(1-\lambda)x_1 + \lambda x_2 \in S \quad \forall \lambda \in [0,1]$

任取两点，则整条线段都在S中.

可以当作 $x_1$ 两个向量的加权项.

$\forall x_1, \ldots, x_n \in S$
$P_1, \ldots, P_n$ s.t. $\sum_{i} P_i = 1 \quad P_i \geq 0$

$\sum_{i=1}^{n} P_i x_i \in S$

下面来证 $M\phi$ is a convex set

$\mu_1 \in M\phi$, $\mu_2 \in M\phi$

证: $(1-\lambda)\mu_1 + \lambda\mu_2 \in M\phi$.

证明:

$$P_1 \longrightarrow \mu_1$$
$$P_2 \longrightarrow \mu_2.$$

$$P' = (1-\lambda)P_1 + \lambda P_2 \quad \text{体晨分布}.$$

$$P' \longrightarrow (1-\lambda)\mu_1 + \lambda\mu_2$$

可以用 $\mu$ 的定义证?

$$C = \{x_1, \cdots, x_6\}$$



$conv(C)$: convex hull

包含了所有的 convex combination of $C$

$$conv(C)$$
$$= \left\{ \sum_{i=1}^{n} p_i x_i \;\middle|\; \sum p_i = 1, \; p_i \geq 0 \right\}$$

if $conv(C)$ is finite it's called a convex polytope.

$\mathcal{X}$: finite space $\{x_1, \cdots, x_{12}\}$

$$P_\theta(x) = \exp\left(\theta^\top \phi(x) - A(\theta)\right)$$

如何 charestic realise mean?

$$M_\phi = conv\left(\phi(x_1), \cdots, \phi(x_n)\right)$$

即由有限个充分统计量的凸组合张子.

$$\sum P_i \phi(x_i)$$

Log Partition Funtion.

很多论文都关注e如何计算 log Partition function

$$P_\theta(x) = h(x)\exp\left(\theta^\top \phi(x) - A(\theta)\right)$$

$$A(\theta) = \log \int_X \exp\left(\theta^\top \phi(x)\right) h(dx)$$

$$\boxed{\nabla_\theta A(\theta) = E_{P_\theta}\left[\phi(x)\right]}$$

canonical
Parameter

mean parameter

下面证明这个重要的等式:

$$\nabla_\theta A(\theta) = \nabla_\theta \log \int_X \exp(\theta^T \phi(x)) h(dx)$$

$$= \frac{1}{\int_X \exp(\theta^T \phi(x)) h(dx)} \nabla_\theta \int_X \exp(\theta^T \phi(x)) h(dx)$$

$$h(dx) = h(x) dx.$$

$$= \frac{1}{Z_\theta} \int_X \nabla_\theta \exp(\theta^T \phi(x)) h(dx)$$

$$= \frac{1}{Z_\theta} \int_X \exp(\theta^T \phi(x)) \nabla_\theta [\theta^T \phi(x)] h(dx)$$

$$= \frac{1}{Z_\theta} \int_X \exp(\theta^T \phi(x)) \phi(x) h(dx)$$

$$= \int_X \frac{\exp(\theta^T \phi(x))}{Z_\theta} \phi(x) h(dx)$$

$$\underbrace{\frac{\exp(\theta^T \phi(x))}{Z(\theta)} h(x)}_{p(x).}$$

$$= \int_X \phi(x) p(x) \nu(dx) = E_p[\phi(x)]$$

$\nabla_\theta A$ : Gradient Map

   Map the canonical parameter

     to mean parameter.

下間 实际 $\nabla_\theta A$, Identifable, 和 Minimal_Rep.

Oct 19. In-class discuss & Presentation.

Exponential Family

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

$$A(\theta) = \int_X \exp(\theta^T \phi(x)) \, h(dx). \quad \boxed{\nabla A}$$

Canonical parameter $\theta$.    Mean parameter $\mu = E_\theta[\phi(x)]$

$M_\phi$. realizable mean (convex) $\Omega$ dataset

---

Gradient Map.

$$\boxed{\nabla_\theta A(\theta) = E_\theta[\phi(x)]}$$

injective?   one-to-one.

$$x_1 \neq x_2 \Rightarrow f(x_1) \neq f(x_2).$$

surjective?

$$\nabla_\theta A(\Omega) \overset{=}{=} M_\phi ?$$

① $\nabla_\theta A$ is injective  iff minimal representation

Proof:
— minimal $\Rightarrow \nabla_\theta A$ injective $\Leftarrow Q$

---

② $\boxed{H = \nabla^2 f \geqslant 0 \quad \text{Semi positive definite}}$

Symmetric matrix $A$. eigenvalue & eigenvector
$\lambda$ $\qquad$ $e$

$$Ae = \lambda e$$

$\lambda_1 \geqslant \cdots \geqslant \lambda_n \geqslant 0 \Rightarrow A$ is ~~posi~~
$\qquad\qquad$ positive Semidefinite

$\downarrow$

$> 0 \Rightarrow A$ is positive
$\qquad$ definite.

---

Hessian Matrix
$$f(x_1, \ldots, x_n)$$

$$\nabla f = \left(\frac{\partial f}{\partial x_1}, \ldots, \frac{\partial f}{\partial x_n}\right) \text{ Gradient.}$$

$$\nabla^2 f \begin{pmatrix} \frac{\partial^2 f}{\partial x_1 \partial x_1} & \frac{\partial^2 f}{\partial x_1 \partial x_2} & \cdots \\ & \ddots & \\ & & \frac{\partial^2 f}{\partial x_i \partial x_j} \end{pmatrix} \Leftarrow \text{Hessian}$$

$$\boxed{\nabla_\theta^2 A(\theta) = \text{Cov}_\theta(\phi(x))}$$
$\qquad\qquad\qquad\qquad \uparrow$
$\qquad\qquad\qquad$ covariance.

$$\text{Cov}(x) = \begin{bmatrix} & C_{ij} & \\ & = E[(x_i - Ex_i)(x_j - Ex_j)] \end{bmatrix}$$

$f: \Omega \to \mathbb{R}$  convex function.



$$x' = \alpha x_1 + (1-\alpha) x_2$$

$$y' = \alpha f(x_1) + (1-\alpha) f(x_2)$$

$$y' \geqslant f(x').$$

$$P(x) = \begin{cases} \alpha & (x = x_1) \\ 1-\alpha & (x = x_2) \end{cases}$$

$$E_P(x) = \alpha x_1 + (1-\alpha) x_2$$

$$f(E_P(x)) \leqslant E_P[f(x)]. \Rightarrow \boxed{f(E(x)) \leqslant E[f(x)]}$$
$\qquad\qquad\qquad\qquad\qquad$ for all convex
$\qquad\qquad\qquad\qquad\qquad$ function

$K\overset{\checkmark}{\overline{\mho}}\overline{\mathfrak{h}}$

**Another view**

$A$ is semi-definite

$\quad x^T A x \geq 0 \quad \forall x$
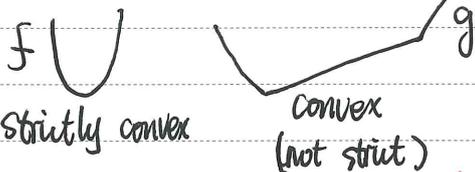
$\quad\quad x = \sum_{i=1}^{n} c_i e_i$

$x^T A x = x^T \left( \sum_{i=1}^{n} c_i A e_i \right)$

$\quad\quad = \left( \sum_{i=1}^{n} c_i e_i \right) \cdot \left( \sum_{i=1}^{n} c_i \lambda_i e_i \right)$

$\quad\quad = \underbrace{\sum_{i=1} \lambda_i c_i^2}_{\geq 0} + \underbrace{\sum_{i \neq j} c_i c_j \lambda_j \, e_i^T \cdot e_j}_{= 0}$

$A$ is definite, $\quad x^T A x > 0 \quad \forall x \neq 0$.

$f$: convex function

$f \; \cup \qquad \diagdown\!\diagup \; g$

strictly convex $\qquad$ convex (not strict)

$H = \nabla^2 f \rightarrow$ $\boxed{\text{positive definite}} \Leftrightarrow f$ is strictly $\Leftrightarrow$ $\boxed{\nabla f \text{ injective}}$
$\qquad\qquad\qquad\qquad\qquad\qquad$ convex $\qquad\qquad$ **Answer.**

$\qquad\qquad\qquad \searrow$ semi-definite $\Leftrightarrow f$ is convex

$\nabla f$ is injective.

$\nabla g$ is not injective.

$\nabla^2 A(\theta) = \text{Cov}_\theta (\phi(X))$.

$\boxed{\text{if overcomplete} \rightarrow A \text{ is } \boxed{\text{not}} \text{ strictly convex}}$
$\boxed{\qquad\qquad\qquad \rightarrow \nabla A \text{ is } \boxed{\text{not}} \text{ injective}}$

$X \quad \cancel{\text{Cov}} \; \text{Cov}(X) = C$.

$Y = a^T X$

$\searrow$ variance $\text{Var}(a^T X) = a^T C a = 0$

overcomplete: $a^T X = b$

$\qquad\qquad\qquad \downarrow$

$\qquad\qquad \text{Var}(a^T X) = 0$.

$\boxed{\text{if minimal} \Rightarrow a^T C a > 0 \rightarrow A \text{ is strictly convex}}$
$\boxed{\qquad\qquad\qquad \rightarrow \nabla A \text{ is injective}.}$

---

$\nabla A$ : Canonical parameter $\rightarrow$ mean parameter

overcomplete $\rightarrow \exists \; \theta_1 \neq \theta_2 \Rightarrow P_{\theta_1} = P_{\theta_2} \Rightarrow \mu_1 = \mu_2 \rightarrow$ not injective
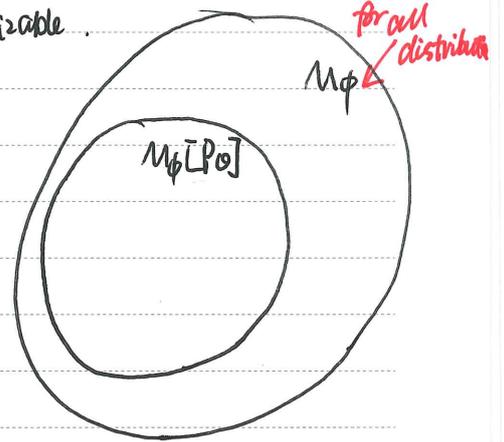
minimal $\rightarrow \forall \theta_1 \neq \theta_2 \Rightarrow P_{\theta_1} \neq P_{\theta_2} \Rightarrow \mu_1 \neq \mu_2 \Rightarrow$ injective.

Surjective : $f : \Omega \rightarrow M \quad M\phi(P_\theta) = M\phi$ ?

$X \cdot \phi(X) \quad$ arbitrary distribution $P$ $\qquad$ True for all regular exponential family

$\qquad\qquad \mu = E_p(X)$

$\qquad\qquad \uparrow$
$\qquad\qquad$ realizable.

$\qquad\qquad\qquad\qquad\qquad\qquad\qquad\qquad\qquad$ for all distribution

$\qquad\qquad\qquad\qquad\qquad\qquad M\phi$

$\qquad\qquad\qquad\qquad M\phi[P_\theta]$

$X \cdot P$

entropy $\rightarrow H(P) = -\int_X P(x) \log P(x) \mu(dx)$.

$\qquad P_1, P_2, \ldots, P_n \quad \mu$

$\qquad E_{P_1}[\phi(X)] = \cdots \quad E_{P_n}[\phi(X)] = \mu$.

Maximum Entropy $\leftarrow$ is the best choice.

$\boxed{\text{Usually, more entropy means less information we have}}$

Maximize $H(P)$ s.t. $E_P[\phi(X)] = \mu$.

$X = \{x_1, x_2, \ldots x_k\}$ finite

$P = (P_1, P_2, \ldots P_k)$.

$\cancel{\text{maximize}} - \sum_{i=1}^{k} P_i \log P_i \quad$ s.t. $\sum_{i=1}^{k} P_i \phi(x_i) = \mu \quad \lambda_i$
minimize.
$\qquad\qquad\qquad\qquad\qquad\qquad\qquad \sum_{i=1}^{k} P_i = 1, \; P_i \geq 0$.
$\qquad\qquad\qquad\qquad\qquad\qquad\qquad\qquad \underset{V}{\smile}$

$L(P, \lambda, V) = \sum_{i=1}^{k} P_i \log P_i$

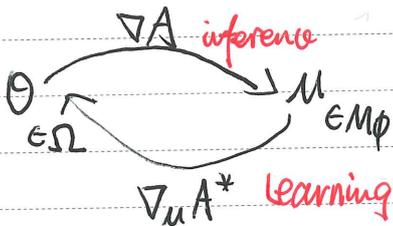$\qquad - \lambda_i \sum_{i=1}^{k} P_i \phi_i - V \sum_{i=1}^{k} P_i$

$$\frac{\partial L}{\partial P_i} = (\log P_i + 1) - \lambda_i \phi_i - \nu = 0$$

$$\log P_i = \lambda_i \phi_i + \nu - 1$$

$$P_i = \exp(\lambda_i \phi_i + \nu - 1)$$

$$P = \frac{\exp(\lambda^T \phi)}{\exp(1 - \nu)} = \frac{1}{Z} \exp(\lambda^T \phi)$$

$$P(x) = \frac{1}{Z} \exp(\lambda^T \phi(x))$$



inference / learning diagram: $\Theta \in \Omega$, $\nabla A$ inference, $\mu \in M\phi$, $\nabla_\mu A^*$ learning

**Convex Conjugate**

$$f : \Omega \to \mathbb{R} \quad \Omega \in \mathbb{R}^d$$

$$\boxed{f^*(y) = \sup_{x \in \Omega}(y^T x - f(x)) \Rightarrow f^*(y) \geqslant y^T x - f(x)}$$

$$y \to \underset{x \in \Omega}{\text{maximize }} y^T x - f(x) \to \hat{x}$$

$$\downarrow$$

$$\boxed{y^T \hat{x} - f(\hat{x})}$$

$$= \sup_{x \in \Omega}(y^T x - f(x))$$

$$\boxed{f^* \text{ is always convex}}$$

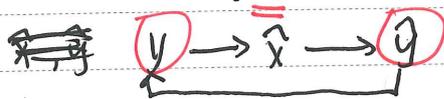**Fenchel's inequality**

$$f(x) + f^*(y) \geqslant y^T x$$

**Fenchel-Moreau Theorem**

$$f^{**} = f \text{ iff } f \text{ is convex and continuous}$$

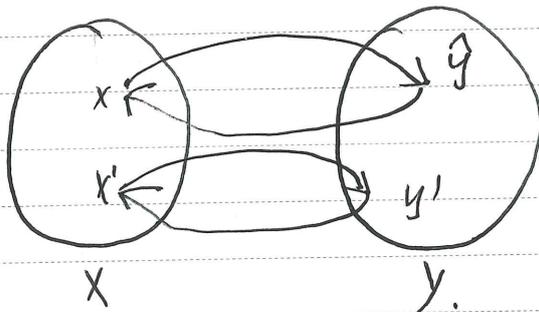$$\begin{cases} f^*(y) = \sup_x (y^T x - f(x)) \\ f(x) = \sup_y (y^T x - f^*(y)) \end{cases}$$

$$f^*(y) = y^T \hat{x} - f(\hat{x})$$

$$f(x) = \hat{y}^T x - f^*(\hat{y})$$

$$\boxed{y} \longrightarrow \hat{x} \longrightarrow \boxed{\hat{y}}$$

$\hat{x}, \hat{y}$ dually coupled

$$f(\hat{x}) + f^*(\hat{y}) = \hat{x}^T \hat{y}$$



$X$ , $Y$

$$A^*(\mu) = \sup_\Theta \{\Theta \mu - A(\Theta)\}$$

$$\text{maximize } \underset{L(\Theta)}{\Theta^T \mu - A(\Theta)}$$
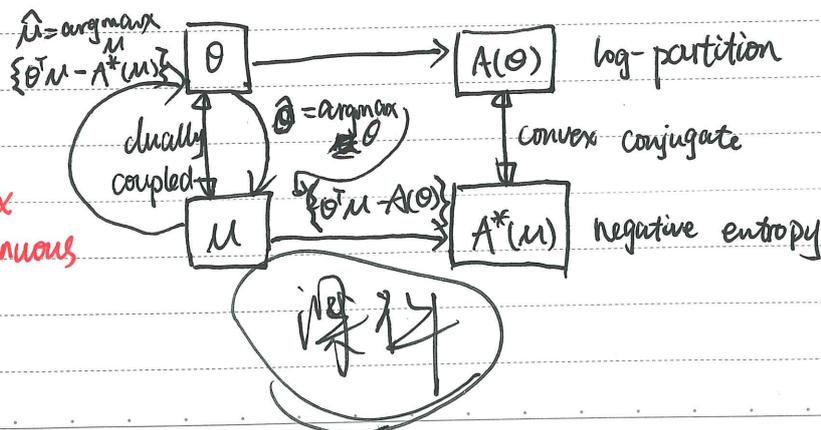
$$\nabla_\Theta L(\Theta) = \mu - \nabla A(\Theta) = 0$$

$$\mu = \nabla_\Theta A(\Theta) = E_\Theta[\phi(x)]$$

$$A(\Theta) = \sup_\mu \{\Theta^T \mu - A^*(\mu)\}$$

$$\hat{\Theta} = \nabla_\mu A^*$$

$$A^*(\mu) = \begin{cases} -H(P_\Theta) & \mu \in M\phi \quad \text{(negative entropy)} \\ +\infty & \text{otherwise} \end{cases}$$



$\hat{\mu} = \underset{\mu}{\text{argmax}} \{\Theta^T \mu - A^*(\mu)\}$, $\Theta \to A(\Theta)$ log-partition

$\hat{\Theta} = \underset{\Theta}{\text{argmax}} \{\Theta^T \mu - A(\Theta)\}$

dually coupled, $\mu$, $A^*(\mu)$ negative entropy, convex conjugate

嗶14

Conjuagate Priors.

$$P(x|\theta)$$
observe, parameter    在Bayesian中 估计 $P(\theta)$ prior

$$P(\theta|x_1,\cdots y_m) = \frac{1}{Z} P(\theta) \prod_{i=1}^{n} P(x_i|\theta)$$   计算Z很复杂.

