

P39.40 Background

| | | | |
|-------------|---|-------------------------------------|--|
| Features | Sufficient Statistic ① | Gaussian | All Exponential Family of Distribution |
| | Conjugate ② | Bernoulli \rightarrow Categorical | |
| | Maximum Entropy (Noninformative priors) | Binomial \rightarrow Multinomial | |
| | | Poisson | |
| Application | Generalised Linear Model | Dirichlet | |
| | Probability Graph | | |
| | Variational Inference | | |

$$\text{Standard Form: } y(x|\eta) = h(x) \exp\{\eta^T f(x) - A(\eta)\}$$

η : parameter vector, $x \in \mathbb{R}^p$ $A(\eta)$: log partition function (log 配分函数) $f(x)$: Sufficient Statistic

$$\therefore P(x|\eta) = \frac{1}{\exp\{A(\eta)\}} h(x) \exp\{\eta^T f(x)\} \quad \because \int P(x|\eta) dx = 1$$

$$\therefore \frac{1}{\exp\{A(\eta)\}} \int h(x) \exp\{\eta^T f(x)\} dx = 1 \quad \therefore \exp\{A(\eta)\} \text{ is partition function (can be simply understood as a normalization factor)}$$

$$\therefore A(\eta) = \log \exp\{A(\eta)\} \quad \therefore A(\eta) \text{ is log partition function}$$

① 包含原样本中全部有用信息, 有①可丢弃原样本

如满足 Gaussian 分布一组样本, $f(x) = \begin{bmatrix} \sum_{i=1}^N x_i \\ \sum_{i=1}^N x_i^2 \end{bmatrix}$ 即为充分统计量

② e.g.: $P(z|x) = P(x|z)P(z)$ If likelihood is EFD, prior and posterior are identically distribution.

③ Initialize Prior

P41 Gaussian Distribution

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

\downarrow

$$P(x|\theta) = h(x) \exp\{\eta^T f(x) + A(\eta)\}$$

$$\# P(x, \mu, \sigma^2) = \frac{1}{\sqrt{2\pi}\sigma\sqrt{2\pi}\sqrt{\sigma^2}} \exp\left\{-\frac{1}{2(1+\sigma^2)}\left[\frac{(x_1-\mu)^2}{\sigma^2} - \frac{2(x_1-\mu)(x_2-\mu)}{\sigma_1\sigma_2} + \frac{(x_2-\mu)^2}{\sigma_2^2}\right]\right\}$$

$$P(x_1, \dots, x_N|\theta) = \frac{1}{\sqrt{2\pi}}^N |\Sigma|^{-\frac{1}{2}} \exp\left\{-\frac{1}{2}(x-\mu)^T \Sigma^{-1}(x-\mu)\right\}$$

$$\begin{aligned} \frac{1}{\sqrt{2\pi}\sigma} \exp\left\{-\frac{(x-\mu)^2}{2\sigma^2}\right\} &= \exp\{\log(\frac{1}{\sqrt{2\pi}\sigma})\} \exp\left\{-\frac{1}{2\sigma^2}(x^2 - 2\mu x + \mu^2)\right\} = \exp\left\{-\frac{1}{2}\log 2\pi\sigma^2\right\} \exp\left\{-\frac{1}{2\sigma^2}(-2\mu)(x) - \frac{\mu^2}{2\sigma^2}\right\} \\ &= \exp\left\{\left(\frac{\mu}{\sigma^2} - \frac{1}{2\sigma^2}\right)x - \left(\frac{\mu^2}{2\sigma^2} + \frac{1}{2}\log 2\pi\sigma^2\right)\right\} \end{aligned}$$

$$\eta = \begin{pmatrix} \eta_1 \\ \eta_2 \end{pmatrix} = \begin{pmatrix} \frac{\mu}{\sigma^2} \\ -\frac{1}{2\sigma^2} \end{pmatrix} \Rightarrow \begin{cases} \sigma^2 = -\frac{1}{2\eta_2} \\ \mu = -\frac{\eta_1}{2\eta_2} \end{cases} \quad A(\eta) = -\left(\frac{\mu^2}{2\sigma^2} + \frac{1}{2}\log 2\pi\sigma^2\right) = -\frac{\eta_1^2}{4\eta_2} + \frac{1}{2}\log\left(-\frac{2}{\eta_2}\right)$$

$$\text{EFD} = \exp\{\eta^T f(x) - A(\eta)\}$$

P42 Log Partition Function

$$P(x|\eta) = \frac{1}{\exp\{A(\eta)\}} h(x) \exp\{\eta^T f(x)\} \quad \because \int P(x|\eta) dx = 1 \quad \therefore \exp\{A(\eta)\} = \int h(x) \exp\{\eta^T f(x)\} dx$$

$$\text{Partial derivative of } \eta \text{ on both side of equation: } A'(\eta) \exp\{A(\eta)\} = \frac{\partial}{\partial \eta} \left(\int h(x) \exp\{\eta^T f(x)\} dx \right) = \int h(x) \exp\{\eta^T f(x)\} f(x) dx$$

