

1. Basics

高维随机变量 $\{$

- 边缘概率 $P(x_i)$
- 条件概率 $P(x_j | x_i)$

$$\text{Sum Rule: } P(x_1, x_2) = \int P(x_1, x_2) dx_2$$

$$\text{Product Rule: } P(x_1, x_2) = P(x_1) \cdot P(x_2 | x_1) = P(x_2) \cdot P(x_1 | x_2)$$

$$\text{Chain Rule: } P(x_1, x_2, \dots, x_p) = \prod_{i=1}^p P(x_i | x_1, x_2, \dots, x_{i-1}) \cdot P(x_1)$$

$$\text{eg: } P(A_1, A_2, A_3) = P(A_3 | A_2, A_1) \cdot P(A_2 | A_1) = P(A_3 | A_2, A_1) \cdot P(A_2 | A_1) \cdot P(A_1)$$

$$\text{Bayesian Rule: } P(x_2 | x_1) = \frac{P(x_1, x_2)}{P(x_1)} = \frac{P(x_1, x_2)}{\int P(x_1, x_2) dx_2} = \frac{P(x_1) \cdot P(x_2 | x_1)}{\int P(x_1) \cdot P(x_2 | x_1) dx_1}$$

困境: 维度高, 计算复杂 $P(x_1, x_2, \dots, x_p)$ 计算量大

$$\downarrow \text{Chain Rule} \rightarrow P(x | y) = \prod_{i=1}^p P(x_i | y)$$

$$\text{简化} \xrightarrow{\text{相互独立}} P(x_1, \dots, x_p) = \prod_{i=1}^p P(x_i) \xrightarrow{\substack{\text{Markov Property} \\ \text{Hull}}} x_i \perp x_{i+1} | x_1, \dots, x_{i-1}$$

(满足Markov假设)

条件独立性

$$x_a \perp x_b | x_c$$

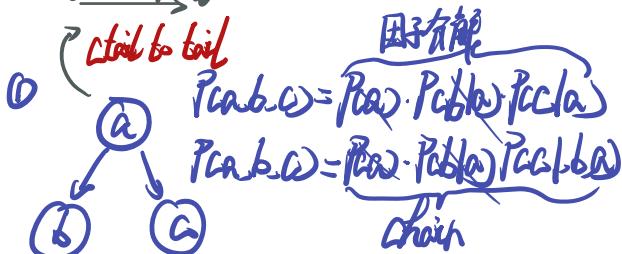
Bayesian Network

$$P(x_1, x_2, \dots, x_p) = \prod_{i=1}^p P(x_i | x_1, x_2, \dots, x_{i-1}) \cdot P(x_1)$$

事件独立性: $x_1 \perp x_B \perp x_C$ ↑ 次序独立性

$$\text{因子分解: } P(x_1, x_2, \dots, x_p) = \prod_{i=1}^p P(x_i | f_{par(i)})$$

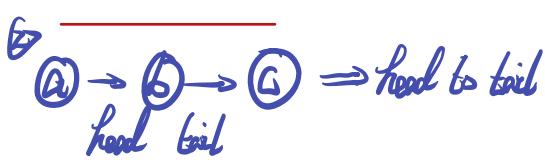
tail to head



$(c \perp b) | a$

↑

若 a 被观测，则路径被阻塞



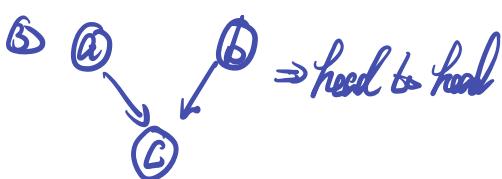
$a \perp c | b$

↑
b 被观测，则路径被阻塞

$$\begin{aligned} P(a,b,c) &= P(c) \cdot P(b|c) \cdot P(a|b) \\ &= P(c) \cdot P(b|c) \cdot P(a|b,c) \end{aligned}$$

$$P(a|b,c) = P(a|b,c)$$

$$\begin{aligned} P(a|b,c) \cdot P(c) &= P(c) \cdot P(a|b,c) \\ &= P(a|b,c) \end{aligned}$$