选择合适的 $P(\theta)$, 简化Z的计算.

参数的贝叶斯估计, 需要选择合适的 prior.

Bernoulli distribution.

$$P(x|\theta) = \begin{cases} \theta & (x=1) \\ 1-\theta & (x=0) \end{cases} = \theta^x(1-\theta)^{1-x}$$

要找一个先验 $P(\theta)$

如果已经有 样本 $x=(x_1 \cdots x_n)$

$$P(\theta|x) = \frac{1}{Z} \prod_{i=1}^{n} P(x_i|\theta) P(\theta)$$

$$Z = \int_{\theta \in \Omega} \prod_{i=1}^{n} P(x_i|\theta) p(\theta) d\theta$$

如果 $P(\theta)$

$$P(\theta|\alpha,\beta) = \frac{1}{B(\alpha,\beta)} \theta^{\alpha-1} \theta^{\beta-1}$$   Beta分布

$$P(\theta|x) = \frac{1}{Z} P(\theta|\alpha,\beta) \prod_{i=1}^{n} P(x_i|\theta)$$

$$= \frac{1}{Z} \cdot \frac{1}{B(\alpha,\beta)} \theta^{\alpha-1}(1-\theta)^{\beta-1} \prod_{i=1}^{n} \theta^{x_i}(1-\theta)^{1-x_i}$$

$$= \frac{1}{Z} \frac{1}{B(\alpha,\beta)} \theta^{\alpha-1+\sum_{i=1}^{n} x_i} (1-\theta)^{\beta-1+\sum_{i=1}^{n}(1-x_i)}$$

$$= \frac{1}{Z} \frac{1}{B(\alpha,\beta)} \theta^{\alpha'-1}(1-\theta)^{\beta'-1}$$

$\underbrace{\qquad\qquad\qquad\qquad}_{\text{posterior distribution}}$

$$\alpha' = \alpha + \sum_{i=1}^{n} x_i \qquad \beta' = \beta + \sum_{i=1}^{n}(1-x_i)$$

→ updating formula.

$$B(\alpha,\beta) = \int_{\Omega} \theta^{\alpha-1}(1-\theta)^{\beta-1} d\theta.$$

把Z合并进 $B(\alpha',\beta')$

$$Z \cdot B(\alpha,\beta) = \int_{\Omega} \theta^{\alpha'-1}(1-\theta)^{\beta'-1} d\theta$$

$$= B(\alpha',\beta')$$

用了 conjugate 后, posterior 和likelihood 形式 相同, 只是改变了参数.

下面 芥式 介绍 conjugate prior.

和 exp family 紧密相关, 如果 likelihood 不是 exp family, 则 没法用 conjugate.

Conjugate prior.

$$P(\theta|\alpha)$$

$$P(x|\theta) \quad \text{likelihood model.}$$

$$D = \{x_1,\cdots, x_n\}$$

$$P(D|\theta) = \prod_{i=1}^{n} P(x_i|\theta)$$

$$P(\theta|D) = P(\theta|\alpha') \quad \text{posteri.}$$

$$\alpha' = \alpha \oplus D \quad \text{用 given samples D 更新 } \alpha.$$

$$= \alpha \oplus \{x_1,\cdots,x_n\}.$$

如何选择 conjugate prior.

$$f(x|\theta) = h(x) \exp(\eta(\theta)^T \phi(x) - \gamma \cdot a(\theta))$$

$$P(\theta|\alpha,\beta) = \exp(\alpha^T \eta(\theta) - \beta \cdot a(\theta) - A(\alpha,\beta))$$
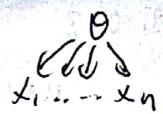
$$P(\theta|D) \propto \quad \overline{\text{exp}(\alpha^T\eta(\theta)-\beta\cdot a(\theta))}$$

prior → $\exp(\alpha^T \eta(\theta) - \beta \cdot a(\theta) - A(\alpha,\beta))$

likelihood · $\prod_{i=1}^{n} \exp(\eta(\theta)^T \phi(x_i) - \gamma \cdot a(\theta))$

$$\propto \exp\left(\underbrace{\left(\alpha + \sum_{i=1}^{n} \phi(x_i)\right)}_{\alpha'}^T \eta(\theta) - (\beta+n\gamma) a(\theta)\right.$$

$$\left. - A(\alpha',\beta')\right)$$

构造公式：直接猜 $\phi(x)$ 的充分统计量.
下面看 bernoulli likelihood 直接来构造

$f(x|\theta) = \theta^x(1-\theta)^{1-x}$
$\qquad = exp(x \cdot log\theta + (1-x)log(1-\theta))$

$P(\theta|\alpha,\beta) \propto exp(\alpha \cdot log\theta + \beta\, log(1-\theta))$

$\qquad = \frac{1}{Z(\alpha,\beta)} \cdot exp(\alpha \cdot log(\theta) + \beta \cdot log(1-\theta))$

$\qquad = \frac{1}{Z(\alpha,\beta)} \cdot \theta^\alpha \cdot (1-\theta)^\beta$

$\therefore P_{beta}(\theta|\alpha,\beta) \propto \theta^{\alpha-1}(1-\theta)^{\beta-1}$

$\theta$
$x_1 \cdots x_n$

$P(x|D) = \int_\Omega P(x|\theta)P(\theta|D)d\theta \longrightarrow P(\theta|\alpha',\beta')$

$\qquad = \int_\Omega exp(\eta(\theta)^T\phi(x) - a\cdot r(\theta))$

$\qquad exp(\alpha'^T\eta(\theta) - \beta'\, s(\theta) - A(\alpha',\beta'))d\theta$

$\qquad = \int_\Omega exp([\alpha'-\phi(x)]^T\eta(\theta) - (\beta'+a)r(\theta) - A(\alpha',\beta'))d\theta$

$\qquad = \frac{1}{exp(A(\alpha',\beta'))}\int_\Omega exp[(\alpha'-\phi(x))^T\eta(\theta) - (\beta'+a)\cdot s(\theta)]d\theta$

$\qquad\qquad exp(A(\alpha'-\phi(x),\beta'+a))$

下一节讲 inference, given model! 讲推断.

Dirichlet Distribution

Beta Distribution 的扩展.

Bernoulli $\xrightarrow{\text{两类→多类}}$ Categorical

$\uparrow$prior $\qquad\qquad \uparrow$ prior

Beta $\longrightarrow$ Dirichlete

~~多项分布~~ $p(x|\pi) = \pi_k$.

$= exp(\sum_{k=1}^K \mathbb{1}(x=k)log\pi_k)$
充分统计量.

prior: $P(\pi|\alpha) = \frac{1}{Z(\alpha)}exp(\sum_{k=1}^K \alpha_k\, log\pi_k)$

$\qquad = \frac{1}{Z(\alpha)} \cdot \prod_{k=1}^K \pi_k^{\alpha_k}$.

规范化:
$P_{prior}(\pi|\alpha) = \frac{1}{B(\alpha)}\prod_{k=1}^K \pi_k^{\alpha_k-1}$

下面看 Dir 的性质.

$P_{prior}(x|\alpha) = \frac{1}{B(\alpha)} = \prod_{k=1}^K x_k^{\alpha_k-1}$

所有的 x 都是一个概率分布 ∴ Dir 是分布的分布

$E[x_k] = \frac{\alpha_k}{\sum \alpha_i}$

# Topic 2  Inference

$$P(x; \theta) \Leftarrow Model$$

$$\downarrow$$

$$(x_1, \cdots, x_n)$$

已知模型参数 $\theta$ 和一部分
变量 $x_B$，求另一部分的条件概率

$$P(\underset{\uparrow}{x_M} | \underset{\uparrow}{x_B}; \underset{\nwarrow}{\theta})$$

query  observation  model, learnt in the past
evidences.

假设变量分为三部分： X Y Z  已知 $p(x, y, z | \theta)$

$$P(Y | X; \theta) = \frac{P(Y, X | \theta)}{P(X | \theta)}$$

本来是 conditional distr.
转化成了 Marginal distr.

要求两部分:

$$P(Y, X | \theta) = \sum_{z \in Z} P(Y, X, z | \theta)$$

$$P(X | \theta) = \sum_{y \in Y} P(X, y | \theta)$$

计算会指数增长.

Evidence Absorption     (一个 trick)

$$P(x,y,z) \propto \psi_x(x)\,\psi_y(y)\,\psi_z(z)$$

$$\phi_{xy}(x,y)\,\phi_{yz}(y,z)\,\phi_{xz}(x,z)$$

$$P(X,Y|\underset{\underset{observation}{\uparrow}}{Z}) = \frac{P(x,y,z)}{\sum_{x,y} P(x,y,z)} = \frac{\psi_x(x)\psi_y(y)\psi_z(z)\phi_{xy}(x,y)\phi_{xz}(x,z)\phi_{yz}(y,z)}{\sum_{x'}\sum_{y'}\psi_x(x')\psi_y(y')\psi_z(z)\phi_{xy}(x',y')\phi_{xz}(x',z)\phi_{yz}(y',z)}$$

$$= \frac{\psi_x(x)\psi_y(y)\phi_{xy}(x,y)\phi_{x|z}(x)\phi_{y|z}(y)}{\sum_{x'}\sum_{y'}\cdots}$$

$$\propto \psi_x(x)\psi_y(y)\phi_{xy}(x,y)\phi_{x|z}(x)\phi_{y|z}(y)$$

∵ z 是 给定 constant, 分子分母的 $\psi_z(z)$ 可以 约掉，出来消掉.

$\phi_{xz}(x,z)$ 实际上是个 矩阵:
∵ x 固定, ∴ 只看其中一行

| x\z | 0 | 1 | 2 | 3 |
|-----|---|---|---|---|
| 0   |   |   |   |   |
| 1   |   |   |   |   |
| 2   |   |   |   |   |
| 3   |   |   |   |   |

∴ $\phi_{xz}(x,z) \rightarrow \phi_{x|z}(x)$

∵ 开始有三个变量 x, y, z, 当经过 evidence Absorption 之后, 变成了只有 x, y 的 MRF, 形式更简单.

下面 看 计算 Marginal probability.

$P(x,y)$ 求 $P(x)$

$$P(x) = \sum_y P(x,y)$$

Complexity : key difference between statistics

and machine learning

核心问题 是 如何 计算 Marginal distr. efficiently

Y 可能 有 很多个变量. $O(|Y|) \rightarrow K^n$

conditional independent is the key to reduce

the complexity of the summation



$$P(x,y,z,w) = \frac{1}{Z} f(x,y) \, g(y,z) \cdot h(y, w)$$

$P(x)?$   marginalization : remove
other variables only keep
marginal distr.

$$P(x) = \sum_y \sum_z \sum_w \frac{1}{Z} f(x,y) g(y,z) h(y,w)$$

$$= \frac{1}{Z} \underbrace{\sum_y \sum_z \sum_w f(x,y) g(y,z) h(y,w)}_{\tilde{P}(x)}$$

一般只关注 $\tilde{P}(x)$ 即 有 $p(x) = \frac{1}{Z} \tilde{P}(x)$; $Z = \sum_x \tilde{P}(x)$

$\hat{p}(x)$ 的复杂度: $O(m_y \cdot m_z \cdot m_w)$ for a single $x$

$\therefore$ Overall complexity: $O(m_x, m_y \cdot m_z, m_w)$    $O(m^4)$

利用 structure 简化计算.

$$\sum_x (\cdot f(x) = C \sum_x f(x)$$ 利用空间结果简化.

$$\hat{p}(x) = \sum_y \sum_z \sum_w f(x,y) g(y,z) h(y,w)$$

$$= \sum_y \sum_z f(x,y) g(y,z) \sum_w h(y,w)$$

$$= \sum_y f(x,y) \sum_z g(y,z) \sum_w h(y,w)$$

下面看如何简化计算的.

令 $h_{/w}(y) = \sum_w h(y,w)$

$$= \sum_y f(x,y) \underbrace{\underbrace{g_{/z}(y)}_{①} \underbrace{h_{/w}(y)}_{②}}_{③}$$

复杂度: ① : $O(m_y \cdot m_z)$

$\qquad + ②: O(m_y, m_z) \Rightarrow O(m^2)$

$\qquad + ③: O(m_x m_y)$

对于不同的 $x$, $h_{/w}(y)$ 都是相同的, $\therefore$ 只一次计算

以上的方法叫 Variable elimination

$$g_z(y) \qquad h_w(v)$$

$$\uparrow \qquad\qquad \uparrow$$

eliminate z $\quad$ eliminate w

条件独立的作用在于 例如 x 和 (z,w) 独立, 对于任意 x, z,w 有影响基独立, 可以一次计算.

# Variable Eliminate

算法描述:

$$Y_0, Y_1, \cdots, Y_n$$

$$\mathcal{F} = \{\phi_1, \cdots, \phi_m\}$$

$$\mathcal{V} = \{Y_0, Y_1, \cdots, Y_n\} \qquad$$ 未消去的变量集合.

$$j = 1, \cdots, n$$

$$i = \pi(j)$$ 决定 Elimination 的顺序.

顺序对 complexity 影响很大.

$$\mathcal{F}, \mathcal{V} = \text{Var Eliminate}(\mathcal{F}, \mathcal{V}, Y_i)$$

$f(Y_i)$: the set of factors involving $Y_i$

上例中 $f(w) = \{h\}$ $f(z) = \{g\}$

$\nu(\phi)$: active variables involved in $\phi$

上例中 $\nu(h) = \{y, w\}$ $\nu(g) = \{y, z\}$

Neighbour of $Y_i$:

$$N_i = \{v \neq Y_i : \exists \phi \in f(Y_i), v \in \nu(\phi)\}$$

上例中: $f(y) = \{f, g, h\}$

$N(y) = \{x, z, w\}$

$f(w) = \{h\}$ $N(w) = \{y\}$

Construct $\psi_i$ on $N_i$

$$\psi_i(z) = \sum_y \prod_{\phi \in f(Y_i)} \phi\left(y, z \mid \nu(\phi)\right)$$

下面看一个例子，应用该算法.

$$\mathcal{F} = \{\psi_x, \psi_y, \psi_z, \psi_w, \phi_{xy}, \phi_{xz}, \phi_{yz}, \phi_{yw}, \phi_{zw}\}$$

$$\mathcal{V} = \{x, y, z, w\}$$

$$\Pi = (w, z, y)$$

① Eliminate $w$

$$\mathcal{F}(w) = \{\psi_w, \phi_{yw}, \phi_{zw}\}$$

$$N_w = \{y, z\}$$

$$\gamma(y, z) = \sum_w \psi(w) \, \phi_{yw}(y, w) \, \phi_{zw}(z, w)$$

$$\Downarrow$$



② Eliminate $z$

$$\mathcal{F}(z) = \{\phi_{xz}, \phi_{yz}, \gamma_{yz}, \psi_z\}$$

$$N(z) = \{x, y\}$$

$$\gamma(x, y) = \sum_z \phi_{xz}(x, z) \, \phi_{yz}(y, z) \, \gamma_{yz}(y, z) \, \psi_z(z)$$

$$\Downarrow$$

③ Eliminate $y$

$$F(y) = \{\phi_{xy}, \psi_y, \gamma_{xy}\}$$

$$N(y) = \{x\}$$

$$\gamma_x(x) = \sum_y \phi_{xy}(x, y)\psi_y(y)\gamma_{xy}(x, y)$$

$$\Downarrow$$

$$\psi_x \underline{\text{III}} \bigotimes - \underline{\text{III}} \gamma_x(x)$$
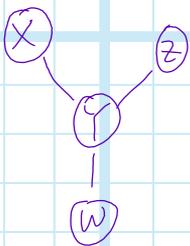
$$\therefore P(x) \propto \psi_x(x)\gamma_x(x)$$

下面分析复杂度：

① : $O(m_y \cdot m_z \cdot m_w)$

② : $O(m_x \cdot m_y \cdot m_z)$

③ $O(m_x \cdot m_y)$

下面来看 order $\pi(j)$ 对 Complexity 的影响。



求 $P(x)$.

$\pi_1 = (y, w, z) \longrightarrow O(m^4)$ $\because y$ 有很多 neighbour.

$\pi_2 = (w, z, y) \longrightarrow O(m^2)$

那么如何选择 optimal order
是一个 NP 问题, 一般靠经验选.
一般优先选 neighbour 少的变量来做假.


例题: 计算复杂度.
A chain of discreate variable

$$\textcircled{x_1} - \textcircled{x_2} - \cdots - \textcircled{x_n} \qquad n > 3$$
求 $P(x_1)$

space conditionality: m 每个 $x_i$ 有 m 个取值.

— direct formulation $O(?)$ $\qquad O(m^n)$
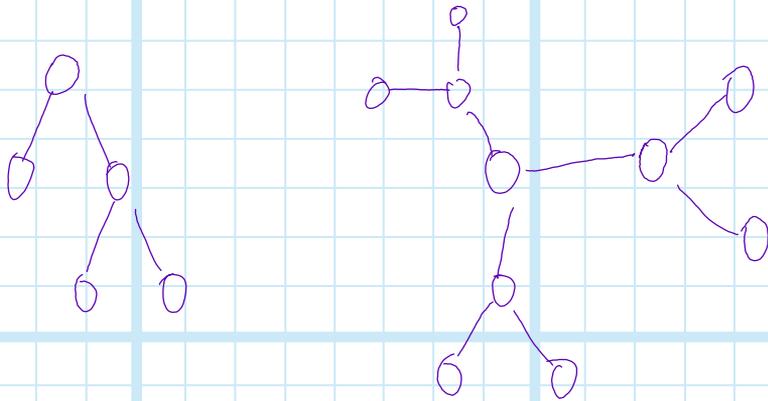
— variable elimination $O(?)$ $\qquad O(n \cdot m^2)$

问题, 如果要对很多个 x 计算 P(x), 里面也有很多相同的计算, 能不能复用? 下面讲 Belief Propagation

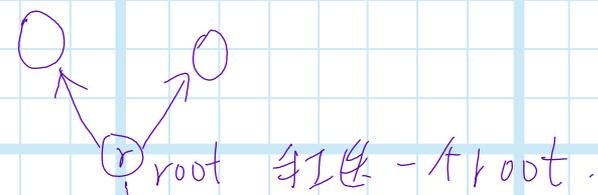# Belief Propagation

如何对所有的 x 来计算 $P(x)$.

