

指数族分布

1. Introduction.

- 指数族分布是一类分布. 如

Gauss Distribution. $f(x) = \frac{1}{\sqrt{2\pi}\sigma} \exp\left(-\frac{(x-\mu)^2}{2\sigma^2}\right)$

Bernoulli $f(x|p) = \begin{cases} p^x q^{1-x} & x=0,1 \\ 0 & x \neq 0,1 \end{cases}$

....

- 指数族分布可以写成统一的形式:

$$\begin{aligned} p(x|\eta) &= h(x) \exp(\eta^T \phi(x) - A(\eta)) \\ &= \frac{1}{\exp(A(\eta))} h(x) \exp(\eta^T \phi(x)) \end{aligned}$$

其中, η 是参数向量, $x \in \mathbb{R}^P$

$A(\eta)$: log partition function (可看作归一化因子)

$\phi(x)$: 充分统计量 (包含样本点的所有信息, 有它可代替样本点, 可理解为“人大代表”)

Explanation:

$$p(x|\theta) = \frac{1}{Z} \hat{p}(x|\theta)$$

对左边式子, $\exp(A(\eta))$ 看作 Z ,
 $h(x) \exp(\eta^T \phi(x))$ 看作 $\hat{p}(x|\theta)$.

$$R|A(\eta) = \log(Z)$$

Example

Gauss Distribution 的指数族分布形式

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

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

$$= \exp \log(2\pi\sigma^2)^{-\frac{1}{2}} \exp\left\{-\frac{1}{2\sigma^2} \underline{(-2\mu - 1)} \begin{pmatrix} x \\ x^2 \end{pmatrix} - \frac{\mu^2}{2\sigma^2}\right\}$$

$$= \exp \left\{ \frac{(\frac{\mu}{\sigma^2} - \frac{1}{2\sigma^2})}{\eta^T} \begin{pmatrix} x \\ x^2 \end{pmatrix} - \frac{\mu^2}{2\sigma^2} + \frac{1}{2} \cdot \log 2\pi\sigma^2 \right\} \quad (1)$$

$$\because \eta^T = \frac{\mu}{\sigma^2} - \frac{1}{2\sigma^2}, \text{ 令 } \eta = \begin{bmatrix} \frac{\mu}{\sigma^2} \\ -\frac{1}{2\sigma^2} \end{bmatrix} \Rightarrow \eta_1 = \frac{\mu}{\sigma^2}, \eta_2 = -\frac{1}{2\sigma^2}$$

$$\Rightarrow \begin{cases} \mu = -\frac{\eta_1}{2\eta_2} \\ \sigma^2 = -\frac{1}{2\eta_2} \end{cases}$$