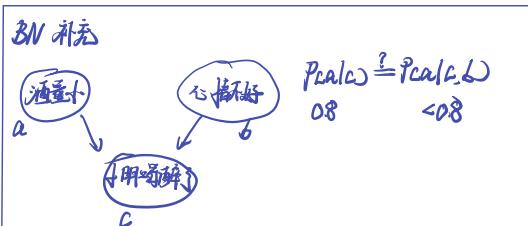
默认 a,b, 路径堵塞

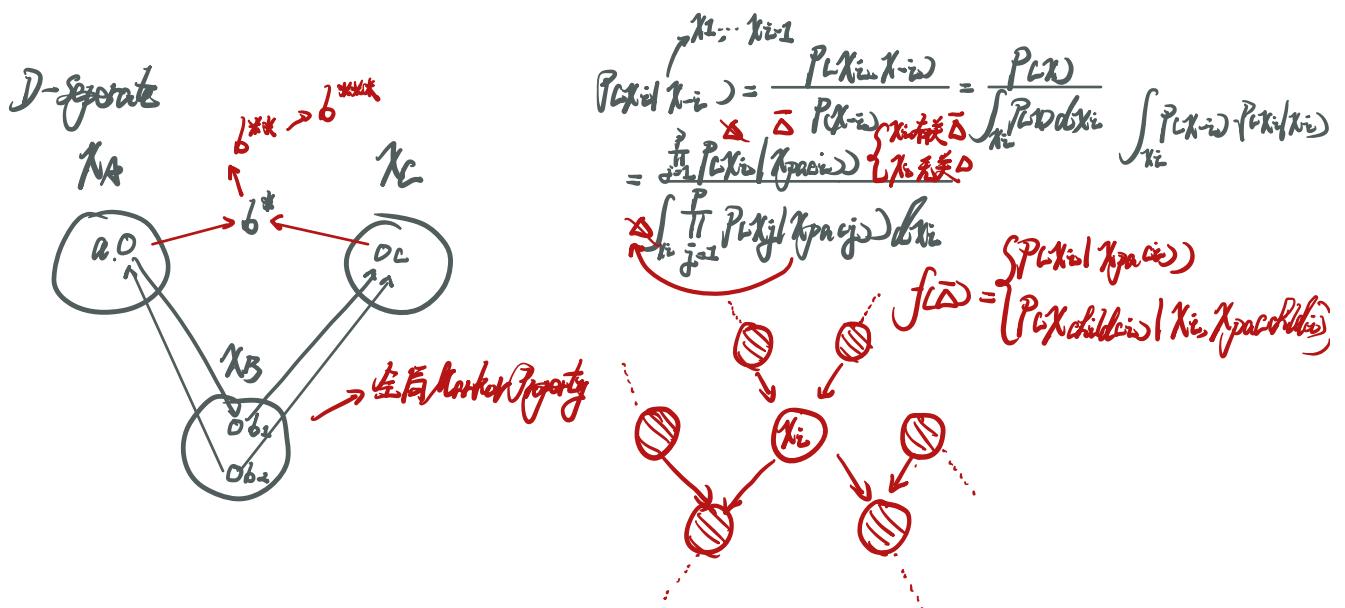
若 c 被观测，则路径能通

$$\begin{aligned} P(a,b,c) &= P(c) \cdot P(b|c) \cdot P(a|c) \\ &= P(c) \cdot P(b|c) \cdot P(b|a) \cdot P(a|c) \end{aligned}$$