Tree-structured Model.



Tree: contains no cycle.

$$P(x) = \frac{1}{Z} \prod_{s \in V} \psi_s(x_s) \prod_{(s,t) \in E} \psi_{st}(x_s, x_t)$$

↑ uniary term    ↑ binary term

← Pairwise MRF

Tree-structured 的优点: 存在很叻(好)的 inference 算法.



root 扫选一个 root.

为 edge 选择方向 一个 parent 和 n children.

$$P(x) = \frac{1}{Z} \prod_{s \in V} \psi_s(x_s) \prod_{s \in V / \{root\}} \phi_s(x_{\pi(s)}, x_s)$$

$\overbrace{\quad}^{n}$     $\overbrace{\quad}^{n-1}$

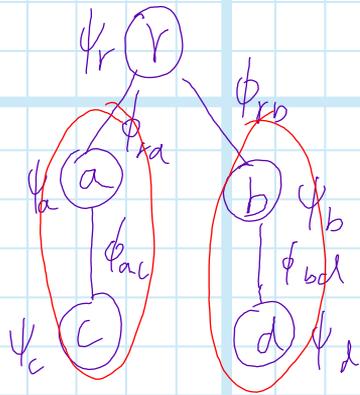$$= \frac{1}{Z} \psi_r(x_r) \prod_{s \in V / \{r\}} \psi_s(x_s) \phi_s(x_{\pi(s)}, x_s)$$

下面看 root 的 marigionul distr, 然后扩展
到所有的结点. 最后扩展到体姜的 graph

# Tree-Structured Model

$$p(x) = \frac{1}{Z} \psi_r(x_r) \prod_{s \in V(T) \setminus r} \psi_s(x_s) \phi_s(x_{\pi(s)}, x_s)$$

每个非根结点和父结点相连

先计算如下的 $P(x_r)$，然后 generalize 到一般模型



$$P(x_r) = \frac{1}{Z} \sum_{x_a} \sum_{x_b} \sum_{x_c} \sum_{x_d} \psi_r(x_r) \psi_a(x_a) \psi_b(x_b) \psi_c(x_c) \psi_d(x_d)$$
$$\phi_{ra}(x_r, x_a) \phi_{rb}(x_r, x_b) \phi_{ac}(x_a, x_c) \phi_{bd}(x_b, x_d)$$

$$= \frac{1}{Z} \psi_r(x_r) \sum_{x_a} \psi_a(x_a) \phi_{ra}(x_r, x_a) \sum_{x_c} (x_c) \phi_{ac}(x_a, x_c) \quad \text{第一行只和}\, a,c\, \text{有关}$$

$$\sum_{x_b} \psi_b(x_b) \phi_{rb}(x_r, x_b) \sum_{x_d} \psi_d(x_d) \phi_{bd}(x_b, x_d) \quad \text{第二行只和}\, b,d\, \text{有关}$$

可以发现是按 sub tree 分解的
decomposition along sub trees.
先介绍几个术语：

$T_s$ : sub-Tree rooted at $s$
$Ch(s)$ : children of $s$      $ch(r) = \{a, b\}$,   $ch(a) = \{c\}$

$V(s): V(T(s))$     例 . $V(a)=\{a,c\}$   all the vertices contained in $T(s)$
$V(r)=V$

$D(s)= V(s)/\{s\}$   decedents

定义 $w_s(x_{V(s)})= \psi_s(x_s)\prod_{t\in D(s)} \psi_t(x_t)\phi_t(x_{\pi(t)},x_t)$

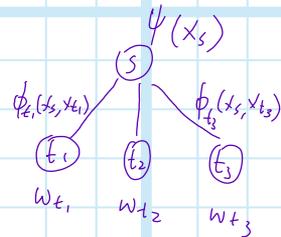$w_r(x_r)= \hat{P}(x_r)$

$P(x) \propto w_r(x_V)$   只差了个 $z$

Leaf $t$:   $w_t(x_t)= \psi_t(x_t)$

Non leaf $s$:

$w_s(x_{V(s)})= \psi(x_s)\prod_{t\in Ch(s)} \phi_t(x_s,x_t) w_t(x_{V(t)})$

$w_r(x)= \psi_r \cdot \psi_a \cdot \psi_b \cdot \psi_c \cdot \psi_d$   递归求解

$\phi_{ra}\cdot \phi_{rb}\cdot \phi_{ac}\cdot \phi_{bd}$

$= \psi_r \cdot (\phi_{ra} \cdot \underline{\psi_a \psi_c \phi_{ac}}) \triangledown w_a$

$(\phi_{rb}\cdot \underline{\psi_b \cdot \psi_d \cdot \phi_{bd}}) \rightarrow w_b$

$= \psi_r (\phi_{ra} w_a)(\phi_{rb} w_b)$

有了 $w_s$, 就可以计算:

$f_s(x_s)= \sum_{x_{D(s)}} w_s(x_s, x_{D(s)})$

然后 $P(x_r) \propto f_r(x_r)$

下面看 $f_s(x_s)$ 的计算

Leaf $t$:

$f_t(x_t)= \psi_t(x_t)$   ∵没有 children

$D(s)= V(s)\backslash s$ , $f_s(x_s)$ 就是 把除 $x_s$ 外 的 所有变量 Marginalize 掉

∴ 计算整个网络用 $w_s$, 而计算 某个变量的 marginal 用 $f_s$

non-leaf $s$:

$$f_s(x_s) = \sum_{\substack{x_{D(s)} \\ \text{定-}X}} w_s(x_s; x_{D(s)})$$

$$= \psi_s(x_s) \cdot \sum_{x_{D(s)}} \begin{array}{c} \phi_{t_1}(s,t_1) \, w_{t_1}(x_{V(t_1)}) \\ \vdots \\ \phi_{t_k}(s,t_k) \cdot w_{t_k}(x_{V(t_k)}) \end{array}$$

$$= \psi_s(x_s) \cdot \sum_{x_{V(t_1)}} \phi_{t_1}(s,t_1) w_{t_1}(x_{V(t_1)})$$
$$\vdots$$
$$\sum_{x_{V(t_k)}} \phi_{t_k}(s,t_k) \cdot w_{t_k}(x_{V(t_k)})$$

$$= \psi_s(x_s) \prod_{t \in Ch(s)} \cdot \sum_{x_{V(t)}} \phi_t(x_s, x_t) \, \omega_t(x_{V(t)})$$

$$= \psi_s \prod_{t \in Ch(s)} \sum_{x_t, x_{D(t)}} \phi(x_s, x_t) \, w_t(x_{V(t)})$$

$$= \psi_s \prod_{t \in Ch(s)} \sum_{x_t} \phi(x_s, x_t) \sum_{x_{D(t)}} w_t(x_t, x_{D(t)})$$

$$= \psi_s \prod_{t \in Ch(s)} \sum_{x_t} \phi_t(x_s, x_t) f_t(x_t)$$

$D(r)$ 可以分解为不相交的

$(a,c)$ 和 $(b,d)$

$$Ch(s) = \{t_1, \cdots, t_k\}$$

$$D(s) = V(t_1) \cup \cdots \cup V(t_k)$$

$$w_s(x_{V(s)}) = \psi_s(x_s) \prod_{t \in D(s)} \psi_t(x_t) \phi_t(x_{\pi(t)}, x_t)$$

Goal:

$$p(x_r) \propto f_r(x_r)$$

Leaf $t$: $f_t(x_t) = \psi_t(x_t)$

Non-leaf:

$$f_s(x_s) = \psi_s(x_s) \prod_{t \in Ch(s)} \cdot \sum_{x_t} \phi_x(x_s, x_t) f_t(x_t)$$
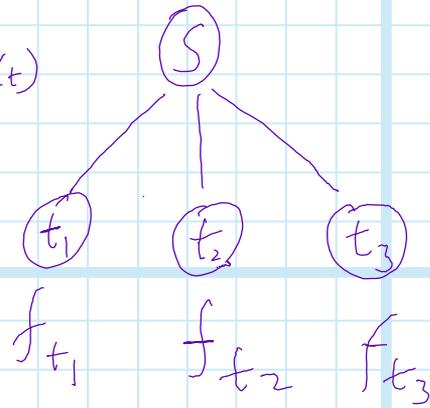
复杂度: $O(m_s) \cdot \left( O(|Ch(s)|) \underbrace{\sum_{t \in Ch(s)} O(m_s \cdot m_t)}_{} \right)$

下面从 Message 的 成来理解这个式子

定义 $M_{t \to s}(x_s) \triangleq \sum\limits_{x_t} \phi_t(x_s, x_t) f_t(x_t)$

那么 $f_s(x_s) = \psi_s(x_s) \prod\limits_{t \in Ch(s)} M_{t \to s}(x_s)$



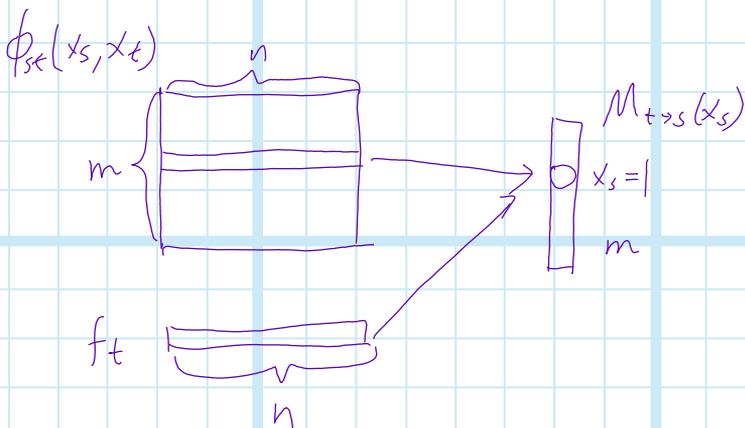$f_{t_1} \quad f_{t_2} \quad f_{t_3}$

下面看如何实现.

s 有 m 个取值
t 有 n 个取值.

$m$   $s$   $f_s = \left[\!\right] \}m$

$\phi_{st_1}\left[\!\right]$   $\left[\right]$   $\left[\right]\}m$   t 固定, s 有 m 种取值.

$n \quad t_1 \qquad t_2 \qquad t_3$

$f_{t_1} \left[\!\right]\}n$

$M_{t_1 \to s}(x_s) = \sum\limits_{x_{t_1}} \phi_{st}(x_s, x_{t_1}) f_t(x_{t_1})$

$f_{t_1}$ 就是个 n 维向量.

$f_s(y_s) = \psi_s(x_s) \cdot \prod\limits_{t \in Ch(s)} M_{t \to s}(x_s)$

下面看 $M_{t \to s}(x_s)$

$\phi_{st}(x_s, x_t)$



$M_{t \to s}(x_s)$

$x_s = 1$

$m$

$f_t$

$n$

# How to read a paper

看 key graph

Basic idea 是什么, 在论文中的作用.
例如 corelated topic model 扯.盖用了个
softmax adapter

key advantage your model can bring:
Boundary of the model.

# How you should formulate the model

Identify the key problem.
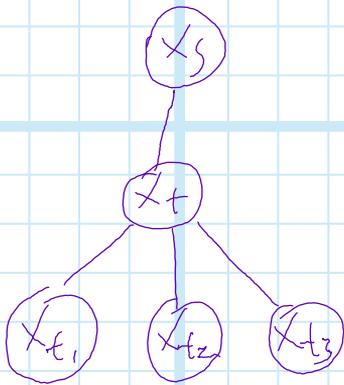Understand the problem, why existing model
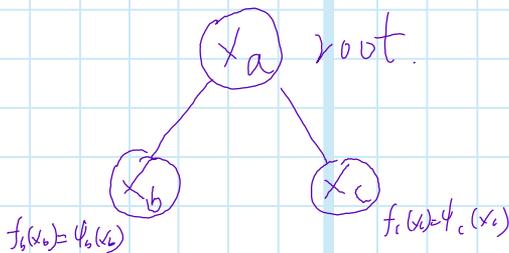not good.
数字分布主要的兔鼠乙类.

Belief propagation 回顾.



Tree based

send message from children to parent.

$M_{t \to s}(x_s)$

$$M_{t \to s}(x_s) = \sum_{x_t} \phi(x_s, x_t) \psi_t(x_t) \prod_{u \in h(t)} M_{u \to t}(x_t)$$

上式只计算 root 的边缘分布, 如何计算所有结点的边缘分布.



root.

$f_b(x_b) = \psi_b(x_b)$       $f_c(x_c) = \psi_c(x_c)$

$$M_{b \to a}(x_a) = \sum_{x_b} \phi_{ab}(x_a, x_b) \psi_b(x_b)$$

$$M_{c \to a}(x_a) = \sum_{x_c} \phi_{ac}(x_a, x_c) \psi_c(x_c)$$

$$f_a(x_a) = \psi_a(x_a) \cdot M_{b \to a}(x_a) \cdot M_{c \to a}(x_a)$$

$$= \psi_a(x_a) \sum_{x_b} \phi_{ab}(x_a, x_b) \psi_b(x_b)$$

$$\sum_{x_c} \phi_{ac}(x_a, x_c) \cdot \psi_c(x_c)$$

$$= \sum_{x_b} \sum_{x_c} \psi_a(x_a) \psi_b(x_b) \psi_c(x_c) \phi_{ab}(x_a, x_b) \phi_{ac}(x_a, x_c)$$

得出了最后的定义.

那么如何计算 b 和 c 的边缘分布?

$$P(x_b) = \sum_{x_a} \sum_{x_c} \psi_a(x_a) \cdot \psi_b(x_b) \cdot \psi_c(x_c) \phi_{ab}(x_a, x_b) \phi_{ac}(x_a, x_c)$$

$$= \psi_b(x_b) \underbrace{\sum_{x_a} \psi_a(x_a) \cdot \phi_{ab}(x_a, x_b) \cdot}_{} \underbrace{\sum_{x_c} \psi_c(x_c) \phi_{ac}(x_a, x_c)}_{定义为\ M_{a \to b}(x_b)}$$

rewrite

$$= \sum_{x_a} \phi_{ab}(x_a, x_b) \psi_a(x_a) M_{c \to a}(x_a)$$

$$= \sum_{x_a} \phi_{ab}(x_a, x_b) \psi_a(x_a) \sum_{x_c} \phi_{ac}(x_a, x_c) \cdot \psi_c(x_c)$$

$$\therefore P(x_b) = \psi_b(x_b) M_{a \to b}(x_b)$$

其实就是把 b 当 root

求谁的概率 就把谁当 root.



这里计算需要考虑谁是 parent, 谁是 children, 下面给出不依赖于 parent-children 选择的方法.

General Form of Message Passing

$$M_{t \to s}(x_s) = \sum_{x_t} \phi(x_s, x_t) \psi_t(x_t) \prod_{u \in N(t) \setminus s} M_{u \to t}(x_u)$$

只需要把邻居 $N(t)$ 中的那个 target
node $s$ 删掉.

$$M_s(x_s) \propto \psi_s(x_s) \prod_{t \in N(s)} M_{t \to s}(x_s)$$



$N(t) = \{s, u_1, u_2, u_3\}$

$N(t) \setminus s = \{u_1, u_2, u_3\}$

下面分析复杂度.

Complexity Analysis.

each edge $(s, t)$    $M_{s \to t}(x_t)$    $M_{t \to s}(x_s)$

                                $|x_t|$         $|x_s|$

多少个 message.

$- \sum\limits_{s \in V} \deg(s) \cdot |x_s|$     space-complexity

                           to store all messages

                                 $|x|$

$2 n_{edges} |x| = 2(|V|-1) \cdot |x|$   总共这么多个 message.

for a tree    $|E| = |V| - 1$

$-$ time complexity 每个 message 的复杂度.

$M_{t \to s}(x_s)$

    compute $|x_s|$ values.

     $|x_t| \cdot$ terms / values

    $\therefore O(|x_s||x_t|) = O(m^2)$

$\therefore$ 总复杂度为 $O(|V| \cdot |x|^2)$

下面看例子.

Star Graph

$\forall t = 1, \cdots, n$

$$M_{t \to 0}(x_0) = \sum_{x_t} \phi_t(x_0, x_t) \psi_t(x_t)$$



$$M_{0 \to t}(x_t) = \sum_{x_0} \phi_t(x_0, x_t) \prod_{u \in N_0 / \{t\}} M_{u \to 0}(x_0)$$

计算量大，下面来简化。

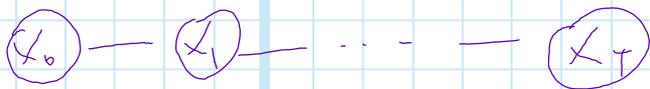$$\hat{\Xi} \; M'(x_0) = \prod_{u \in N_0} M_{u \to 0}(x_0)$$

$$M = \frac{M'(x_0)}{M_{0 \to u}(x_0)} \quad u \in N_0/\{t\}$$

$$M_0(x_0) \propto \psi_0(x_0) \cdot \prod_{u \in N_0} M_{u \to 0}(x_0)$$

$$M_t(x_t) \propto \psi_t(x_t) \cdot M_{0 \to t}(x_t)$$

Chain



$$M_{i_0 \to i_1} = \sum \phi(x_{i_0}, x_{i_1}) \psi_{i_0}(x_{i_0}) M_{i_0-1, i_0}(x_{i_0})$$

$i_0 + 1 = i_1$

$$M(x_i) = \psi_i(x_i) M_{i-1 \to i}(x_i) M_{i+1 \to i}(x_i)$$

下周讲带 cycle 的。

Bethe Interpretation
考虑 MRF:

$$P_\theta(S) = \frac{1}{Z(\theta)} \prod_{S \in V} \psi_S(x_S) \cdot \prod_{(S,t) \in E} \phi_{st}(x_s, x_t)$$

$$\forall s, t, \quad |\mathcal{X}_s| \text{ finit space.}$$

下面用 exp family 来表示.