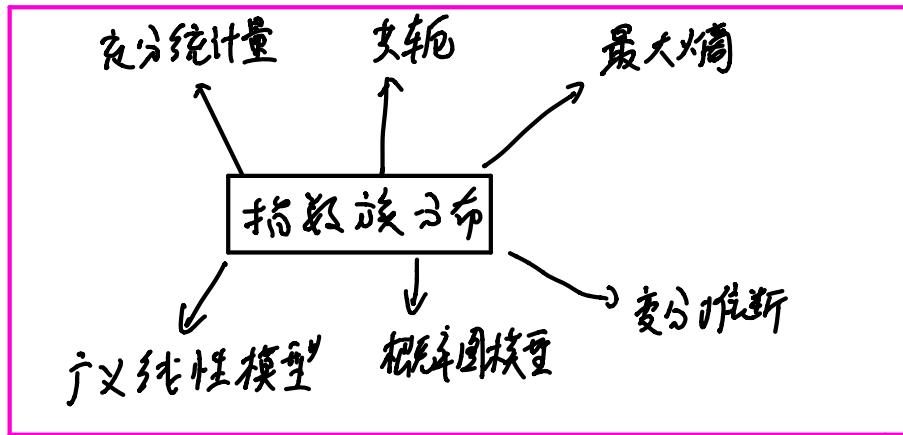
$$e^a \cdot e^b = e^{a+b}$$

$$\begin{aligned}
 A(\eta) \text{ 可由 -\frac{\eta^2}{2\sigma^2} + \frac{1}{2} \cdot \log 2\pi\sigma^2 \text{ 写成:}} \\
 &= -\frac{\eta_1^2}{4\eta_2^2} - \frac{1}{2} \log (2\pi \cdot \frac{1}{2\eta_2}) \\
 &= \frac{\eta_1^2}{4\eta_2} - \frac{1}{2} \log (\frac{\pi}{\eta_2})
 \end{aligned}$$

\therefore Gauss Distribution 的对数似然分布可写成:

$$\begin{aligned}
 \eta = \begin{pmatrix} \eta_1 \\ \eta_2 \end{pmatrix}, \quad \phi(x) = \begin{pmatrix} x \\ x^2 \end{pmatrix}, \\
 A(\eta) = \frac{\eta_1^2}{4\eta_2} - \frac{1}{2} \log \left(\frac{\pi}{\eta_2} \right) \\
 \Rightarrow \exp(\eta^T \phi(x) - A(\eta)) \\
 \boxed{p(x|\eta) = h(x) \exp(\eta^T \phi(x) - A(\eta))}
 \end{aligned}$$

· 指数族分布的性质和应用



2. $p(x|\eta) = h(x) \exp(\eta^T \phi(x) - A(\eta))$ 中, $\phi(x)$ 与 $A(\eta)$ 的关系:

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

$$= \frac{1}{\exp(A(\eta))} h(x) \exp(\eta^T \phi(x))$$

explanation:

$$\int p(x|\eta) dx = 1$$

$$\exp(A(\eta)) = \int h(x) \exp(\eta^T \phi(x)) dx$$

对 η 进行求导：

$$\begin{aligned} \exp(A(\eta)) \cdot A'(\eta) &= \frac{\partial}{\partial \eta} \left(\int h(x) \exp(\eta^T \phi(x)) dx \right) \\ &= \int h(x) \underbrace{\exp(\eta^T \phi(x))}_{\text{对模型部分的 } \eta \text{ 进行求导!}} \cdot \phi(x) dx \\ \Rightarrow A'(\eta) &= \frac{\int h(x) \underbrace{\exp(\eta^T \phi(x))}_{\text{对模型部分的 } \eta \text{ 进行求导!}} \cdot \phi(x) dx}{\exp(A(\eta))} \\ &= \int \underbrace{h(x) \exp(\eta^T \phi(x) - A(\eta))}_{\text{该类型是期望定义: } E(x) = \int p(x) f(x) dx, \text{ 其中 }} \cdot \phi(x) dx \\ &= E[\phi(x)] \end{aligned}$$

即有 $A'(\eta) = E[\phi(x)]$

$A''(\eta) \approx \text{Var}[\phi(x)]$

3. 极大似然估计与充分统计量

$$D = \{x_1, x_2, \dots, x_n\}$$

$$p(x|\eta) = h(x) \exp\{\eta^T \phi(x) - A(\eta)\}$$

对 $p(x|\eta)$ 进行极大似然估计：

$$\begin{aligned} \eta_{MLE} &= \arg \max_{\eta} \log p(x|\eta) \\ &= \arg \max_{\eta} \sum_{i=1}^N \log p(x_i|\eta) \\ &= \arg \max_{\eta} \sum_{i=1}^N \log [h(x_i) \exp\{\eta^T \phi(x_i) - A(\eta)\}] \\ &= \arg \max_{\eta} \sum_{i=1}^N [\log h(x_i) + \eta^T \phi(x_i) - A(\eta)] \\ &= \arg \max_{\eta} \sum_{i=1}^N [\eta^T \phi(x_i) - A(\eta)] \end{aligned}$$

$$\frac{\partial p(x|\eta)}{\partial \eta} = 0 \Rightarrow \frac{\partial}{\partial \eta} \sum_{i=1}^N [\eta^T \phi(x_i) - A(\eta)]$$