$$\Rightarrow P(b) = P(b|a)$$

$a \perp b$

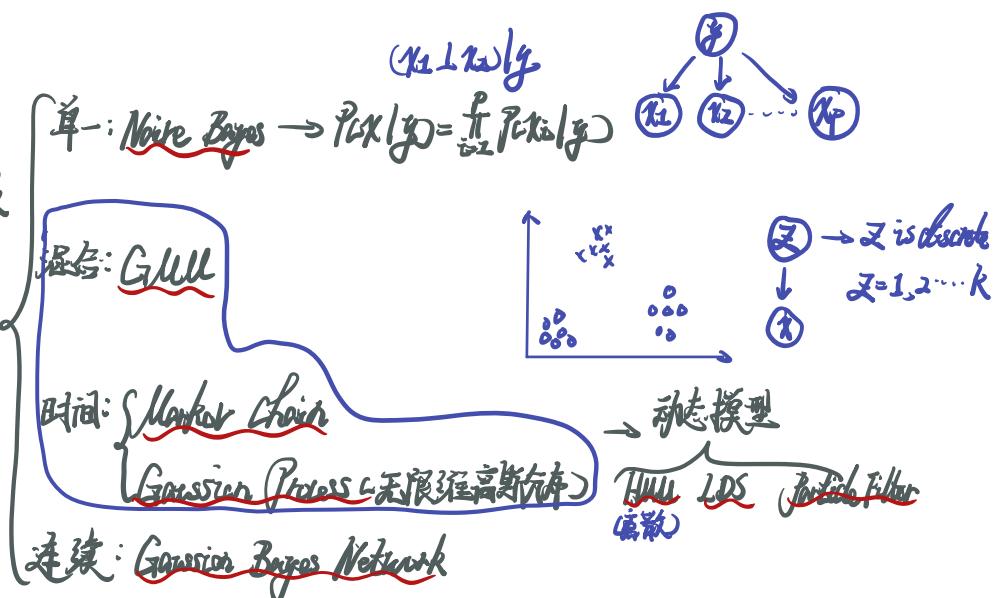




Ex:

Bayesian Network

从单一到动态
从有限到无限
① 空间 (离散 → 连续)
② 时间

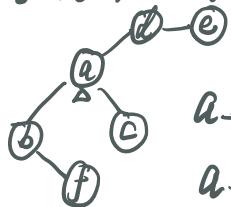


Markov network

$$\textcircled{1} X_A \perp X_C | X_B$$

Global Markov

$$\textcircled{2} \text{ local Markov}$$



$A \perp \text{集合}-a-\text{邻居}$

$a \perp \{e, f\} | \{b, c, d\}$

③ 局部 Markov

$$X_i \perp X_j | X_{-i-j} (\text{即 } j)$$

条件独立性体现在

三个方面: ① ② ③

$$\textcircled{1} \leftrightarrow \textcircled{2} \leftrightarrow \textcircled{3}$$

因子分解

图: 因, 最大因

$$P(X) = \frac{1}{Z} \prod_{i=1}^k \phi(X_{ci}) = \frac{1}{Z} \prod_{Q \in C} \psi_Q(X_Q)$$

$$Z = \sum \prod_{i=1}^k \phi(X_{ci})$$

$$= \sum_{Q \in C} \prod_{i \in Q} \phi(X_{ci})$$

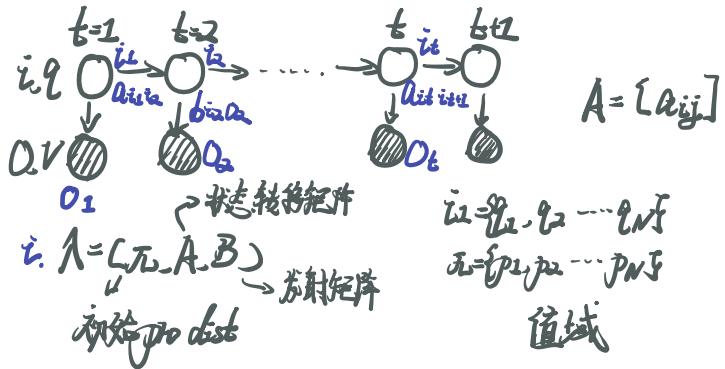
$$P(X) = \frac{1}{Z} \prod_{i=1}^k \phi(X_{ci})$$

C_i : 最大因, X_{ci} : 最大因随机变量集合

$$\phi(X_{ci}): \text{势函数, 必须为} + \rightarrow \phi(X_{ci}) = \exp\{-E(X_{ci})\}$$

$$Z = \sum \prod_{i=1}^k \phi(X_{ci}) = \sum_{Q_1} \sum_{Q_2} \dots \sum_{Q_k} \prod_{i \in Q} \phi(X_{ci}) \quad \hookrightarrow \text{Perceptron Distribution}$$