— Index trick

$$\psi_S(x_S) \quad x_S \in \{0, \cdots, m_{s-1}\}$$

$$= exp\left(\log \psi_S(x_S)\right)$$

$$= exp\left(\sum_{k \in \mathcal{X}_S} \mathbb{1}(x_S = k) \log \psi_S(k)\right)$$

$$= exp\left(\sum_{i \in \mathcal{X}_S} \theta_s^i \mathbb{1}(x_S = k)\right)$$

同理:

$$\phi_{st}(x_s, x_t) = exp\left(\sum_{i \in \mathcal{X}_s} \sum_{j \in \mathcal{X}_t} \theta_{st}^{ij} \mathbb{1}(x_s = i) \mathbb{1}(x_t = j)\right)$$

∴ 用 index trick 来把 概率分布写成指数
和的形式:

$$P_\theta(x) = \frac{1}{Z(\theta)} exp\left(\sum_{S \in V} \sum_{i \in \mathcal{X}_s} \theta_s^i \mathbb{1}(x_s = i) + \sum_{(S,t) \in E} \sum_{i \in \mathcal{X}_s} \sum_{j \in \mathcal{X}_t} \theta_{st}^{ij} \mathbb{1}(x_s = i) \mathbb{1}(x_t = j)\right)$$

其中 $\theta = \{ (\theta_s)(\theta_{st}) \}$ 是 Canonical parameters

$\mu_s \in R^{|X_s|}$ is mean parameter.

$\mu_{st} \in R^{|X_s| \cdot |X_t|}$

Inference Problem:

$\theta$ are given

mean parameters should be get.

$\mu_s$ and $\mu_{st}$.

— Global consist 记作 $M(G)$

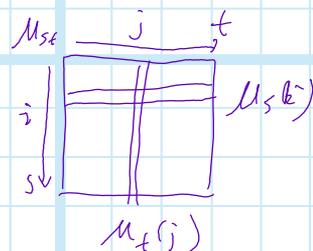$\mu_s, \mu_{st}$ are constist with some drawn distr. It's hard to calc, Instead we use relaxed local const.

— Local consitancy 记作 $L(G)$

$\mu_s(x_s) = P(x_s = i)$

$\mu_{st}(i,j) = P(x_s = i, x_t = j)$

观察: ① $\sum\limits_{i \in X_s} \mu_s(i) = 1$

② $\sum\limits_{j \in X_t} \mu_{st}(i,j) = \mu_s(i) \quad \forall i \in X_s$

③ $\sum_{i \in X_s} \mu_{st}(i,j) = \mu_t(j) \quad \forall j \in X_t$

$\sum_{i,j} \mu_{st}(i,j) = 1$ (这个已经被第1个表示了,i不用写)

$\mu = \{(\mu_s)_{s \in V}, (\mu_{st})_{st \in E}\}$  $\mu$ 表示所在 solution

那么 global 和 local consistency 的关系是:

$$\mu \in M(G) \Rightarrow \mu \in L(G)$$

$$M(G) \subseteq L(G)$$

$$\boxed{\begin{array}{l} \text{if } G \text{ is a tree} \\ \text{then } M(G) = L(G) \end{array}}$$

如何把 $\mu$ 和 tree-structured model 联系起来

Tree-structured

$$P(x) = P_r(x_r) \prod_{v \in V_r} P_v(x_v | x_{\pi(v)}) \nearrow \frac{P(x_v, x_{\pi(v)})}{P(x_{\pi(v)})}$$

$$= \mu_r(x_r) \prod_{s \in V_r} \frac{\mu_{\pi(s),s}(x_{\pi(s)}, s)}{\mu_{\pi(s)}(x_{\pi(s)})}$$

$$= \mu_r(x_r) \prod_{s \in V_r} \frac{\mu_{\pi(s),s}(x_{\pi(s)}, s)}{\mu_{\pi(s)}(x_{\pi(s)}) \mu_s(x_s)} \underbrace{\mu_s(x_s)}_{\text{放到前面去}}$$

$$= \prod_{v \in V} \mu_v(x_v) \cdot \prod_{(s,t) \in E} \frac{\mu_{st}(x_s, x_t)}{\mu_s(x_s) \cdot \mu_t(x_t)} \qquad (*)$$

以上说明了如何用 mean parameter 来表示 free-structured node

只要我们有一个 locally consistency $\mu$, tree-struc model 才就能写成这种形式

下面看 exp family

$$\hat{\mu} = \arg\max_{\mu \in M(G)} \{\theta^T \mu - A^*(\mu)\} \quad \xleftarrow{\text{convex conject of}} \text{log parcul function}$$

$\longrightarrow$ all set of realizable mean

$$\boxed{\hat{\mu} = \arg\max_{\mu \in M(G)} \{\theta^T \mu + H(\mu)\}}$$

两个值先: $M(G)$: realizable, global constancy are the same thing.

非常难计算，甚至验证是都很难，

值差区的是, tree-structure $M(G) = L(G)$

$H(\mu) =$ entropy

由信息论:

$$H(\mu) = -\sum_x p \log P$$

代入 $(*)$ 得:

$$= -\sum_x P(x) \cdot \left\{ \sum_{v \in V} \log \mu_v(x_v) + \sum_{(s,t) \in E} \log \frac{\mu_{st}(x_s, x_t)}{\mu_s(x_s) \mu_t(x_t)} \right\}$$

$$H(\mu) = \sum_{v \in V} H_v(\mu_v) - \sum_{(s,t) \in E} I_{st}(\mu_{st})$$

only applies to free-structured models.

entropy          mutual information

$$H_s(\mu_s) = -\sum_{x \in X_s} \mu_s(x) \log \mu_s(x)$$

$$I_{st}(\mu_{st}) = \sum_{(x_s, x_t) \in X_s \times X_t} \mu_{st}(x_s, x_t) \log \frac{\mu_{st}(x_s, x_t)}{\mu_s(x_s) \cdot \mu_t(x_t)}$$

由上, decomposed the entropy along the tree
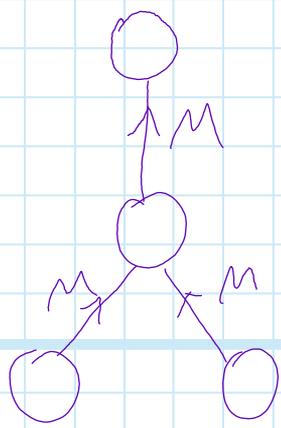
为什么要研究 tree-structured model?

好处:

① $M(G) = L(G)$

② Entropy can be factorize into several terms which only contains one or two terms.

这样, tree-structured model 间接解决了 $M(G)$ 和 $H$ 的问题.

那么能否推广到任意的图?

可以用① approximation.

Belief propagation

但实际写程序的时候，直接把这个算法扔到任意图上去跑就行了。效果还可以。后来人们开始找理论支持。

For Loopy-structured

$$M(G) \approx L(G)$$

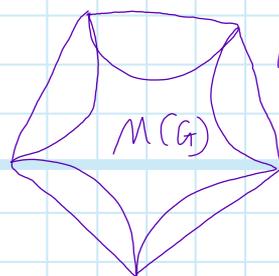$$H_{Be}(\mu) \approx \sum_s H_s(\mu_s) - \sum_{(s,t)\in E} I_{st}(\mu_{st})$$

Bethe Entropy

$$\hat{\mu} = \arg\max_{\mu \in L(G)} \left\{ \theta^T \mu + H_{Be}(\mu) \right\}$$
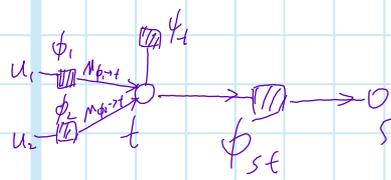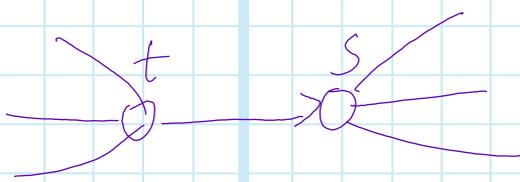
Bethe Variational

$\Downarrow$ optimization

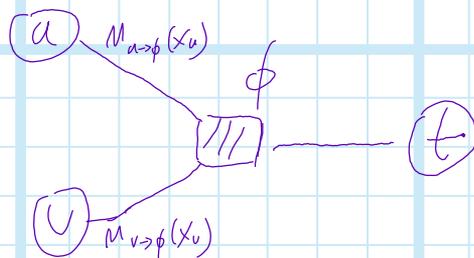Loopy Belief Propagation
(LBP)



L(G): convex bound.

M(G)

Loop Belief Propagation 也可以扩展到不止两个 factor

$$M_{t \to s}(x_s) = \sum_{x_t} \phi(x_s, x_t) \psi_t(x_t) \prod_{u \in N_t / \{s\}} M_{u \to t}(x_t)$$

$$= \sum_{x_t} \phi(x_s, x_t) M_{t \to \phi}(x_t)$$

$$M_{t \to \phi}(x_t) = M_{\psi_t \to t}(x_t) \cdot M_{\phi_1 \to t}(x_t) M_{\phi_2 \to t}(x_t)$$

$$M_{\phi \to s}(x_s) = \sum_{x_t} \phi(x_s, x_t) M_{t \to \phi}(x_t)$$

以上 $\frac{1}{2}$ binary factor, 下面 $\frac{1}{3}$ 3 factors



$$M_{\phi \to t}(x_t) = \sum_{x_u} \sum_{x_v} \phi(x_u, x_v, x_t) M_{u \to \phi}(x_u) M_{v \to \phi}(x_v)$$

下面继续 Inference 的问题.

已知 $\theta$ 求 $M$.

① $L(G) \approx M(G)$ 把它固起图②.

② New approximation of the Entropy.

即为: $\hat{\mu} = \underset{\mu \in M(G)}{\arg\max} \{ \theta^T \mu - A^*(\mu) \}$

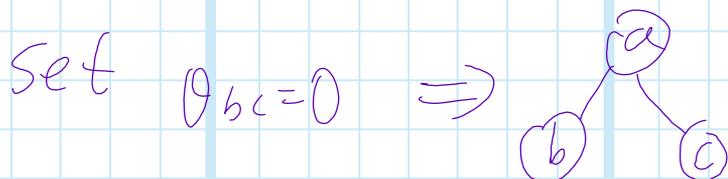$$A^*(\mu) = \sup_\theta \{\theta^T\mu - A(\theta)\},$$

also intractable.

$A(\theta)$ is easy to compute for tree-structured models.

对其意图，可以分解为另个树的组合。



$$\exp( \theta_a \cdot \mathbb{1}(x_a) + \theta_b \mathbb{1}(x_b) + \theta_c \cdot \mathbb{1}(x_c) +$$
$$\theta_{ab}\mathbb{1}(x_a,x_b) + \theta_{ac}\mathbb{1}(x_a,x_c) + \theta_{bc}\mathbb{1}(x_b,x_c)$$

set $\theta_{bc}=0 \Rightarrow$



project the parameters into some subspace

it will becomes to tree.

① projection $\begin{cases} \theta \xrightarrow{\theta_{bc}=0} \theta' \\ \theta \xrightarrow{\theta_{ac}=0} \theta'' \end{cases}$

② combination: $\alpha\theta' + (1-\alpha)\theta'' = \theta$

$$A(\theta) = A(\alpha \cdot \theta' + (1-\alpha)\theta'')$$
$$\leq \alpha \underbrace{A(\theta')}_{easy} + (1-\alpha)\underbrace{A(\theta'')}_{easy} \qquad \because A \text{ is convex.}$$

在找出 θ 的上组合的地址去估计 $A(\theta)$

上节#] :
Inference on exp family

$$exp(\theta^T \phi(x) - A(\theta))$$

图 样 : $\theta \to \mu$ .

法① $\mu = E_\theta[\phi(x)]$

法② $A(\theta) = \sup_\mu \{\theta^T \mu - A^*(\mu)\}$

$$= \sup_\mu \{\theta^T \mu + H(\mu)\}$$

estimate $A(\theta)$, is related to

$$\theta^T \mu + H(\mu)$$

$A(\theta)$ for tree-structure is tractable,
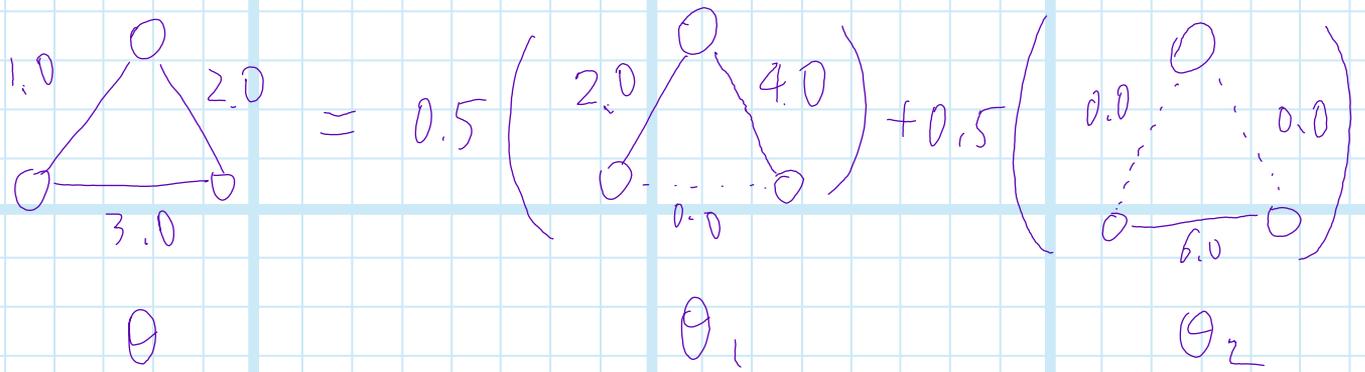
We can decompose non-tree to tree.

A is convex.

$$\theta = \alpha_1 \theta_1 + \alpha_2 \theta_2 \quad (\alpha_1 + \alpha_2 = 1)$$

$$A(\theta) \leq \alpha_1 A(\theta_1) + \alpha_2 A(\theta_2)$$

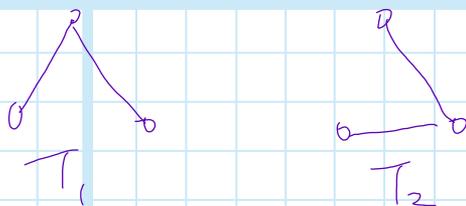通过分解 $\theta$ , 让其变式 tree-structured.

例:



$$\theta = 0.5\,\theta_1 + 0.5\,\theta_2.$$

通过不同的分解 能得到 不同的 upper bound.
能否找到一个最小的上界?

Find the best upper bound.



$$\theta(T_1) = \begin{pmatrix} 0.2 \\ 0.4 \\ 0.0 \end{pmatrix} \qquad \theta(T_2) = \begin{pmatrix} 0.0 \\ 6.0 \\ 6.0 \end{pmatrix}$$

$$\sum p(T)\,\theta(T) = E_p\big[\theta(T)\big]$$

$$\min_{T \in \Sigma} \sum p(T)\,A(\theta(T)) \longrightarrow E_p\big[A(\theta(T))\big]$$

$$\text{s.t.} \sum_T p(T)\,\theta(T) = \bar{\theta} \checkmark \quad \text{target parameter}$$

给定的 $\theta$, 先 各个 T 的
$\theta$ 组合起来应该等于 $\bar{\theta}$

consistent constraint 等价 $E_p(\theta(T))$

用拉格朗日求解：

$$L(\theta,\mu) = E_p[A(\theta(T))] + \langle \mu, \bar{\theta} - E_p(\theta(T)) \rangle$$

$\mu$ 和 $\bar{\theta}$ 同维度规模

$$= \mu^T\bar{\theta} + E_p[A(\theta(T)) - \mu^T\theta(T)]$$

$$= \mu^T\bar{\theta} + \sum P(T)\left(A(\theta(T)) - \mu^T\theta(T)\right)$$

$$\frac{\partial L(\theta,\mu)}{\partial \theta(T)} = P(T)\left[\nabla A(\theta(T)) - \mu\right] = 0$$

$$E_{\theta(T)}[\phi_\alpha] = \mu_\alpha$$

$\alpha$ 表示 和某个材料相关
的参数。对某个材料，
不是所有参数都要求导。

例：

$$\theta(T_1) \quad \overset{a}{\underset{b \quad c}{\bigwedge}} \quad \rightarrow \alpha_1 = \begin{pmatrix} ab \\ ac \end{pmatrix} \qquad \theta(T_2) \quad \overset{a}{\underset{b \quad c}{\diagdown}} \quad \rightarrow \alpha_2 = \begin{pmatrix} ac \\ bc \end{pmatrix}$$

$$\begin{bmatrix} \mu_{ab} \\ \mu_{ac} \\ \mu_{bc} \end{bmatrix} \quad \begin{array}{l} \rightarrow = E_{\theta(T_1)}[\phi_{\alpha_1}(x)] \\ \\ \rightarrow = E_{\theta(T_2)}[\phi_{\alpha_2}(x)] \end{array}$$

In order to calc $\mu$, we have to move
to the dual problem.

$$A^*(\Pi_T(\hat{\mu})) = \langle \hat{\theta}(T), \hat{\mu} \rangle - A(\hat{\theta}(T))$$

$$L(\hat{\theta}, \hat{\mu}) = \hat{\mu}^T \bar{\theta} + E_p[-A^*(\Pi_T(\hat{\mu}))]$$

$$= \hat{\mu}\,\bar{\theta} - E_p[A^*(\Pi_T(\hat{\mu}))]$$

$$= \hat{\mu}\,\bar{\theta} + E_p[H_T(\Pi_T(\hat{\mu}))]$$

↑ projection of entropy to
each tree.

$$\sum_{s \in V} H_s(\mu_s) - \sum_{(s,t) \in E(T)} I_{st}(\mu_{st})$$

等于每个 tree 的 entropy 之和
减去互信息.

$$E_p[H_T(\Pi_T(\hat{\mu}))] = E_p\left[\underbrace{\sum_{s \in V} H_s(\mu_s)}_{\text{对所有 tree 都相同, 移出去.}} - \sum_{(s,t) \in E(T)} I_{st}(\mu_{st})\right]$$

$$= \sum_{s \in V} H_s(\mu_s) - \sum_{(s,t) \in E} \boxed{\rho_{st}}\, I_{st}(\mu_{st})$$

here don't use $E(T)$, it mean for all
pairs in the graph.
edge appearance probability

$$= H_{Taw}(\mu)$$

$$H_{se}(\mu) = \sum_{s \in V} H_s(\mu_s) - \sum_{(s,t) \in E} I_{st}(\mu_{st})$$

假设 edge appearance probability.