$$\therefore A'(\eta) = \int h(x) \exp\{\eta^T f(x) - A(\eta)\} f(x) dx$$

Derivative and integral are interchangeable

$$= \int P(x|\eta) f(x) dx = E[f(x)]$$

$$A''(\eta) = \int h(x) \exp\{\eta^T f(x) - A(\eta)\} f(x) [f(x) - A'(\eta)] dx$$

$$= \int P(x|\eta) f(x)^2 dx - A'(\eta) \int P(x|\eta) f(x) dx \quad ? \int P(x) x^2 dx = E(x^2)?$$

$$= E[f(x)^2] - E[f(x)]^2 = \text{Var}[f(x)]$$

P43 MLE & Sufficient Statistics

$$\eta_{\text{MLE}} = \arg \max_{\eta} \sum_{i=1}^N \log P(x_i|\eta) = \arg \max_{\eta} \sum_{i=1}^N \log [h(x_i) \exp\{\eta^T f(x_i) - A(\eta)\}]$$

$$= \arg \max_{\eta} \sum_{i=1}^N [\eta^T f(x_i) - A(\eta)]$$

$$\therefore \frac{\partial}{\partial \eta} \left\{ \sum_{i=1}^N [\eta^T f(x_i) - A(\eta)] \right\} = \sum_{i=1}^N f(x_i) - N A'(\eta) = 0$$

$$\therefore A'(\eta_{\text{MLE}}) = \frac{1}{N} \sum_{i=1}^N f(x_i) \Rightarrow \eta_{\text{MLE}} = A'^{-1}(\eta_{\text{MLE}}) \text{ Inverse Function}$$

P44/45 Maximum Entropy Perspective

$$H(P) \text{ (Entropy)} = \begin{cases} -\int P(x) \log P(x) dx \\ -\sum P(x_i) \log P(x_i) \end{cases}$$

$$\begin{cases} \max \sum_{i=1}^N -P_i \log P_i = \min \sum_{i=1}^N P_i \log P_i \\ \text{s.t. } \sum_{i=1}^N P_i = 1 \end{cases}$$

$$\begin{cases} \max \sum_{i=1}^N -p_i \log p_i = \min \sum_{i=1}^N \bar{p}_i \log \bar{p}_i \\ \text{s.t. } \sum_{i=1}^N p_i = 1 \end{cases}$$

Lagrange Multiplier Method: $L(p_i, \lambda) = \sum_{i=1}^N p_i \log p_i + \lambda (1 - \sum_{i=1}^N p_i)$

$$\frac{\partial L}{\partial p_i} = \log p_i + p_i \frac{1}{p_i} - \lambda = 0 \Rightarrow p_i = \exp\{1 - \lambda\} \quad \lambda \text{ is hyperparameter } \therefore p_i \text{ is constant}$$

$$\therefore p_1 = p_2 = \dots = p_N = \frac{1}{N}$$

Maximum Entropy + known facts (constraints) \Rightarrow Principle of Maximum Entropy

(经验分布)

Data: $\{x_1, x_2, \dots, x_N\}$ obey Empirical Distribution $\hat{p}(x=c) = \hat{p}(x) = \frac{\text{count}(c)}{N}$

$\therefore E_{\hat{p}}[x]$, $\text{Var}_{\hat{p}}[x]$ are known # P is the overall distribution; \hat{p} is the sample distribution

$\therefore E_{\hat{p}}[g(x)]$ is known # $g(x)$ is an arbitrary vector of functions of x & $g(x)$ is a sufficient statistic of X

let $E_{\hat{p}}[g(x)] = \Delta$

$$\therefore \begin{cases} \min \sum_{x \in \mathcal{X}} P(x) \log P(x) \\ \text{s.t. } \sum_{x \in \mathcal{X}} P(x) = 1 \\ E_P[g(x)] = E_{\hat{p}}[g(x)] = \Delta \end{cases}$$

Lagrange: $L(P, \lambda_0, \lambda) = \sum_{x \in \mathcal{X}} P(x) \log P(x) + \lambda_0 [1 - \sum_{x \in \mathcal{X}} P(x)] + \lambda^T [\Delta - E_P[g(x)]]$ # $E_P[g(x)] = \sum_{x \in \mathcal{X}} P(x) g(x)$

$$\frac{\partial L}{\partial P(x)} = \sum_{x \in \mathcal{X}} [\log P(x) + 1 - \lambda_0 - \lambda^T g(x)] = 0 \quad \therefore \log P(x) + 1 - \lambda_0 - \lambda^T g(x) = 0 ?$$

$$\therefore \log P(x) = \lambda^T g(x) + \lambda_0 - 1 \Rightarrow P(x) = \exp \left\{ \underbrace{\lambda^T g(x)}_{f(x)} - (1 - \lambda_0) \right\}$$