Hidder Markov Model (HMM)



观测变量 $O = O_1, O_2, \dots, O_t \rightarrow V = \{V_1, V_2, \dots, V_N\}$

状态变量 $i_1, i_2, i_3, \dots, i_t \rightarrow Q = \{q_1, q_2, \dots, q_{N+1}\}$

$$A = [a_{i_l i_j}], a_{i_l i_j} = P(i_{l+1} = q_j | i_l = q_i)$$

$$B = [b_{i_l O_t}], b_{i_l O_t} = P(O_t = V_k | i_l = q_i)$$

ii 两个假设:
 ①齐次Markov
 ②观测独立
 ↓
 无后效性

$$\textcircled{1} P(i_{t+1} | i_t, \dots, i_1, O_t, O_{t-1}, \dots, O_1) = P(i_{t+1} | i_t)$$

$$\textcircled{2} P(O_t | i_t, i_{t-1}, \dots, i_1, O_{t-1}, \dots, O_1) = P(O_t | i_t)$$

Evaluation: Given λ 求 $P(O|\lambda)$

$$P(O|\lambda) = \sum_I P(I, O|\lambda) = \sum_I P(O|I, \lambda) \cdot P(I|\lambda)$$

$$P(I|\lambda) = P(i_1, \dots, i_T|\lambda) = \underbrace{P(i_T | i_1, \dots, i_{T-1}, \lambda)}_{P(i_T | i_{T-1})} \cdot P(i_2, \dots, i_{T-1}|\lambda) = a_{i_{T-1} i_T} \cdot a_{i_{T-2} i_{T-1}} \cdots a_{i_1 i_2} \cdot \pi(i_1)$$

$$= \pi(i_1) \prod_{t=2}^T a_{i_{t-1} i_t}$$

$$P(O|I, \lambda) = \prod_{t=1}^T b_{i_t}(O_t)$$

$$\therefore P(O|\lambda) = \sum_I \pi(i_1) \prod_{t=2}^T a_{i_{t-1} i_t} \prod_{t=1}^T b_{i_t}(O_t) = \underbrace{\sum_{i_1} \sum_{i_2} \cdots \sum_{i_T} \pi(i_1) \prod_{t=2}^T a_{i_{t-1} i_t} \prod_{t=1}^T b_{i_t}(O_t)}$$

iii 三个问题

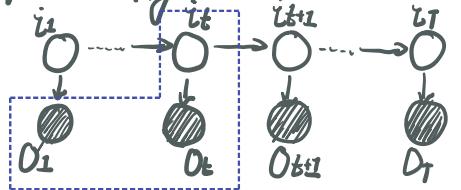
① Evaluation $\rightarrow P(O|\lambda) \rightarrow$ 前向向

② Learning \rightarrow 学习何求? \rightarrow EM
 Baum Welch

③ Decoding $\rightarrow \arg \max_{i_t} P(I|O)$

\hookrightarrow S 预测 $\rightarrow P(i_{t+1}|O_1, O_2, \dots, O_t)$
 滤波 $\rightarrow P(i_t|O_1, O_2, \dots, O_t)$

Forward Algorithm



$$\alpha_i(i_t) = P(O_1, \dots, O_t, i_t = q_i | \lambda)$$

$$\alpha_T(i_T) = P(O_1, \dots, O_T, i_T = q_i | \lambda)$$

$$P(O | \lambda) = \sum_{i=1}^N P(O, i_t = q_i | \lambda) = \sum_{i=1}^N \alpha_i(i_t)$$

$$\alpha_{t+1}(j) = P(O_1, \dots, O_t, O_{t+1}, i_{t+1} = q_j | \lambda)$$

$$= \sum_{i=1}^N P(O_1, \dots, O_t, O_{t+1}, i_{t+1} = q_j, i_t = q_i | \lambda)$$

$$= \sum_{i=1}^N P(O_{t+1} | O_1, \dots, O_t, i_t = q_i, i_{t+1} = q_j, \lambda)$$

$$P(O_1, \dots, O_t, i_t = q_i, i_{t+1} = q_j | \lambda)$$

$$= \sum_{i=1}^N P(O_{t+1} | i_{t+1}) \cdot P(O_1, \dots, O_t, i_t = q_i, i_{t+1} = q_j | \lambda)$$