 $= 0.4$ $\bigwedge$ $+ 0.2$ $\angle$ $+ 0.4$ $>$

$$P_{ab} = 0.6$$
$$P_{bc} = 0.6$$
$$P_{ac} = 0.8$$

总结:

max $\mu^T \theta + H_{approx}(\mu)$ is a common
way to do inference.

另一种方法 就是 Variational Inference.

Learning
 — Parameter estimation.
      Data → Parameter
      方法: variational estimation
          variational inference

Space of distributions.
    in typical space, we have distance.
    how to measure the distance between
    two distribution?
       Use K-L Divergence.
       two distribution p q.

$$D_{KL}(p\|q) = E_p\left(\log\frac{P(x)}{q(x)}\right)$$

    some basic property:

① not symmetric

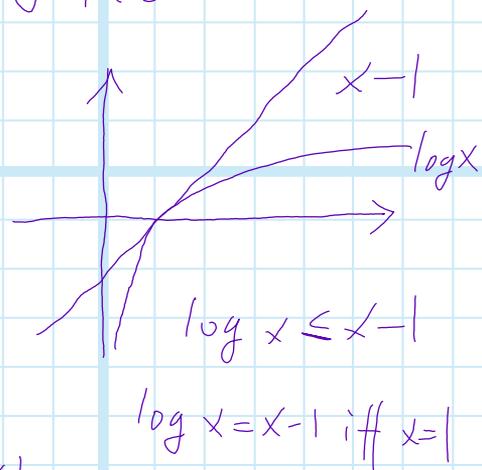$$D_{KL}(p\|q) \neq D_{KL}(q\|p)$$

② non-negative
$$D_{KL}(p\|q) \geq 0$$

proof:

$$-D_{KL}(p\|q) = E_p\left[\log\frac{q(x)}{p(x)}\right]$$

$$= \int P(x)\cdot\log\frac{q(x)}{P(x)}\mu(dx)$$
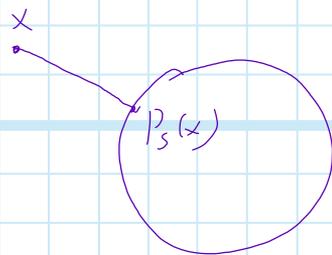
$$\leq \int P(x)\left(\frac{q(x)}{P(x)}-1\right)\mu(dx)$$

$$= \int q(x)\mu(dx) - \int P(x)\mu(dx) = 0$$

$$D(P\|q) = 0 \quad only \quad when \quad \frac{P(x)}{q(x)}=1 \quad i.e. \quad P(x)=q(x)$$

$x-1$

$\log x$

$\log x \leq x-1$

$\log x = x-1 \ iff \ x=1$

有了基代 Distance 后，就 可以引入 projection

projection

$$P_s(x) = \underset{x'\in S}{\arg\min}\ d(x',x)$$

$x$

$P_s(x)$

$S:$ convex set

可以对 distribution 做同样的定义, 把某个分布投影到某个分布族上.

∵ KL Divergence 非对称, ∴ 有 两种投影方法.

设有一个分布族 P

$$I_{proj Q}(p) = \underset{q\in Q}{\arg\min}\ D_{KL}(q\|p) \qquad information\ projection$$

$$M_{proj Q}(p) = \underset{q\in Q}{\arg\min}\ D_K(p\|q) \qquad moment\ projection$$

# 下面图形 Model Estimation

$P_\theta(x)$

已知 $D = \{x_1, \cdots, x_n\}$ $\xrightarrow{\text{estimate}}$ $\theta$

basic idea: max likelihood Estimation.



$P_1$ is good.

Likelihood

likelihood vs. density

$P(x; \theta)$ $\nearrow$ $P_\theta(x)$ density

$\searrow$ $L_x(\theta) = P_\theta(x)$ likelihood

when we talk about density, we assume we already know $\theta$

$$L_D(\theta) = \prod_{i=1}^{n} L_{x_i}(\theta)$$

$$= \prod_{i=1}^{n} P_\theta(x_i)$$

This formula implicitly assume samples are independent.

注用乘法会 overflow, 改用 log likelihood

log - likelihood.
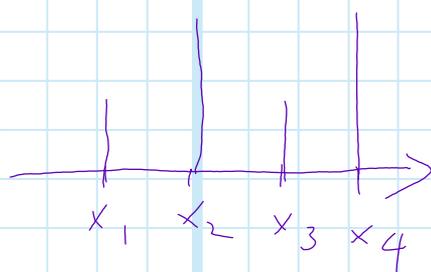
$$\log L_D(\theta) = \sum_{i=1}^{n} \log P_\theta(x_i)$$

Max log $l_D(\theta)$ is MLE.

下面给出几何解释

Emprical Distribution

Given $D = \{x_1, \cdots, x_n\}$

$$\tilde{P}_D(x) = \frac{1}{n} \sum_{i=1}^{n} \delta_{x_i}(x)$$



$$E_{\tilde{P}_D}(f) = \frac{1}{n} \sum_{i=1}^{n} f(x_i)$$

$$D_{KL}(\tilde{P}_D \| P_\theta) = E_{\tilde{P}_D}\left[ \log \frac{\tilde{P}_D}{P_\theta} \right]$$

$$= \underbrace{E_{\tilde{P}_D}[\log \tilde{P}_D]}_{\text{与参数无关.}} - E_{\tilde{P}_D}[\log P_\theta]$$

$$= \text{constant} - \underbrace{\frac{1}{n} \sum_{i=1}^{n} \log P_\theta(x_i)}_{\text{object of MLE}}$$

∴ maximise the likelihood $\iff$

minimize the KL divergence.

$$P_{\hat{\theta}} = \underset{P_\theta \in P}{\text{argmin}} \; D_{KL}(\hat{P}_\theta \| P_\theta) = M\text{proj}_P(\tilde{P}_D)$$

MLE 的几何意义 我 是 一个 经验分布。
M-projection 距离意下, 向 分布族 P 投射
当 这个 分布 是 指数分布族 时:
exp family:

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

maximize

$$E_{\tilde{P}_D}[\log P_\theta(x)]$$

$$\log L_D(\theta) = E_{\tilde{P}_D}(\theta^T \phi(x) - A(\theta))$$

$$= \theta^T \left(\frac{1}{n} \sum_{i=1}^{n} \phi(x_i)\right) - A(\theta)$$

$$= \theta^T \tilde{\mu}_D - A(\theta) \quad \text{with } \tilde{\mu}_D = E_{\tilde{P}_D}[\phi(x)]$$

$$\nabla_\theta \log L_D(\theta) = \tilde{\mu}_D - \nabla_\theta A(\theta) = 0$$

$$\nabla_\theta A(\theta) = E_{\tilde{P}_D}[\phi(x)]$$

$$\parallel$$

$$E_\theta[\phi(x)] =$$

M projection

try to align the
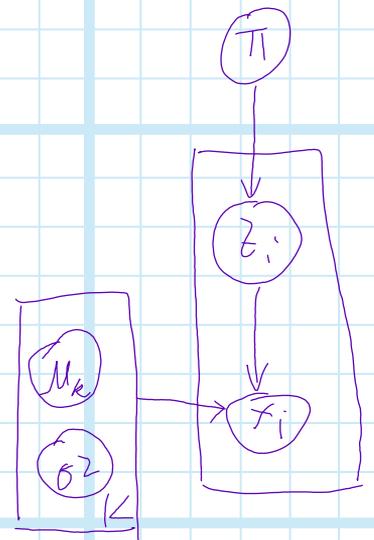moment of the data
and parameter.

下面用高斯混合来作了例。

GMM: $\qquad$ $\sigma$ 已知。

$$p(x_i, z_i | \pi, \{\mu_k\})$$

$$= \pi(z_i) \cdot \frac{1}{\sqrt{2\pi\sigma^2}} \exp\left(-\frac{(x-\mu_{z_i})^2}{2\sigma^2}\right)$$

$$= \exp\left(\log\pi(z_i) - \frac{(x-\mu_{z_i})^2}{2\sigma^2}\right)$$

$$= \exp\left(\sum_{k=1}^{K} \delta_k(z_i) \cdot \log\pi(k) - \sum_{k=1}^{K} \delta_k(z_i)\frac{(x-\mu_k)^2}{2\sigma^2}\right) \quad \binom{\text{index}}{\text{trick}}$$

$$D = \sum_{i=1}^{n}\left(\sum_{k=1}^{K}\delta_k(z_i)\log\pi_k - \sum_{k=1}^{K}\delta_k(z_i)\frac{(x-\mu_k)^2}{2\sigma^2}\right)$$

求 $\pi_k$ 和 $\mu_k$，分别成两个问题。

1. Solve $\pi_k$.

$$J(\pi) = \sum_{i=1}^{n}\sum_{k=1}^{K}\delta_k(z_i)\log\pi_k$$

$$= \sum_{k=1}^{K} n_k \log\pi_k \qquad n_k = \#\{i : z_i = k\}$$

$$s.t. \quad \sum_k \pi_k = 1 \quad \pi_k \geq 0 \quad \forall k = 1, \cdots, K$$

$$\Rightarrow \pi_k \propto n_k = \frac{n_k}{\sum_{l=1}^{K} n_l} = \frac{n_k}{n} \qquad \text{投标统归来了.}$$

2. Solve $\mu_k$

$$I(\mu_k) \quad \leftarrow \text{minimize}$$

$$= \sum_{i \in S_k} \frac{(x_i - \mu_k)^2}{2\sigma^2}$$

$$\mu_k = \frac{\sum_{i \in S_k} x_i}{|S_k|}$$

找一个点到其它的
距离之和最短.
找 mean.

这里我们假设了 $z_i$ 是已知的. $z_i$ 未知的
话不能这么解. 要用 EM. 它我们 partial
observed model.

$$\underset{\underset{\text{observed}}{\uparrow}\quad\underset{\text{latent}}{\uparrow}}{P_\theta(x,z)}$$

$$= g(x) h_x(z) \exp\left( \theta^T \phi(x,z) - A(\theta) \right)$$

$x$ is observed    $z$ is unknown.

$$P(z|x) = \frac{P(x,z)}{P(x)} = \frac{h_x(z) \exp(\theta^T \phi(x,z) - A(\theta))}{\int_z h_x(z) \exp(\theta^T \phi(x,z) - A(\theta)) \, dz}$$

这个条件分布仍型是 exp family. 可以写成:

$$h_x(z) \exp(\theta^T \phi(x) - A(\theta|x)) \text{ 的形式}.$$

$$P(z|z) = \frac{h_x(z) \exp(\theta^T \phi(x, z))}{\int_z h_x(z) \exp(\theta^T \phi(x, z)) dz}$$

$$\therefore A(\theta|x) = \log \int_z \exp(\theta^T \phi(x, z)) h_x(z) dz$$

called conditional log-partition function

$$P(x) = \int h_x(z) \exp(\theta^T \phi(x, z) - A(\theta)) dz$$

$$= \frac{1}{\exp(A(\theta))} \int h_x(z) \exp(\theta^T \phi(x, z)) dz$$

$$= \frac{\exp(A(\theta|x))}{\exp(A(\theta))} = \exp(A(\theta|x) - A(\theta))$$

$$\log P(x) = A(\theta|x) - A(\theta)$$

这样我们写成了和隐变量 $z$ 无关
的形式.

$$P_\theta(x, z) = g(x) h_\theta(z) \exp(\theta^T \phi(x, z) - A(\theta))$$

$$\log L(\theta | x) = A(\theta | x) - A(\theta)$$
$$P(x | \theta)$$

目标: $\max \sum_i \log P(x_i | \theta)$

但是 $A(\theta | x) = \log \int_Z \exp(\theta^T \phi(x, z)) h_x(z) dz$

可能 z 空间很大, 可以用 EM 来求解

今天用 EM 去解这项,

EM : find lower bound of $\log L(\theta | x)$

$\log(\theta | x)$ is difficult to compute, we
can max its lower bound. which is
eaiser to compute.


$A(\theta | x)$: conditional log partition function

$A(\theta) = \sup_\mu \{ \theta^T \mu - A^*(\mu) \}$ duality between
parameter domain and
canonical mean domain

$A(\theta | x) = \sup_\mu \{ \theta^T \mu - A^*(\mu | x) \}$
$\underset{E_{P(z | x)}[\phi(x, z)]}{\downarrow}$

$\underset{\text{lower bound}}{\Downarrow}$

$A(\theta | x) \geq \theta^T \mu - A^*(\mu | x) \quad \forall \mu$ . $A^*(\mu | x)$ is conditional
entropy

$$\log L(\theta|x) \geq \theta^T \mu - A^*(\mu|x) - A(\theta) = Q_x(\theta, \mu)$$

得到了 $\log L(\theta|x)$ 的下界. 当 $\mu$ 取到 $\arg\max$ 时, 等号成立

左边只有一个参数, 右边有两个参数 $\mu, \theta$,

但是可以保证任取 $\theta$ 都使不等式成立.

用 coodinate asend 来优化 $Q_x(\theta, \mu)$

Coodinate assent

$\max f(x, y)$

  fix $x^{(t)}$

  $y^{(t+1)} \leftarrow \arg\max_{y} f(x^{(t)}, y)$

  $x^{(t+1)} \leftarrow \arg\max_{x} f(x, y^{(t+1)})$

回到之前那个 $Q_x(\theta, \mu)$, EM 算法,

分为 E-step and M-step.

  $\hat{\mu} \leftarrow \arg\max_{\mu} Q(\theta; \mu)$

   $= \arg\max_{\mu} \theta^T \mu - A^*(\mu|x)$

    $\uparrow \hat{\mu} = E_{\theta}[\phi(x, z)]$

    实质上就是在计算 expection

M-tep:

  $\hat{\theta} \leftarrow \arg\max_{\theta} Q(\theta; \mu)$

   $= \arg\max_{\theta} \theta^T \mu - A(\theta)$

这里 $\theta$ 和 $\mu$ 的形式都一样, 但是计算

可以解释为计算 mean, 然后 do the model estimation

下面看为什么 EM 能 work.

$$\log L_x(\theta^{(t)})$$
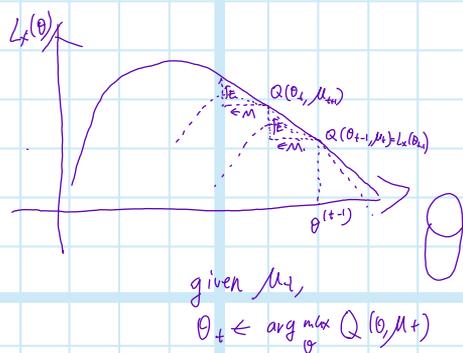$\hookleftarrow$ x is constant

$$\mu^{(t+1)} = \arg\max_{\mu} Q(\theta^{(t)}, \mu)$$

$$Q(\theta^{(t)}, \mu^{(t+1)}) = L_x(\theta^{(t)}) \quad \mu 取到最优, 等号到达.$$

$$\geq Q(\theta^{(t)}, \mu^{(t)})$$

$$\geq Q(\theta^{(t-1)}, \mu^{(t)}) = L_x(\theta^{(t-1)})$$



$L^{(\theta)}$

$Q(\theta_t, \mu_{t+1})$

$Q(\theta_{t-1}, \mu_t) = L_x(\theta_t)$

$\theta^{(t-1)}$

$\theta$

given $\mu_t$,
$\theta_t \leftarrow \arg\max_{\theta} Q(\theta, \mu_t)$

E: close the gap ($\mu, \theta$ are coupled)

M: move to another procedure.

EM in distribution space.

try to get the KL divergence
between two distr. via which
one is parameterize by mean parameter
$\mu$

another is parameter by canonical
parameter $\theta$

$$KL_x(\mu \| \theta)$$

$$KL_x(\mu \| \theta)$$

$$= A(\theta | x) + A^*(\mu | x) - \theta^T \mu$$

$$Q(\theta; \mu)$$

$$= \mu^T \theta - A^*(\mu | x) - A(\theta)$$

$$L(\theta|x) - Q(\theta, \mu)$$

$$= (A(\theta|x) - A(\theta)] - [\mu^T\theta - A^*(\mu|x) - A(\theta)]$$

$$= A(\theta|x) + A^*(\mu|x) - \mu^T\theta = KL_x(\mu||\theta)$$

这个 KL距离就是那个 gap.

以上 只考虑了一个 sample $x$, 对于一个 sample
集合 $D$

$$L_D(\theta) = \sum_i (A(\theta|x_i) - A(\theta)]$$

$$= \sum_i A(\theta|x_i) - nA(\theta)$$

$$= \sum_i (\mu_i^T\theta - A^*(\mu_i|x_i)) - nA(\theta)$$

$$= \sum_i (\mu_i^T\theta - A^*(\mu_i|x_i) - A(\theta))$$

注意对于不同的 sample 有相同的 parameter $\theta$
和各自的 mean $\mu$.

∴

E - Step:

$$\hat{\mu}_i \leftarrow \arg\max_{\mu} \mu^T \theta - A^*(\mu | x_i)$$

respetively for $i = 1, \cdots, n.$

M - step

$$\hat{\theta} = \arg\max_{\theta} \{ \theta^T \bar{\mu} - A(\theta) \} \qquad \bar{\mu} = \frac{1}{n} \sum_{i=1}^{n} \mu_i$$

与单个的 EM 的区别：上免自锁，M 取平均

Variational Inference
very closely related to partially
observed model, so us EM.