$$= \sum_{i=1}^N \phi(x_i) - A'(\eta)$$

$$= \sum_{i=1}^N \phi(x_i) - N A'(\eta) = 0$$

$$\Rightarrow A'(\eta) = \frac{1}{N} \sum_{i=1}^N \phi(x_i)$$

结论：为了估计参数，只需要知道充分统计量就可以。

4. 最大熵角度

• 几个概念

① 信息量： $-\log P$

② 熵： $E_{\text{P(x)}}(-\log P) = \int -P(x) \cdot \log P(x) dx$ (连续型)

$$\downarrow \quad = - \sum_x P(x) \log P(x) \quad (\text{离散型})$$

$$H[P] = - \sum_x P(x) \log P(x)$$

③ **最大熵 \Leftrightarrow 等分布**

• Example (最大熵)

假设 x 是离散的

x	1	2	...	K	
P	P_1	P_2	...	P_K	, $\sum_{i=1}^K P_i = 1$

求最大熵，则有如下离散 (x_i) 问题

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

写成拉格朗日函数

$$L(P, \lambda) = \sum_{i=1}^K P_i \log P_i + \lambda \left(1 - \sum_{i=1}^K P_i \right)$$

$$\frac{\partial L}{\partial P_i} = \log P_i + P_i \cdot \frac{1}{P_i} - \lambda = 0$$

$$\Rightarrow \hat{P}_i = e^{\lambda-1} \quad (\text{常数})$$

$$\therefore P_1 = P_2 = \dots = P_K = \frac{1}{K}$$

参考书：《统计学习方法》P94

《李航》读书指导 P6 30分钟

$p(x)$ 是均匀分布 \rightarrow 等可能

最大熵原理：满足约束条件的模型集合中选取熵最大的模型。

④ 在实际问题中，求最大熵需满足已知事实（数据）

$$Data = \{x_1, x_2, \dots, x_n\}$$

满足：经验分布 $\hat{p}(x=x) = \hat{P}(x) = \frac{\text{count}(x)}{N}$

则易求得 $E_{\hat{p}}[x]$, $\text{Var}_{\hat{p}}[x]$. 特别地 $f(x)$ 关于经验分布 \hat{p} 的期望值
设 $f(x)$ 是任意关于 x 的函数. $E_{\hat{p}}[f(x)] = \Delta \rightarrow$ 已知

加上约束后，最大熵问题成：

$$f(x) = \begin{bmatrix} x_1 \\ x_2 \\ \vdots \\ x_n \end{bmatrix}, \quad \Delta = \begin{bmatrix} \Delta_1 \\ \Delta_2 \\ \vdots \\ \Delta_n \end{bmatrix}$$

$$\left\{ \begin{array}{l} \min_x p(x) \log p(x) \\ \text{s.t. } \sum_x p(x) = 1 \end{array} \right.$$

松弛

$p(x)$ 是拉格朗日分布

$$E_p[f(x)] = E_{\hat{p}}[f(x)] = \Delta$$

真实分布 = 训练集采样分布

拉格朗日形式：

$$L(p(x), \lambda_0, \lambda) = \sum_x p(x) \log p(x) + \lambda_0 (1 - \sum_x p(x)) + \lambda^T (\Delta - E_p[f(x)])$$

$$\frac{\partial L}{\partial p(x)} = \sum_x \log p(x) + p(x) \cdot \frac{1}{p(x)} - \sum_x \lambda_0 - \sum_x \lambda^T f(x) = 0$$

$$\Rightarrow \sum_x (\log p(x) + 1 - \lambda_0 - \lambda^T f(x)) = 0$$

$$\Rightarrow \log p(x) + 1 - \lambda_0 - \lambda^T f(x) = 0$$

$$\Rightarrow p(x) = \exp(\lambda^T f(x) + \lambda_0 - 1)$$

拉格朗日分布