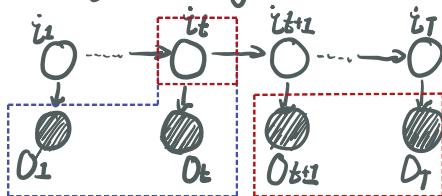
$$= \sum_{i=1}^N P(O_{t+1} | i_{t+1}) \cdot P(i_{t+1} = q_j | O_1, \dots, O_t, i_t = q_i, \lambda)$$

$$\cancel{P(O_1, \dots, O_t, i_t = q_i | \lambda)} \cdot P(i_{t+1} = q_j | i_t = q_i)$$

$$= \sum_{i=1}^N P(O_{t+1} | i_{t+1}) \cdot P(i_{t+1} = q_j | i_t = q_i) \cdot \alpha_i(i_t)$$

$$\alpha_j(O_{t+1}) \cdot \alpha_{ij} \cdot \alpha_i(i_t)$$

Backward Algorithm



$$\beta_{t+1}(i_t) = P(O_{t+1}, \dots, O_T | i_t = q_i, \lambda)$$

$$\beta_{t+1}(i_t) = P(O_2, \dots, O_T | i_t = q_i, \lambda)$$

$$P(O | \lambda) = P(O_1, \dots, O_T | \lambda)$$

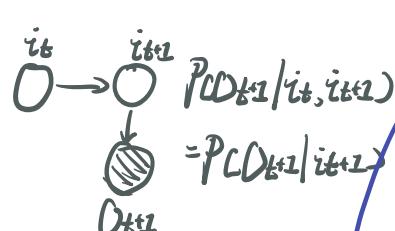
$$= \sum_{i=1}^N P(O_1, \dots, O_T, i_1 = q_i) \cancel{\pi_i}$$

$$= \sum_{i=1}^N P(O_1, \dots, O_T | i_1 = q_i) \cdot \cancel{P(i_1 = q_i)}$$

$$= \sum_{i=1}^N P(O_1 | i_1 = q_i) \cdot P(O_2, \dots, O_T | i_1 = q_i) \cdot \cancel{\pi_i}$$

$$= \sum_{i=1}^N P(O_1 | i_1 = q_i) \cdot \beta_{i+1}(i_1) \cdot \cancel{\pi_i}$$

$$= \sum_{i=1}^N b_{i+1}(O_1) \pi_i \beta_{i+1}(i_1)$$



$$\beta_{t+1}(i_t) = P(O_{t+1}, \dots, O_T | i_t = q_i, \lambda)$$

$$= \sum_{j=1}^N P(O_{t+1}, \dots, O_T, i_{t+1} = q_j | i_t = q_i) \alpha_{ij}$$

$$= \sum_{j=1}^N P(O_{t+1}, \dots, O_T | i_{t+1} = q_j, i_t = q_i) \cdot P(i_{t+1} = q_j | i_t = q_i)$$

$$= \sum_{j=1}^N P(O_{t+1} | O_{t+2}, \dots, O_T, i_{t+1} = q_j) \cdot P(O_{t+2}, \dots, O_T | i_{t+1} = q_j) \cdot \alpha_{ij}$$

$$= \sum_{j=1}^N P(O_{t+1} | i_{t+1} = q_j)$$

$$= \sum_{j=1}^N b_{j+1}(O_{t+1}) \cdot \alpha_{ij} \cdot \beta_{t+2}(q_j)$$

Learning

$$\lambda_{MLE} = \underset{\lambda}{\operatorname{argmax}} P(D|\lambda) \rightarrow \text{Bau-Welch}$$

$$\theta^{(t+1)} = \underset{\theta}{\operatorname{argmax}} \int_Z \log P(X, Z | \theta) \cdot P(Z | X, \theta^{(t)}) dZ$$