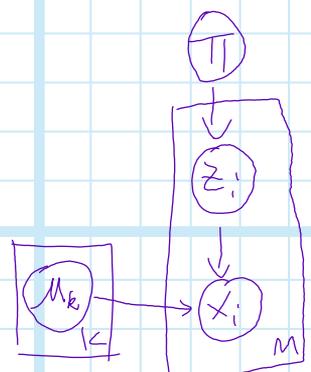$$P(x, z) = g(x) h_z(z) \exp(\theta^T \phi(x, z) - A(\theta))$$

E - M :

    origional : $L_x(\theta)$

    alternative : $\geqslant Q(\theta, \mu) = \mu^T \theta - A^*(\mu|x) - A(\theta)$

    E - step : update $\mu$ : $\hat{\mu} \leftarrow \underset{\mu}{\arg\max} \{\mu^T \theta - A^*(\mu|x)\}$

    M - step : update $\theta$ : $\hat{\theta} \leftarrow \underset{\theta}{\arg\max} \{\mu^T \theta - A(\theta)\}$

下面看 EM 如何在图模型上应用.

Gaussian Mixture Model.

$$P(x_i, z_i | \{\mu_k\}, \pi) \propto \pi(z_i) \exp\left(-\frac{(x_i - \mu_{z_i})^2}{2\sigma^2}\right)$$

$$= \exp\left(\sum_{k=1}^{K} 1_k(z_i) \log \pi_k - \sum_{k=1}^{K} 1_k(z_i) \frac{(x - \mu_k)^2}{2\sigma^2}\right)$$



$x_i$ : observed   $\sigma^2$: known
$z_i$ : latent.

$$= \exp\left(\sum_{k=1}^{K} 1_k(z_i) \log \pi_k + \sum_{k=1}^{K} 1_k(z_i) \frac{x_i \mu_k}{\sigma^2} - \sum_{k=1}^{K} 1_k(z_i) \frac{x_i^2}{2\sigma^2} + \dots\right)$$

写成指数形式后, 充分统计量有 $1_k(z_i), 1_k(z_i)x_i, 1_k(z_i)x_i^2$

E-step: infer the expection

$$E[1_k(z_i)] = q_i^k = P_r(z_i = k)$$

$$E[1_k(z_i) x_i] = q_i^k \cdot x_i$$

$$E(1_k(z_i)x_i^2) = q_i^k \cdot x_i^2 \dots$$

∴ 只需计算 $q_i^k$

$$q_i^k = P_r(z_i = k) \propto \pi_k \cdot \exp\left(-\frac{(x_i - \mu_k)^2}{2\sigma^2}\right)$$

$$\uparrow \quad E\text{-step}$$

M-step

$q_i^k$ 表似于一个 类型的 soft-assignment.

$$\mu^T \theta - A(\theta)$$

$$= \sum_{i=1}^{n} \sum_{k=1}^{K} q_i^k \log \pi_k - A(\pi)$$

$$s.t, \quad \pi \in S_K$$

$$\pi_k \propto \sum_{i=1}^{n} q_i^k$$

$$\mu_k = \frac{\sum_{i=1}^{n} q_i^k x_i}{\sum_{i=1}^{n} q_i^k}$$

In variational inference, $A^*(\mu|x)$ is very hard to compute. We relay on Entroy Approximation of it.
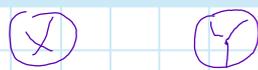
$$A^*(\mu|x)$$

$$\downarrow$$

$$A_f^*(\mu|x)$$

$f =$ factorized variational distribution.



$p(x,y)$     ⇑ approximate.

$$p(x,y) \approx q_x(x) \, q_y(y)$$

$q \leftarrow$ parameter: $\lambda$   $q_\lambda$ 所控制的范族.

$$\hat{\lambda} \leftarrow \arg\min_\lambda D_{KL}(q_\lambda \| p)$$
估计

因此 MLG 中, 有 $\arg\min D_{KL}(\tilde{P}_p \| P_\theta)$
估计

用右边估计左边, 叫做 M-projection

而这里用左边估计右边, 这是 I-projection
call mean field approximate.

下面看一个例子.

Hidden MRF

$p(x_1, x_2)$ prior on latent space

$x_1, x_2 \in \{0, 1\}$

会用到很多 tricks.


← latent
← observed

$$p(x_1, x_2) \, p(y_1 | x_1) \, p(y_2 | x_2)$$

$$\propto \exp\Big( \psi_1(x_1) + \psi_2(x_2) + \phi(x_1, x_2) \Big) \exp\big( f(x_1, y_1) \big) \exp\big( f(x_2, y_2) \big)$$

$$= \exp\Big( \sum_{i \in \{0,1\}} \delta_i(x_1) \cdot \theta_1^i + \sum_{j \in \{0,1\}} \delta_j(x_2) \, \theta_2^j + \sum_{i,j} \theta_{12}^{ij} \delta_i(x_1) \delta_j(x_2)$$

$$+ \sum_i f_1^i \delta_i(x_1) + \sum_j f_2^j \delta_j(x_2) \Big)$$

$$\theta_1^i = \psi_1(i) \quad i \in \{0, 1\} \qquad \theta_2^j = \psi_2(j) \quad j \in \{0, 1\}$$

$$f_1^i = f(i, y_1) \quad i \in \{0, 1\} \qquad \theta_{12}^{ij} = \phi(i, j) \qquad 田$$

| canonical | suff. stats |
|---|---|
| $\delta_i(x_1)$ | $\theta_1^i + f_1^i$ |
| $\delta_j(x_2)$ | $\theta_2^j + f_2^j$ |
| $\delta_i(x_1)\delta_j(x_2)$ | $\theta_{12}^{ij}$ |

$E[\delta_i(x_1)] = P_r(x_1 = i)$

$E[\delta_j(x_2)] = P_r(x_2 = j)$

$E[\delta_i(x_1)\delta_j(x_2)] = P_r(x_1 = i, x_2 = j)$

∴ 可以变成三项.

$$= exp\left(\sum_{i \in \{0,1\}} \delta_i(x_i)\cdot\theta_1^i + \sum_{j \in \{0,1\}} \delta_j(x_j)\theta_2^j + \sum_{i,j}\theta_{22}^{ij}\delta_i(x_1)\delta_j(x_2)\right)$$

就是 Ising Model.

在 E-step 中需要计算 $E[\delta_i(x_i)]$  $E[\delta_j(x_2)]$

$E[\delta_i(x_i)\delta_j(x_2)]$ 很麻烦. 可以用

Variational inference 来计算.

问题:

$$p(x_1,x_2) = exp\left(\sum_i \theta_1^i \delta_i(x_1) + \sum_j \theta_2^j \delta_j(x_2) + \right.$$

$$\left. \sum_{ij} \theta_{12}^{ij} \delta_i(x_1) \delta_j(x_2)\right)$$

用 Mean field Approx.    ↗ approx.

$$q(x_1,x_2) = q_1(x_1) q_2(x_2)$$

$$\hat{q} = \arg\min_{q} D_{KL}(q \| p)$$

$$D_{KL}(q \| p)$$

$$= E_q\left[\log\frac{q}{p}\right]$$

$$= E_q\left[\log q - \log p\right)$$

$$= E_{q_1, q_2}\left[\log q_1 + \log q_2 - \log p\right]$$

用 coordinate descent

fix $q_2$, update $q_1$. 先去掉和 $q_1$ 无关的项

$$E_{q_1}\left[\log q_1 - \log p\right)$$

注意到:

$$E_{q_1 q_2}\left[\log p\right]$$

$$= E_{q_1 q_2}\left[\sum_i \theta_1^i \delta_i(x_1) + \sum_j \theta_2^j \delta_j(x_2) + \right.$$

$$\left. \sum_{ij} \theta_{12}^{ij} \delta_i(x_1)\delta_j(x_2)\right] \longrightarrow \text{this is the key to simplify.}$$

$$= \sum_i \theta_1^i q_1^i + \sum_j \theta_2^j q_2^j + \sum_{ij} \theta_{12}^{ij} q_1^i q_2^j$$

化成了一个很简单的优化(凸)题.

回顾:

① E-M:

$$L(\theta) \geqslant Q(\theta, \mu)$$

$$\Downarrow$$

$$KL(\mu \| \theta) \longleftarrow \begin{matrix} close \\ minimize \\ gap \end{matrix} \rightarrow \begin{matrix} variational \\ inference. \end{matrix}$$

② Ising Model

$$G = (V, E)$$

uniary term ↓        binery term ↓

$$P(x) \propto exp\left( \sum_{v \in V} \theta_v x_v + \sum_{(u,v) \in E} \theta_{uv} x_u x_v - A(\theta) \right)$$

if it's loppy graph.

use Mean-Field Approximation

$$q_\lambda(x) = exp\left( \sum_{v \in V} \lambda_v x_v - B_v(\lambda) \right)$$

use $q_\lambda(x)$ to approximate $P(x)$.

$$D_{KL}(q_\lambda \| P_\theta) = E_q\Big[ \log q - \sum_{v \in V} \theta_v x_v - \sum_{(u,v) \in E} \theta_{uv} x_u x_v$$

$$+ A(\theta) \Big]$$

← not related to optimization

$$= -\sum_{v \in V} H_v(q_v) - \sum_{v \in V} \theta_v E_{q_v}[x_v] - \sum_{(u,v) \in E} \theta_{uv} E_{q_u q_v}[x_u, x_v]$$

$$= -\sum_{v \in V} H_v(q_v) \sum_{v \in V} \theta_v q_v - \sum_{(u,v) \in E} \theta_{u,v} q_u q_v$$

下面看 LDA 的例子.

## Latent Dirichlet Allocation

$$P(\underbrace{\theta_d}_{latten}, \{z_{di}, \underbrace{w_{di}}_{observe}\} | \alpha, \beta)$$

$$Dir: \ P(\theta | \alpha) = \frac{1}{B(\alpha)} \prod_{m=1}^{m} \theta_m^{\alpha_m - 1}$$

$$= exp\left( \sum_{m=1}^{m} (\alpha_m - 1) \log \theta_m - \log B(\alpha) \right)$$

$$\therefore \ P \propto exp\left( \sum_{k=1}^{k} (\alpha_k - 1) \underbrace{\log \theta_d^k}_{①} \right) \quad \leftarrow Dir$$

$$+ \sum_{i=1}^{nd} \sum_{k=1}^{m} \underbrace{1_k(z_{di}) \cdot \log \theta_d^k}_{②}$$

$$+ \sum_{i=1}^{nd} \sum_{k=1}^{M} \underbrace{1_k(z_{di})}_{③} \log \beta_k(w_{di}) \Big)$$

Follow Mean Field Approximation.

$$q(\theta_d, \{z_{di}\}) = \underbrace{q_\gamma(\theta_d)}_{Dir(\gamma)} \prod_{i=1}^{nd} \underbrace{\rho_{di}(z_{di})}_{Categorical \ distr.}$$

① $E_q[\log \theta_d^k] = E_{q_\gamma}[\log \theta_d^k] = \psi(\gamma_d^k) - \psi(1^T \gamma_d)$

$\left( E[\phi(x)] = \nabla_\theta A(\theta) = \psi(\alpha^k) - \psi(\sum_h \alpha^h) \right)$ ↖ digamma function

② $E_q[1_k(z_{di}) \cdot \log \theta_d^k]$ 　　　② 是①和③ 乘起来

$= E_q[1_k(z_{di})] \cdot E_{q_\gamma}[\log \theta_{di}^k]$

③ $E_q[1_k(z_{di})] = P_\gamma(z_{di}=k) \leftarrow$ w.r.t. $\rho_{di}$

$\quad = \rho_{di}(k)$

回顾.

Inference.

$$\theta \longrightarrow E_\theta[f(x)]$$

基础式: $E_\theta[f(x)] = \int_x f(x)p(x)\mu(dx)$ 复杂!

对于 exp family:

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

所以利用这有最重要的式子:

$$A(\theta) = \sup_\mu \{\theta^T \mu - A^*(\mu)\}$$

① 变成一个优化问题.

$$\hat{\mu} \leftarrow \arg\max_\mu \theta^T\mu - A^*(\mu)$$

$$\|$$

$$\theta^T\mu + H(\mu)$$

剩下的就是对 $H(\mu)$ entropy 来估计.

② $\hat{\mu} = E_\theta[\phi(x)] = \nabla_\theta A(\theta)$

EM 算法:

E-step: inference

M-step: $\theta = \arg\max_\theta \{\theta^T\mu - A(\theta)\}$

if have complete observe.

$$\mu = \frac{1}{n} \sum_{i=1}^{n} \phi(x_i)$$

has unobservable

利用 IOPt optimization.

本节开始讲蒙特卡罗

目标: 计算 $E_\theta[f(x)] = \int_x f(x) p(x) \, d\Lambda(dx)$

$$= \begin{cases} \sum_{x \in X} f(x) p(x) & discrete. \\ \\ \int_x f(x) p(x) \, d(x) & continous. \end{cases}$$

计算量很大.

蒙特卡罗利用 the law of large number

Law of Large Number (L.L.N)

$$x_1, \cdots, x_n \overset{iid}{\sim} d$$

$$\frac{1}{n} \sum_{i=1}^{n} f(x_i) \xrightarrow{a.s.} E[f(x)]$$

expection.

almost surely converge.

given a distri., an expection is also
fixed no mattere how to calculate.

Motte Carlo sample mean

$$\frac{1}{n} \sum_{i=1}^{n} f(x_i) \quad x_i \sim d.$$

How many sample is enough to give

Such approximation? use CLT.

Central Limit Theorem

denote $I_n(f) = \frac{1}{n} \sum_{i=1}^{n} f(x_i)$

- $I_n(f)$ is also random variable.

- $E[I_n(f)] = E\left[\frac{1}{n} \sum_{i=1}^{n} f(x_i)\right]$

$$= \frac{1}{n} \sum_{i=1}^{n} E[f(x_i)] \quad (\because i.i.d)$$

$$= E[f(x)]$$

$\therefore$ this estimation is unbiased.

- $\sqrt{n}\left(I_n(f) - E\right) \xrightarrow{d} N(0, \sigma_f^2)$

$$\sigma_f^2 = Var(f(x))$$

$$Var(I_n(f)) \sim \frac{\sigma_f^2}{n}$$

when n ↑ the var ↓, usually
set a tolerance and choose n to
satisfy the tolerance.

The most difficult part of MC is
how to get the samples.

下面介绍 random sampling.

一般计算机这里都提供 rand 函数，且
好是利用这个函数来产生新的分布.

如何评估 random number generator 的好坏?
run a generator long time, it will repeat.

Linear Congrential Generator (LCG)
110010110011······· random bits
$$m = 2^{32} / 2^{64}$$

rand() function in C/C++/Jave is in this way

Mersenne Twister (MT) 更好的 generator.

C++11 <random>
MATLAB/Numpy/Julia.
$$m = 2^{19937}$$

现在可以用 generator 做的:
- generate integer [a, b]
- uniform的 distribution [0, 1) real number
- normal distribution (randn)

下面看如何生成指定分布.
Discrete Distribution

categorical distribution:
$$p = (p_1, \cdots, p_k) \quad p_1 + \cdots + p_k = 1$$

例如 (0.2, 0.5, 0.3)
stick breaking



```
        0.2        0.5         0.3
    ⌢⌢⌢   ⌢⌢⌢⌢⌢   ⌢⌢⌢
   |___|_____|_____|
       0.2          0.7
```

利用 uniform 生成一个数看落在哪个区间.

① Linear search   $O(k)$
    遍历所有分段, 看到底落在哪里.
    not efficient.

② Sorted search   $O(k)$
    加快速度: 把概率从大到小排序.

```
   |_____|____|__|
         0.5    0.8
```

③ Binary search.
    按概率选个数.



$$O(\log_2 k)$$

可以 构造个 Hoffman Search Tree,
来让这个树最优.

$$O(Entropy) \leq O(log_2 k)$$

④ Alins Table. 构造复杂, 但使用 很快

$$O(1).$$

现在看如何以任意分布采样.
已知 $p(x)$.

Transform Sampling

$$u \sim U[0, 1]$$

$$f(u) \sim P \quad \text{how can we get f?}$$

Sampling   P(x)   

Transform  Sampling

$U \sim Uniform([0, 1])$

$X \sim T(u)$   通过变换 T, 从均匀分布生成任意分布.

$X = F^{-1}(u)$

$F$: cumulation distribution function
(cdf)

$$F(x) = Pr(X \leq x) = \int_{-\infty}^{x} P(v) dv.$$

以下证明 x 的分布 也是我们需要的.

Proof:

$X = F^{-1}(u)$

$Pr(X \leq x) = F(x)$
⇑
$Pr(F^{-1}(u) \leq x) = F(x)$

$Pr(u \leq F(x)) = F(x)$

uniform 即 $Pr(u \leq x) = x$

$Pr(u \leq 0.6) = 0.6$

例: Exponential distribution $P(x) = exp(-x)$

$$F(x) = 1 - exp(-x)$$

$$y = 1 - e^{-x} \Rightarrow x = -log(1-y)$$

$$\therefore F^{-1}(x) = -log(1-x)$$

$u \sim U[0, 1] \qquad x = -log(1-x)$

实际上 $u \sim U[0,1] \Leftrightarrow 1-u \sim U[0,1]$

$\therefore$ 也可以 $x = -log\ u$.

上面实现了单一变量的采样, 下面看多变量.

$$N(\mu, \Sigma)$$

$$P(x) = \frac{1}{\sqrt{(2\pi)^d |\Sigma|}} exp\left(-\frac{1}{2}(x-\mu)^T \Sigma^{-1}(x-\mu)\right)$$

cdf 很难求. 可以从单变量的正态简单解释

$$u \sim N(0, 1)$$

$$\begin{pmatrix} u_1 \\ \vdots \\ u_d \end{pmatrix} \sim N(0, I)$$

正态分布中，若 $x \sim N(\mu, \Sigma)$

$$x + b \sim N(\mu + b, \Sigma)$$

$$Ax \sim N(A\mu, A\Sigma A^T)$$

∴ 生成步骤：

① $\begin{pmatrix} u_1 \\ \vdots \\ u_d \end{pmatrix} \sim N(0, I)$

② Get $A$, s.t. $AA^T = \Sigma$     Cholesky $\checkmark$ Decomposition

$$A \cdot u \sim N(0, AA^T) = N(0, \Sigma)$$

③ $Au + \mu \sim N(\mu, \Sigma)$

上面依赖于分布的性质.
下面看 Rejection Sampling

# Rejection Sampling

target distribution $P(x)$.

proposal distribution $q(x)$    (easy to sample)

Steps:

① $x \sim q$

② Accept $x$ with $\dfrac{P(x)}{M \cdot q(x)}$

$M$: s.t. $M q(x) \geq P(x)$



$Mq(x)$

$P(x)$: hard to sample

$q(x)$: easy to sample.