X: 观测 Z: 隐变量 θ: 参数
 ↓ ↓ ↓
 O I → 离散 λ

$$\lambda^{(t+1)} = \underset{\lambda}{\operatorname{argmax}} \sum_I \log P(D, I | \lambda) \cdot P(I | D, \lambda^{(t)})$$

$\frac{P(I | D, \lambda^{(t)})}{P(D | \lambda^{(t)})} \rightarrow \text{constant}$

$$\lambda^{(t+1)} = \underset{\lambda}{\operatorname{argmax}} \sum_I \log P(D, I | \lambda) \cdot P(I | \lambda^{(t)})$$

$$\lambda^{(t+1)} = (\pi^{(t)}, A^{(t)}, B^{(t)})$$

$$Q(\lambda | \lambda^{(t)}) = \sum_I \log P(D, I | \lambda) \cdot P(D, I | \lambda^{(t)})$$

$$= \sum_I \left[\log \pi_i^{(t)} + \sum_{t=1}^T \log a_{it} q_{it} + \sum_{t=1}^T b_{it} (O_{it}) \right] \cdot P(D, I | \lambda^{(t)})$$

$$\pi^{(t+1)} = \underset{\pi}{\operatorname{argmax}} Q(\lambda | \lambda^{(t)})$$

$$= \underset{\pi}{\operatorname{argmax}} \sum_I \left[\log \pi_i^{(t)} \cdot P(D, i_1 = q_1 | \lambda^{(t)}) \right]$$

$$= \underset{\pi}{\operatorname{argmax}} \sum_{i_1} \sum_{i_2} \cdots \sum_{i_T} \left[\log \pi_i^{(t)} \cdot P(D, i_1 = q_1, \dots, i_T = q_T | \lambda^{(t)}) \right] \xrightarrow{\text{将 } \sum_{i_2} \cdots \sum_{i_T} P(D, i_1 = q_1, \dots, i_T = q_T | \lambda^{(t)}) \text{ 消掉}}$$

$$= \underset{\pi}{\operatorname{argmax}} \sum_{i_1} \left[\log \pi_i^{(t)} \cdot P(D, i_1 = q_1 | \lambda^{(t)}) \right] \text{ s.t. } \sum_{i_1} \pi_i^{(t)} = 1$$

$$L(\pi, \eta) = \sum_{i_1} \left[\log \pi_i^{(t)} \cdot P(D, i_1 = q_1 | \lambda^{(t)}) \right] + \eta \left(\sum_{i_1} \pi_i^{(t)} - 1 \right)$$

$$\frac{\partial L}{\partial \pi} = \frac{1}{\pi^{(t)}} P(D, i_1 = q_1 | \lambda^{(t)}) + \eta = 0$$

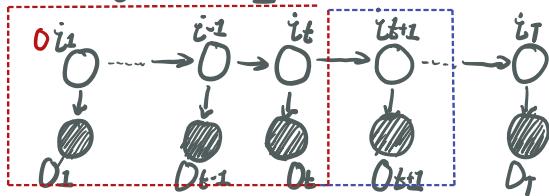
$$\sum_{i_1} \left[P(D, i_1 = q_1 | \lambda^{(t)}) + \pi_i^{(t)} \eta \right] = 0$$

$$P(D | \lambda^{(t)}) + \eta = 0$$

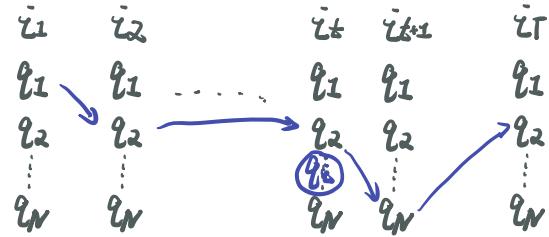
$$\therefore \eta = -P(D | \lambda^{(t)}) \Rightarrow \pi_i^{(t+1)} = \frac{P(D, i_1 = q_1 | \lambda^{(t)})}{P(D | \lambda^{(t)})}$$

$$\pi^{(t+1)} = (\pi_1^{(t+1)}, \dots, \pi_N^{(t+1)})$$

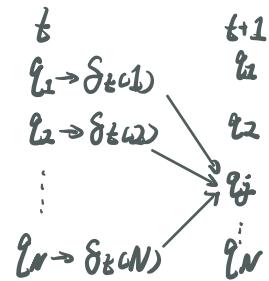
Decoding: $\hat{I} = \underset{I}{\operatorname{argmax}} P(I|D, A) \rightarrow \text{Viterbi}$



$$\delta_t(i_t) = \max_{i_1, i_2, \dots, i_{t-1}} P(D_1, D_2, \dots, D_t, i_1, i_2, \dots, i_{t-1}, i_t = q_{i_t})$$



$$\begin{aligned} \delta_{t+1}(i_j) &= \max_{i_1, i_2, \dots, i_t} P(D_1, D_2, \dots, D_t, D_{t+1}, i_1, i_2, \dots, i_{t-1}, i_t = q_{i_t}) \\ &= \max_{1 \leq i \leq N} \delta_t(i) \cdot a_{ij} \cdot b_j(D_{t+1}) \end{aligned}$$



$$\psi_{t+1}(j) = \underset{1 \leq i \leq N}{\operatorname{argmax}} \delta_t(i) \cdot a_{ij}$$

Markov Chain & Monte Carlo Method

Monte Carlo Method: 基于采样的随机近似法

$$\downarrow \begin{array}{c} P(z|x) \\ \text{latent data} \end{array} \xrightarrow{\text{observed data}} E_{z|x}[f(z)] = \int P(z|x) \cdot f(z) dx \\ \approx \frac{1}{N} \sum_{i=1}^N f(z_i)$$

Variational Inference

频率角度 → 优化问题 $\left\{ \begin{array}{l} \text{回归: } f(w) = w^T X, \text{ Loss function: } L(w) = \sum_{i=1}^N \|w^T X_i - y_i\|^2 \\ \text{模型: } \hat{w} = \arg \min_w L(w) \end{array} \right.$

$$D = \{(x_i, y_i)\}_{i=1}^N \sim \mathcal{N}(0, I)$$

$\left\{ \begin{array}{l} \text{解法: } \frac{\partial L}{\partial w} = 0 \Rightarrow \hat{w} = (X^T X)^{-1} X^T Y \\ \text{Algo: } \begin{cases} \text{① 解析解: } \\ \text{② 数值解: GD/SGD} \end{cases} \end{array} \right.$

SKM $\left\{ \begin{array}{l} f(w) = \text{Sgn}(w^T X + b) \\ \text{损失: } \dots \end{array} \right.$

$\left\{ \begin{array}{l} \text{② loss function: } \min \frac{1}{2} w^T w, \text{ s.t. } \sum_{i=1}^N y_i w^T X_i \geq 1 \end{array} \right.$

$\left\{ \begin{array}{l} \text{③ QP/Lagrange 对偶} \end{array} \right.$

$Z(\theta) = \theta = \arg \max_{\theta} \log P(X|\theta)$

$\theta^* = \arg \max_{\theta} \int \log P(X, Z|\theta) \cdot P(Z|X, \theta) dz$

Likelihood $P(X|\theta) \cdot P(\theta)$ prior

贝叶斯角度 → 损失问题 $P(\theta|X) = \frac{P(X|\theta) \cdot P(\theta)}{\int_\theta P(X|\theta) \cdot P(\theta) d\theta}$

Decision $X \rightarrow M$ 样本

$$\hat{x}, P(\bar{x}|X) = \int_\theta P(\bar{x}, \theta|X) d\theta = \int_\theta P(\bar{x}|\theta) \cdot \underbrace{P(\theta|X)}_{\text{posterior}} d\theta$$

$$= \mathbb{E}_{\theta|X}[P(\bar{x}|\theta)]$$

贝叶斯 Inference $\left\{ \begin{array}{l} \text{精确推断} \\ \text{近似推断} \end{array} \right.$

$\left\{ \begin{array}{l} \text{确定性近似} \rightarrow VI \\ \text{随机近似} \rightarrow MCMC, MH, Gibbs \end{array} \right.$

$\left\{ \begin{array}{l} \text{随机近似} \rightarrow MCMC, MH, Gibbs \end{array} \right.$