$a(x_1) = \dfrac{P(x_1)}{M q(x_1)}$ 很小.

$a(x_2) = \dfrac{P(x_2)}{M q(x_2)}$ 几乎等于1

Auxiliary Variable.

① $x \sim q$

② $u | x \sim Bernoulli \left( \dfrac{P(x)}{M q(x)} \right)$

$q(x) \cdot a(u|x)$

$p(x|u=1)$ 因为只有 $u=1$ 时，才接收，∴这个条件概率才是 x 的分布。

下面证明 $p(x|u=1)=p(x)$.

$$p(x|u=1) \propto q(x) \cdot a(u=1|x)$$

$$= q(x) \cdot \frac{p(x)}{Mq(x)}$$

$$= \frac{p(x)}{M}$$

优点：
① Simple, 只要知道了 $p(x)$ 就能采样.

缺点：
M 可能会特别大.



这样接收率就特别低

Overall acceptance rate

$$P_q(u=1) = \int q(x) \cdot \frac{p(x)}{Mq(x)} d(x)$$

$$= \int \frac{p(x)}{M} d(x) = \frac{1}{M}$$

∴ 接收率为 $\frac{1}{M}$

为了！减小 $M$，可以找一个和 $P(x)$ 很接近的分布 $q(x)$.

## Importance Sampling

difficult

$$E[f(x)] \approx \frac{1}{n} \sum_{i=1}^{n} f(x_i) \quad \text{with } x_i \sim p$$

proposal distribution $q$

$$E_p[f(x)] = \int p(x) f(x) dx = \int f(x) \frac{P(x)}{q(x)} \cdot q(x) \, dx$$

$$= E_q \left[ f(x) \frac{P(x)}{q(x)} \right]$$

$$\approx \frac{1}{n} \sum_{i=1}^{n} f(x_i) w(x_i) \quad \text{with } x_i \sim q$$

$w(x_i) = \frac{P(x_i)}{q(x_i)}$ is called importance weight

$w(x_i)$ 有时也很难计算，例如：

$$P(x) = \frac{1}{Z_p} \exp(A(x))$$

$$q(x) = \frac{1}{Z_q} \exp(B(x))$$

$$Z_p, \quad Z_q \text{ 很难计算}.$$

只能计算 $\dfrac{\tilde{p}(x)}{\tilde{q}(x)} = exp(A(x) - B(x))$

只能算出一个 scale, 不能算出具体值.

self-normalized $z_p, z_q$ 这个 normalize 项不

计算 $\tilde{w}(x_i) = \dfrac{\tilde{p}(x_i)}{\tilde{q}(x_i)}$ 好计算, 所以直接对 $w$ 进行 normalize.

$$w(x_i) = \dfrac{\tilde{w}_i}{\sum\limits_{j=1} \tilde{w}_i}$$ 可以证明是无偏估计

和 Rejection Sampling 类似, 也存在相同问题.



会采到很多 weight 很小的样本, 会不稳定

下面介绍更有效的方式 MCMC

(实际上 GAN 也是一种 sampling, generator
采样, discriminator 拒绝.)

# MCMC (Markov Chain Monte Carlo)

- Markov Chain

$$X_0, \quad X_1 \cdots, X_T, \cdots$$

$$p(X_t | X_{t-1}) = p(X_t | X_{t-1}, X_{t-2}, \cdots)$$

$$\text{Markov}(\pi_0, P)$$

↑ initial distribution    ↖ transition probability matrix

如果知道当前状态 $X$

$$X \to X' \quad P(X'|X) \leftarrow \quad P$$

如果知道的是当前状态的概率分布 $\mu$.

$$P(X') = \sum_X \mu(x) \, p(X'|X)$$

$$= \sum_X \mu(x) P(X, X')$$

$$\mu' = \mu P$$

如果 $\mu = \mu P$ 则称 $\mu$ 为 invariant distribution

if $P$ is irreducible & aperiodic

there exists a unique $\mu$

s.t. $\mu = \mu P$

reducible    irreducible

irreducible & aperiodic 叫做 ergodic

MCMC idea:
target distribution $\mu$  hard to sample.
$\Downarrow$
construct a MC with P
    s.t. $\mu = \mu P$

$$x_0 \to x_1 \to \cdots \to x_{100000} \overset{\sim \mu}{\to} X \to X \cdots$$

target distribution $\pi$.
ergodic Markov chain P  s.t. $\pi = \pi P$

Ergodic Throem:

$$\frac{1}{n} \sum_{i=1}^{n} f(\tau + i \cdot m) \to E_\pi [f(x)]$$

Burning          Sampling



$\tau$: burning time
    开始时的分布 指定 还不是 $\mu$，运行一段时间
    间让 分布 接近 $\mu$.

m: 一般采样都是假设的样本间相互独立，
如果采了个马氏链，那么这样个样本
相关性太强。

m used to reduce dependence.

下面看如何构造 P.

P: $\pi P = \pi$

$$\sum_x \pi(x) P(x,y) = \pi(y) \quad \forall y$$

Detail Balance 介类

$$\pi(x) P(x,y) = \pi(y) P(y,x) \quad \forall x,y$$

这个式子更简单，没有求和，积分 (更形)

Proof:
$$\sum_x \pi(x) P(x,3) = \sum_3 \pi(y) P(x,y)$$

$$= \sum_x \pi(y) P(x|y) = \pi(y) \underbrace{\sum_x P(x|3)}_{=1} = \pi(3)$$

Matropolis-Hasting (M-H algorithm)

构造 MC 的算法 类似于 reject sampling, 先由 proposal kernal
产生一个样本, 然后来拒绝.

1. $Q(x,y)$ given $x$, next step $y \sim Q(x,3)$
proposal kernal.

$$Q \left. \begin{array}{|c|} \hline \phantom{xxxx} \\ \hline x \phantom{|\pi|\pi|\pi|\pi|} \\ \hline \end{array} \right\}_{Q_x}$$

2. Every step:

1) $y \sim Q_x$   $x$是当前状态

2) accept $y$ with chance $a(x,y) = \min\{r(x,y), 1\}$
where $r(x,y) = \dfrac{\pi(y) \, Q_y(x)}{\pi(x) \, Q_x(y)}$

$$= \frac{\pi(y) \, Q(y \to x)}{\pi(x) \, Q(x \to y)}$$

$$= \frac{\frac{1}{z} h(y) \, Q(y \to x)}{\frac{1}{z} h(x) \, Q(x \to y)}$$

$$= \frac{h(y) \, Q(y \to x)}{h(x) \, Q(x \to y)}$$



$\pi -$ 函数桥 correctness 和 efficiency.

— correctness
— efficiency $\begin{cases} \text{acceptance (high)} \\ \text{mixing rate} \\ \text{(convergence rate)} \, |\mu_t - \pi| \end{cases}$ trade off

如果样本都在一起,
conservative → $\begin{cases} \text{high acceptance} \\ \text{low mixing rate.} \end{cases}$

如果样本能得跳很远.
aggressive → $\begin{cases} \text{low acceptance} \\ \text{accderate mixing rate.} \end{cases}$

回顾：
Markov Chain Monte Carlo

target distribution $\pi$ 很难采样

构造 Markov chain $P$: $\pi P = \pi$

如何构造 $P$?

通过 Detailed Balance

$$\pi(x) P(x,y) = \pi(y) P(y,x)$$

Metroplis-Hastiny Alg.

$Q$: proposal kernal

$$Q(\underset{\underset{\text{current}}{\uparrow}}{x}, \underset{\underset{\text{next}}{\uparrow}}{y})$$ give proposal for next based on current.

$$x \Rightarrow y \sim Q(x,y)$$
$$q_x(y)$$

acceptance rate $a(x,y) = \min\left\{\dfrac{\pi(y) q_y(x)}{\pi(x) q_x(y)}, 1\right\}$

不宜证明 至 $a(x,y)$ 满足 Detailed Balance.

proof:

$$\pi(x)\,Q(x,y)\min\left\{\dfrac{\pi(y)Q(y,x)}{\pi(x)Q(x,y)},1\right\}$$   proposal * acceptance ratio

$$=\min\left\{\pi(x)Q(x,y)\cdot\dfrac{\pi(y)Q(y,x)}{\pi(x)Q(x,y)},\ \pi(x)Q(x,y)\right\}$$

$$=\min\left\{\pi(y)Q(y,x),\ \pi(x)Q(x,y)\right\}$$

$$\dfrac{\pi(y)\,q_y(x)}{\pi(x)\,q_x(y)}$$     很难计算, ∵ $\pi(y)=\dfrac{1}{Z}h(y)$

$$=\dfrac{h(y)\cdot q_y(x)}{h(x)\,q_x(y)}$$     ↑ normalize term 很难算.

抵消了 normalize term

假设



$x_1$

$$x_{t+1}\leftarrow x_t+\varepsilon_t \qquad \varepsilon_t\sim N(0,\sigma^2)$$

也是 thermal Q

$$q_x(y)\sim N(x,\sigma^2)$$

例如, 开始的 $x_1$ 超图, 概率很小.

用这个 kernal : $q_x(y) \sim N(x, \sigma^2)$

注意到它是一个对称的 kernal 即:

$q$ 只与 $x, y$ 的距离有关.

symmetric kernal

$$q_x(y) = q_y(x)$$

$$\therefore \frac{\pi(y) q_y(x)}{\pi(x) q_x(y)} = \frac{\pi(y)}{\pi(x)}$$

$\therefore$ 跳到概率高的 地方, 其接收率就高,
跳到概率低的地方接收率低.

那么如何设计一个好的 proposal kernal.
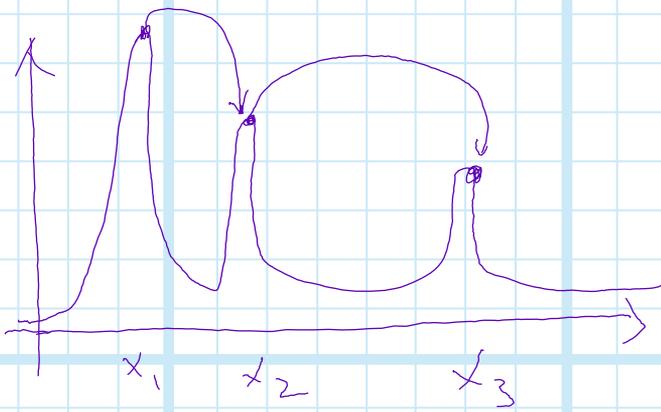
要求:

— high acceptance ratio

— explore the space efficient.



接收率低　　　　　　不能很好地探索到

好的 proposal:
Jump from peak to peak.

Gibbs Sampling

$$p(x_1, \cdots, x_n) \quad 多变量采样,$$

每次只改变一维.

$$S_1 = (x_1, \boxed{x_2}, x_3)$$

$$\downarrow$$

$$S_2 = (x_1, \boxed{x_2'}, x_3) \qquad x_2' \sim P(x_2 | x_1, x_3)$$

下面来看正确性:

$$P(x, y)$$

$$S = (x, y)$$

$$S_2 = (x', y) \qquad Q(S \to S') \propto p(x' | y)$$

$$r(S \to S') = \frac{P(x', y)}{P(x, y)} \cdot \frac{P(x | y)}{P(x' | y)}$$

$$= \frac{P(y) \, P(x' | y)}{P(y) \, P(x | y)} \cdot \frac{P(x | y)}{P(x' | y)} = 1$$

$$a(S \to S') = \min\{r(S \to S'), 1\}$$

∵ Gibbs Sampling 的接收率恒等于1
效率很高.

如何选变化的坐标.

Cycling Scheme

— fix scheme  $1 \to 2 \to 3 \to 4 \to 1 \to 2 \to 3 \to \cdots$

— random cycle.  $1 \to 3 \to 4 \to 2 \to 4 \to 3 \to 1 \cdots$

Gibbs Sampling 在参数空间很大的时候 不是很高效
很多算法发明出来提升 efficiency.

Collapsed Gibbs Sampling



$$g \sim N(\mu, \sigma_0^2)$$

$$x_i \sim N(g, \varepsilon^2) \quad \varepsilon << \sigma$$

如果用 Gibbs Sampling

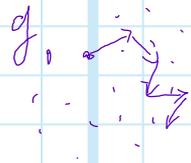$$x_i | g \sim N(g, \varepsilon^2) \quad \forall i = 1, 2, 3, 4. \quad (x_i 间相互独立)$$

$$g | x_1, x_2, x_3, x_4 \sim N(\mu, \sigma'^2)$$

$$\mu' = \frac{\frac{1}{\sigma_0^2} \cdot \mu_0 + \sum_{i=1}^{n} \frac{1}{\varepsilon^2} x_i}{\frac{1}{\sigma_0^2} + \sum_{i=1}^{n} \frac{1}{\varepsilon^2}}$$

$$(\sigma'^2)^{-1} = \frac{1}{\sigma_0^2} + \sum_{i=1}^{h} \frac{1}{\varepsilon^2} \gg \frac{1}{\sigma_0^2}$$
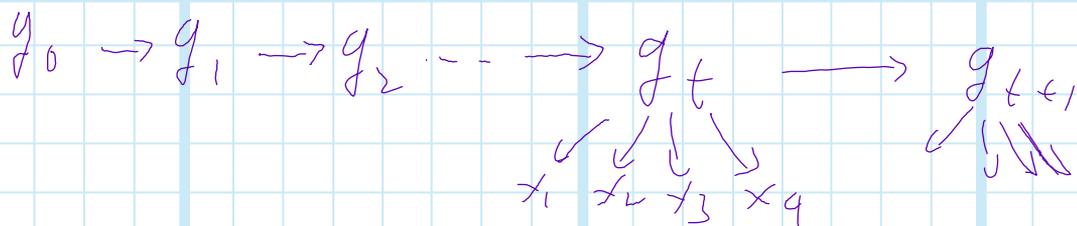
$$\sigma'^2 \ll \sigma_0^2$$

$g_0 \dashrightarrow$

$\overset{\bullet}{\mu_0}$

mutual locking: 当前结点 conditioned on other nodes

∴ mutual locking, $g_0$ 到 $\mu_0$ 连接很少度.

解决方法:

marginalize out

$$y_0 \rightarrow g_1 \rightarrow g_2 \cdots \rightarrow g_t \longrightarrow g_{t+1}$$

$x_1 \ x_2 \ x_3 \ x_4$

Collapsed Gibbs Sampling

$$N(\mu, \sigma^2) \rightarrow$$



迭代时，变量相互影响，move 很少

当采样 $\theta$ 时，不需考虑 $x_1, x_2, x_3$ 影响。

Rao-Blackwell Theorem

① $P(X, Y)$

$(x_1, y_1), (x_2, y_2) \cdots, (x_n, y_n) \sim P$

$$E[h(x, y)] \approx \frac{1}{n} \sum_{i=1}^{n} h(x_i, y_i) \rightarrow \bar{h}_1$$

以上是标准的 MC 过程。

② $P(X, Y)$

$\downarrow$ margional out $Y$
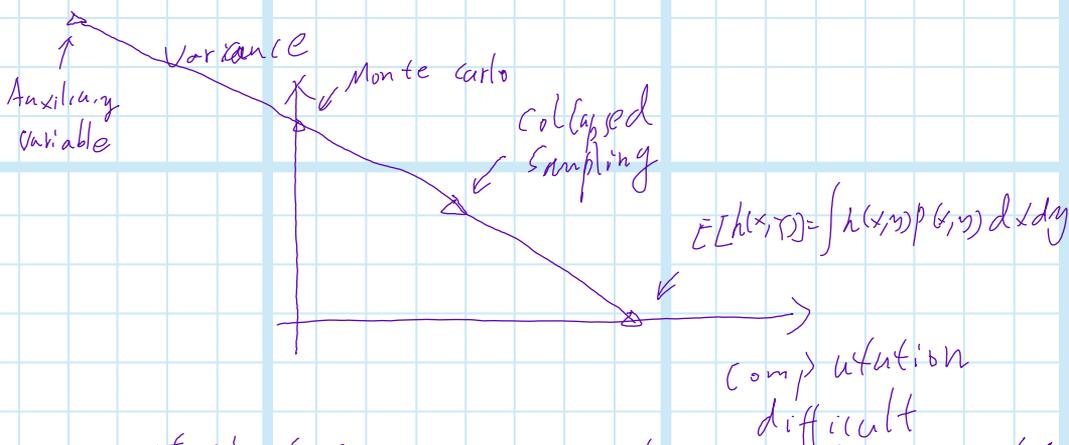
$P(X)$

$x_1, x_2, \cdots, x_n \sim P(x)$

$$E[h(X, Y)] \approx \frac{1}{n} \sum_{i=1}^{n} E_{Y|x_i}[h(x_i, Y)]$$

$$\rightarrow \bar{h}_2$$

正确性:

$$E[h(X,Y)] = \int_{X\times Y} h(x,y)\, p(x,y)\, \mu(dx\,dy)$$

$$= \iint_{X,Y} h(x,y)\, p(y|x)\, p(x)\, \mu(dx)\mu(dy)$$

$$= \int_X \int_Y h(x,y)\, p(y|x)\, dy\, p(x)\, dx$$

$$= \int_X \underbrace{E_{Y|X}[h(x,Y)]}_{f(x)}\, p(x)\, dx$$

$$\sim \frac{1}{n}\sum_{i=1}^{n} f(X_i)$$

Rao-Blackwell Theorms

$$Var(\bar{h_1}) \geqslant Var(\bar{h_2})$$



MC 直接采样, 不以求算做积分计算.
直接计算没有 Variance.
Collapsed Sampling 是一个折衷的方式.