X : Observed Data

Z : Latent Variable + parameter

(X, Z) : Complete Data

$$\log P(X) = \log P(X, Z) - \log P(Z|X)$$

$$= \log \frac{P(X, Z)}{P(Z|X)} = \log \frac{P(X, Z)}{P(Z)}$$

左边: $\int_Z (\log P(X) q(Z)) dZ$

$$= \log P(X) \underbrace{\int_Z q(Z) dZ}_1$$

$$= \log P(X)$$

VI (mean field) \rightarrow Classical VI

Assumption: $q(Z) = \prod_{i=1}^n q_i(Z_i) \rightarrow$ independent

$$\log q_j(Z_j) = \sum_{i=1}^n q_i(Z_i) [\log P(X, Z_i | \theta)] + C$$

两边同时取期望

$$\text{左边: } \underbrace{\int_Z (\log P(X) q(Z)) dZ}_{\text{ELBO evidence lower Bound}} - \underbrace{\int_Z q(Z) \log \frac{P(X, Z)}{q(Z)} dZ}_{KLD q || P}$$

$$= L(q) + \underbrace{KL(q || P)}_{\text{变分}} \geq 0$$

$$\tilde{q}(Z) = \arg \max_{q(Z)} L(q) \Rightarrow \tilde{q}(Z) \approx P(Z | X)$$

$$q(Z) = \prod_{i=1}^n q_i(Z_i)$$

$$L(q) = \underbrace{\int_Z q(Z) \log P(X, Z) dZ}_1 - \underbrace{\int_Z q(Z) \log q(Z) dZ}_2$$

$$① = \int_Z \prod_{i=1}^n q_i(Z_i) \log P(X, Z) dZ_1 dZ_2 \dots dZ_n$$

$$= \int_{Z_j} q_j(Z_j) \left(\prod_{i \neq j} q_i(Z_i) \log P(X, Z) dZ_1 dZ_2 \dots dZ_n \right) dZ_j$$

$$= \int_{Z_j} q_j(Z_j) \left[\int_{Z_{j+1}} \dots \int_{Z_n} \log P(X, Z) \prod_{i \neq j} q_i(Z_i) dZ_i \right] dZ_j$$

$$= \underbrace{\int_{Z_j} q_j(Z_j) \cdot \sum_{i \neq j} q_i(Z_i) [\log P(X, Z)] \cdot dZ_i}_{\text{red}}$$

$$② = \int_Z q(Z) \log q(Z) dZ \quad \log \tilde{P}(X, Z)$$

$$= \int_Z \prod_{i=1}^n q_i(Z_i) - \sum_{i=1}^n \log q_i(Z_i) dZ$$

$$= \int_Z \prod_{i=1}^n q_i(Z_i) [\log q_1(Z_1) + \log q_2(Z_2) + \dots + \log q_n(Z_n)] dZ$$

$$= \sum_{i=1}^n \int_{Z_i} q_i(Z_i) \log q_i(Z_i) dZ_i$$

$$= \underbrace{\int_{Z_j} q_j(Z_j) \log q_j(Z_j) dZ_j}_{\text{red}} + C$$

$$\int_{Z_1} \prod_{i=2}^n q_i(Z_i) \log q_1(Z_1) dZ$$

$$= \int_{Z_1} q_1(Z_1) \log q_1(Z_1) dZ$$

$$= \int_{Z_1} \dots \int_{Z_n} q_1(Z_1) \log q_1(Z_1) dZ_1 \dots dZ_n$$

$$= \int_{Z_1} q_1(Z_1) dZ_1 \cdot \int_{Z_2} q_2(Z_2) \dots \int_{Z_n} q_n(Z_n) \dots$$

$$= \int_{Z_1} q_1(Z_1) dZ_1$$

$$\text{LHS} = \int_{z_j} q_j(z_j) \cdot \left(\log \frac{\hat{P}(x, z_j)}{q_j(z_j)} \right) dz_j + C$$
$$= -k L C q_j(z_j) \hat{P}(x, z_j) \leq 0$$

$$\therefore q_j(z_j) = \hat{P}(x, z_j)$$