Collapsed Sampling 先积分掉一些变量,
再来 Sampling, 这种方式叫做 Rao-Blackditation

Collapsed Sampling 积掉了一些变量, 一种
相反的方法是引入更多的辅助变量,
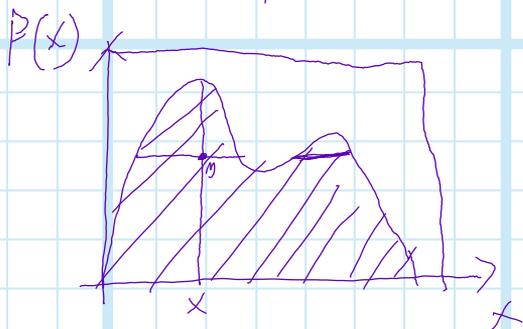有些情况下, 需要让 Sampling 比 MC
更简单

Sampling with Auxiluary variable.

$$target : P(x)$$
$$\downarrow$$
$$P(x, u)$$
$$(x_1, u_1) \cdots , (x_n, u_n)$$

① Slice Sampling



每次不以 $x$ 采样, 而是以二维的阴影
部分采样

可以用 Gibbs Sampling 来从阴影部分采样

$$y|x \sim U[0, f(x)]$$
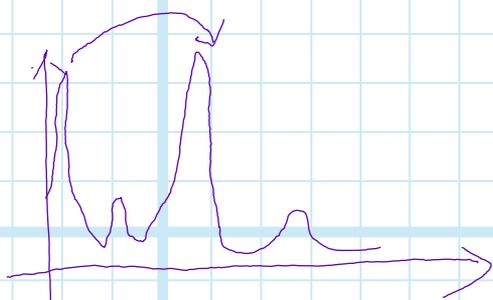$$x|y \sim U[\{x : f(x) \geq y\}]$$

Property:
~ very efficient
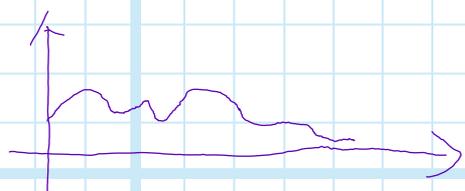  对 Gibbs Sampling, 所有的样本都被接收.

~ $\{x : f(x) \geq y\}$ 不好计算

② Tempering

introduce temperature to certain distribution



目标: Jump from peak to peak.

很困难, 所以考虑一个 smooth 的分布



引入 temperage

$$p(x) = \frac{1}{Z} \exp(-E(x))$$ ← Gibbs Distribution

$$P_\alpha(x) = \frac{1}{Z_T} \exp\left(-\frac{E(x)}{T}\right)$$

T temperature.

温度越高, 越 smooth

低 peaky

$T \gg 1$ : much easier to sample

$T = 1$ : Our target distribution

# Basic Idea

introduce $\tau$ as an auxiliary varible.

## Simulated Tempering

$$p(x, k) = \frac{\pi_k}{Z_k} \exp\left(-\frac{E(x)}{T_k}\right)$$

构造了一个不同温度 $\{T_1, \cdots, T_n\}$ 的链接

$\pi$ 是比例

$$P_1 = \frac{1}{Z_1} \exp(-E(x)) \qquad \begin{array}{c} \pi \\ 0.3 \end{array} \quad \pi_1 \qquad a_1$$

$$P_2 = \frac{1}{Z_2} \exp\left(-\frac{E(x)}{T_2}\right) \qquad 0.4 \quad \pi_2 \qquad a_2$$

$$P_3 = \frac{1}{Z_3} \exp\left(-\frac{E(x)}{T_3}\right) \qquad 0.3 \quad \pi_3 \qquad a_3$$

MCMC

一、transition Proposal

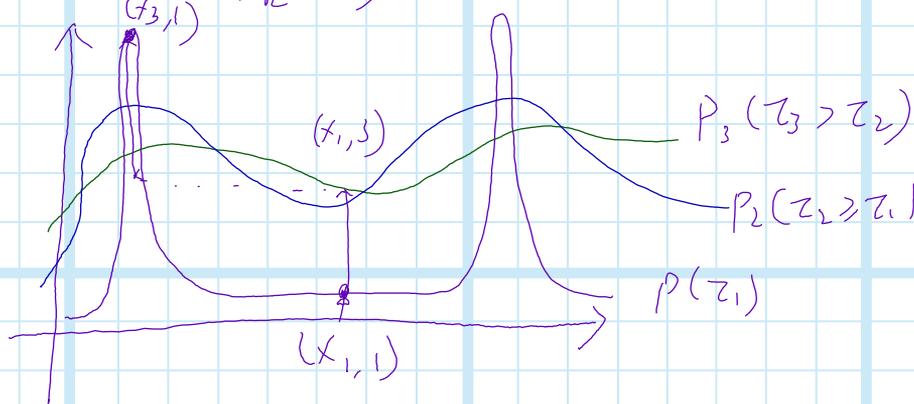$(x, k) \longrightarrow (x', k)$  $x' \sim q_k(x \rightarrow x')$

二、Temperature switch

$q_T(k \rightarrow k')$

$a(k, k'|x) = \min\{1, \frac{P(x, k')}{p(x, k)}\}$    温度转移的
接收率。

另一种方法

$p(k|x) \propto \pi_k \cdot P_k(x)$



在温度高的地方, 容易 explore space, 在某一
另一个分布大的时候 就切换分布进行采样

三、高温分布的作用是 bridge for large move.
低温分布负责采样, 高温分布负责移动.

# Parallel Tempering (告会并行计算)
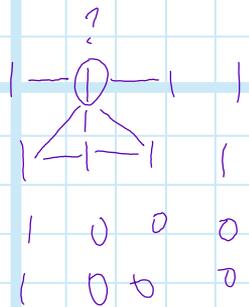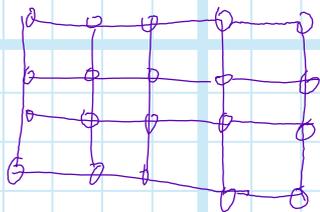
## Croplica switch

$$\text{target} \rightarrow P_1 \rightarrow x_1^{(1)}, \cdots, x_n^{(1)}$$
$$P_2 \rightarrow x_1^{(2)}, \cdots, x_n^{(2)}$$
$$P_3 \rightarrow x_1^{(3)}, \cdots, x_n^{(3)}$$

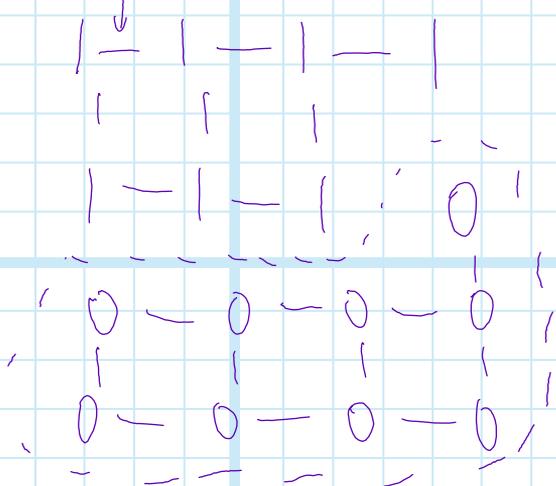$$(x_{t_1}^{(1)}, x_{t_2}^{(2)}, x_{t_3}^{(3)}) \rightarrow (x_t^{(1)}, x_t^{(2)}, t_t^{(3)})$$

## Swendsen - Wang



## Gibbs Sampling

$$p(x) \propto \exp\left(-\sum_{ij} w_{ij} \mathbb{1}(x_i = x_j)\right)$$

$g_{ij} := 0$时有边, $=0$时无边.



Basic: cast node into groups and update each group together

$$P(x|\theta) = \frac{1}{Z} \prod_{(i,j) \in E} f_{ij}(x_i, x_j)$$

$$= \frac{1}{Z} exp\left(\sum_{i,j \in E} \theta_i x_i x_j\right)$$

$U_{ij}$ : for each edge turn on or off

$$P'(x, u) = \frac{1}{Z'} \prod_{(i,j) \in E} g_{ij}(x_i, x_j, u_{ij})$$

$$g(x_i, x_j, u_{ij}) = \begin{cases} exp(-\theta_{ij}) & u_{ij} = 0 \\ \mathbb{1}(x_i = x_j)(e^{\theta_{ij}} - e^{-\theta_{ij}}) & u_{ij} = 1 \end{cases}$$

$u_{ij} = 1 \Rightarrow x_i = x_j$ ( all variable in the

$x_i \neq x_j \Rightarrow u_{ij} = 0$ same component should have the same value)

$$g(x_i, x_j, u_{ij})$$

$$= f_{ij}(x_i, x_j) \, q_{ij}\left(u_{ij}|x_i, x_j\right)$$

when $x_i \neq x_j$ ; $u_{ij} = 0$

when $x_i = x_j$ $P(u_{ij} = 0 | x_i = x_j) = exp(-\theta_{ij})$

算法步骤是

1. Clustering step

$u|x \quad \sim \quad P(u_{ij} = 0 | x_i = x_j) = exp(-\theta_{ij})$

2. mapping $x \sim P(x|u)$

$$p'(x, u) = \prod_{(i,j) \in E} g(x_i, x_j, u_{ij})$$

$$= \prod_{(i,j) \in E} f_{ij}(x_i, x_j) \, q(u_{ij} | x_i, x_j)$$

$$= \left( \prod_{(i,j) \in E} f_{ij}(x_i, x_j) \right) \cdot \prod_{(i,j) \in E} q(u_{ij} | x_i, x_j)$$

$$\xrightarrow{\text{marginalize}} P(X)$$

注意!

## Common ideas of auxiliary variables

$$p(x) \xrightarrow{\quad aux \quad} p(x, u)$$

$\Bigg\downarrow$ Gibbs Sampling

Given $u$: $p(x|u)$

Given $x$: $p(u|x)$

| | Slice Sampling | Simulated Tempering | Swendson - Wang |
|---|---|---|---|
| Aux var $u$ |  | $T_k$: Temperature | use $u$ as connections $u_{ij}$ $(i,j) \in E$ |
| Gibbs Sampling $\begin{cases} \\ \\ \end{cases}$ $x \sim p(x|u)$ | $x \sim U[\{x_i : f(x) \geq y\}]$ | $x \sim P(x|T)$ | Given $u$, update $x$ by components |
| $u \sim p(u|x)$ | $u \sim U[0, f(x)]$ | $p(z|x) \propto T(z) p(x|\pi)$ | $u_{ij} \sim q(u_{ij} | x_i, x_j)$ |

可以给出 $\mu \sim P(\mu|x)$ 它一步的给出是
switch the environment
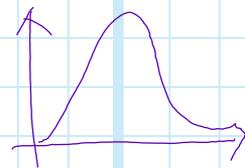$x \sim P(x|\mu)$ 你可以 Sampling


Variational Auto - Encoder

Generative Adverial Learning.

是连续的方法.

Fundamental Problem: How to character
a distribution

— Descriptive way
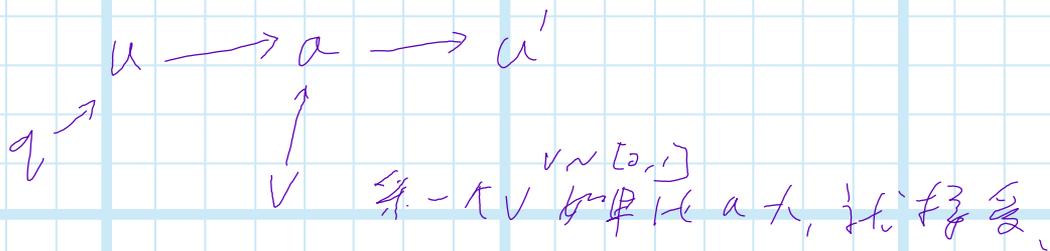density function $P(x)$



— Construtive way (Generative Way)
focus on how to get samples
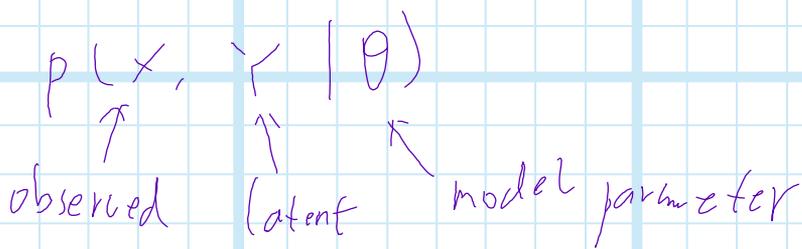It's actually a sampler.

○ Transform sampler : 寻找形式的算法

$u \overrightarrow{\pi} T(u) \sim P$

△ Rejection Sampling

$$q \nearrow \quad u \longrightarrow a \longrightarrow u'$$

$$V \qquad \underset{}{V \sim [0,1]}$$

第一次V如果比 $a$ 大，就拒绝.

△ MCMC

EM

$$p(X, Y | \theta)$$

observed    latent    model parmeter

$$\hat{\theta} = \arg\max_{\theta} p(X|\theta) \qquad Y 滴 marglize 掉$$

很难积分. 实际上:

$$p(X|\theta) = E_{q(\cdot)}[p(X,Y)|\theta]$$

$$\rightsquigarrow p(Y|X; \theta) \quad \Leftarrow difficult$$

Variational E-M

$$q = q_{y_1} \cdot q_{y_2} \cdots \quad 假设 q 可分解.$$

Mean field Approximation

这个估计可以直接用 NN 来做

$p(Y|X; \theta)$ 用一个神经网络

也可以用 law of large Number

GAN:

$\qquad$ D $\rightarrow$ Descriptive way

$\qquad$ G $\rightarrow$ Constructive way

Revisit

$P(x, z)$
$\theta_\rho \quad q$
observe latent

$$L(\theta) = \frac{1}{n} \sum_{i=1}^{n} \log P_\theta(x_i) \quad \leftarrow \begin{array}{l} \text{margional density} \\ \text{hard to compute.} \\ \text{consider lower bound} \end{array}$$

$$E_q[p(x,z)] = E_q[\log P(x) + \log P(z|x)]$$

$$= \log P_\theta(x) + E_q[\log P(z|x)]$$

① $\log P_\theta(x) = E_q[\log P(x,z)] - E_q[\log P(z|x)]$

② $L_x(\theta,q) = E_q[\log P(x,z)] - E_q[\log q(z)]$

∴ ① — ② :

$\log P_\theta(x) - L_x(\theta, q)$

$= E_q[\log q(z) - \log P(z|x)]$

$$= E_q\left[\frac{\log q(z)}{\log P(z|x)}\right] = D_{KL}[q(z) \| P(z|x)] \geq 0$$

那么 q(z)如何建模.
— mean field approximation

$$q(z_1, z_2) = q_1(z_1) q_2(z_2) \quad \text{直接假设取以分解}$$

— neural network
VAE

目标: 建模 $L_x(\theta, q) = E_q[\log P(x,z)] - E_q[\log q(z)]$

假设 $q$ 的参数为 $\phi$

$$L(\theta, \phi) = \underbrace{E_\phi[\log P_\theta(x,z)]}_{\Downarrow} - E_{q_\phi}[\log q_\phi(z)]$$

monte carlo

sampler $g$ $\quad \varepsilon \sim P_0 \quad \overset{\text{某种随机}}{g(\varepsilon, x)} \sim q_x(\cdot)$

$\underset{\text{randomness}}{\nearrow} \quad \underset{\text{condition}}{\nwarrow}$

$$E_q[f(z)] \approx \frac{1}{n} \sum_{i=1}^{n} f(z_i) \quad z_i \sim q_\phi(\cdot|x)$$

$$= \frac{1}{n} \sum_{i=1}^{n} f(g_\phi(\varepsilon_i, x))$$

代入目标得:

$$L(\theta, \phi)$$

$$= \frac{1}{L} \sum_{\ell=1}^{L} \left\{ \log P_\theta(x, z^\ell) - \log q_\phi(z^\ell|x) \right\}$$

where $z^\ell = g_\phi(q^\ell, x) \quad \varepsilon^\ell \sim P_0$

neural network
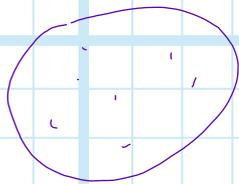
$P_\theta(x, z)$ 是对分布的 descriptive 指生式.

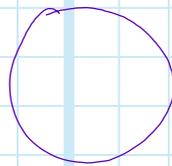而 $q(z, x)$ 则是 generative way.

VAE 同时使用了 descriptive 和 generative way
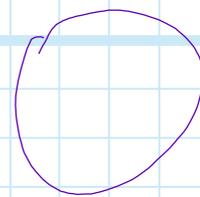
# GAN

GAN 的 generator 和 desxriminator 是一种直观的解释

real data

$G(z)$

generator 实质是一个 sampler

D

(fake)

现在想直接用 generator 来做采样, 需要一个 descriptive way of real data, 这其实 就是 Descriminator

$$\min_{G} \max_{D} E_{x \sim P_{data}} \left[ log D(x) \right] + E_{x \sim P(z)} \left[ log(1 - D(G(z))) \right]$$

$$x_{real} \quad x_{fake}$$

$$D: \overset{max}{\quad} log D(x_{real}) + log(1 - D(x_{fake}))$$

$$G: \overset{min}{\quad} log(1 - D(G(z)))$$
discriminative training of
a generation

capture the density
in a transformed way

$$D^*(x) = \frac{P_{data}(x)}{P_{data}(x) + P_g(x)}$$

$$= \frac{1}{1 + \frac{P_g(x)}{P_{data}(x)}} = \frac{1}{1 + \left( \frac{P_{data}(x)}{P_g(x)} \right)^{-1}}$$

real data is non-parametric distribution

$D^*(x)$ convert it to a parametric way.

因为：
$$\max P_\theta(G(\varepsilon_i))$$

$$log(1 - D(G(\varepsilon)))$$ transformed way

NCE $\overset{estimation}{\rightarrow}$ 提供了 GAN 的理论基础!

NCE 没有 generator, 而是用 fake example.

G-M 阶段：inference

estimator

两个指标很重要

consistuncy：当足够多的样本,能收敛.

efficiency：收敛速度.