

指数族分布

Exponential Family Distribution

- ① 背景
- ② 高斯分布
- ③ 对数配分函数
- ④ 极大似然估计
- ⑤ 最大熵角度

背景 (Background)

高斯、伯努利、类别、二项、多项式、Beta、Dirichlet、Gamma

上述的分布都属于指数族分布。

指数族分布的标准形式：

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

其中， $x \in \mathbb{R}^p$

η : 参数，是一个向量

$A(\eta)$: log partition function
配分函数。

配分函数：

一般地， $P(x|\theta) = \frac{1}{Z} \hat{p}(x|\theta)$

其中 $\hat{p}(x|\theta)$ 是关于 x 的函数，其积分不为 1。

所以为了满足概率密度的定义，加上一个归一化因子 Z 。

此时，
$$\int \hat{p}(x|\theta) dx = \int \frac{1}{Z} \hat{p}(x|\theta) dx$$

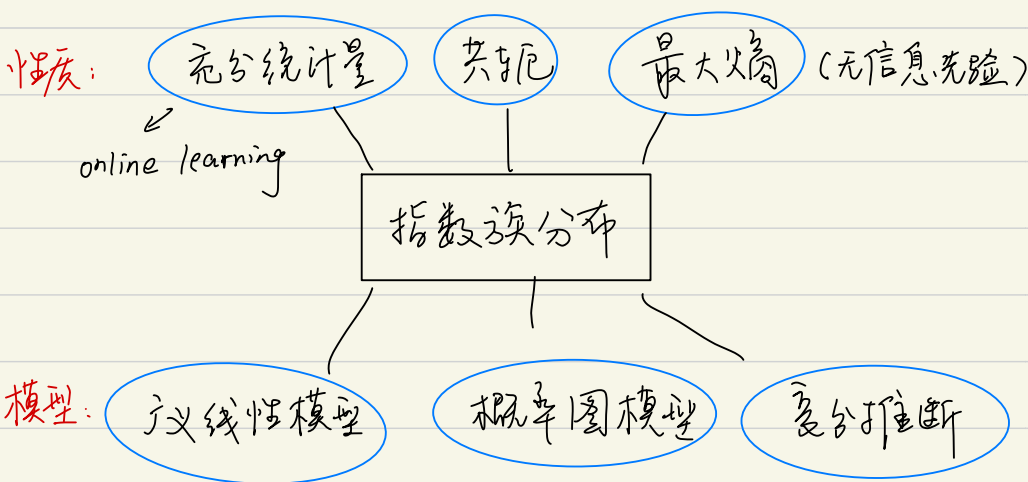
则
$$Z = \int \hat{p}(x|\theta) dx$$
 称 Z 为配分函数。

(有点类似于能量模型)

例|对于指数族的标准形式:

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

例| $Z = \exp\{A(\eta)\} \Rightarrow A(\eta) = \log Z$. 故称之为 log partition function



充分统计量: 指的就是标准形式中的 $\phi(x)$.

充分统计量的作用是压缩数据, 不需要储存样本, 只需要储存充分统计量就可以表示分布.

共轭:
$$p(z|x) = \frac{p(x|z)p(z)}{\int_z p(x|z)p(z)dz}$$

往往上式这个后验是很难求的。

原因在于 $\int_z p(x|z)p(z)dz$ 很难求。

主要是由于分布复杂或数据维度过高

所以需要用到近似的方法。

有变分推断、MCMC 等。

但是也可以用共轭的方法来使得后验可计算。

$$p(z|x) \propto p(x|z)p(z)$$

如果后验与先验具有相同的分布形式

那么称先验与似然, 是共轭的

例如

$$p(z|x) \propto p(x|z)p(z)$$

Beta = 项式 Beta.

最大熵 (无信息先验): 熵越大, 表示不确定性越大。

假设当不知道参数的信息时, 所有可能发生的
情况都是等可能的, 即熵最大。

确定先验分布的几种方法:

① 共轭 \rightarrow 计算方便

② 最大熵 \rightarrow 无信息先验

③ Jeffrey 先验 \rightarrow 无信息先验

广义线性模型：三个最基本的概念：

① 线性组合 $w^T x$

② Link function 激活函数的反函数。

③ 指数族分布： $y|x$ 满足指数族分布。

例子：线性回归： $y|x \sim \mathcal{N}(\mu, \Sigma)$

分类： $y|x \sim$ 伯努利分布。

泊松回归： $y|x \sim$ 泊松分布

概率图模型：无向图：RBM (限制玻尔兹曼机) \rightarrow 需要满足指数族分布

变分推断：指数族分布可以简化变分推断的运算。

高斯分布 (Gaussian Distribution)

标准指数族分布形式:

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

以一维高斯分布为例, 将其化为标准指数族分布的形式:

$$\theta = (\mu, \sigma^2)$$

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

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

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

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

$$\text{令} \left\{ \begin{array}{l} h(x) = 1 \\ \eta = \begin{pmatrix} \eta_1 \\ \eta_2 \end{pmatrix} = \begin{pmatrix} \frac{\mu}{\sigma^2} \\ -\frac{1}{2\sigma^2} \end{pmatrix} \Rightarrow \begin{cases} \mu = -\frac{\eta_1}{2\eta_2} \\ \sigma^2 = -\frac{1}{2\eta_2} \end{cases} \\ \phi(x) = \begin{pmatrix} x \\ x^2 \end{pmatrix} \\ A(\eta) = -\frac{\mu^2}{2\sigma^2} - \frac{1}{2}\log(2\pi\sigma^2) = \frac{\eta_1^2}{4\eta_2} - \frac{1}{2}\log\left(-\frac{\pi}{\eta_2}\right) \end{array} \right\} \Rightarrow \underline{h(x) \exp\{\eta^T \phi(x) - A(\eta)\}}$$

对数配分函数

(Log Partition Function)

标: 指数族分布形式.

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

两边同时对 x 求积分:

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

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

两边同时对 η 求导:

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

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

$$A'(\eta) = \int h(x) \exp\{\eta^T \phi(x) - A(\eta)\} \cdot \phi(x) dx.$$

$$A'(\eta) = E_{x \sim p(x|\eta)} [\phi(x)]$$

由上式，易得：

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

两边对 η 求导：

$$\exp\{A(\eta)\} [A'(\eta)]^2 + \exp\{A(\eta)\} A''(\eta) = \int h(x) \exp\{\eta^T \phi(x)\} \phi^2(x) dx$$

$$[A'(\eta)]^2 + A''(\eta) = \int h(x) \exp\{\eta^T \phi(x) - A(\eta)\} \cdot \phi^2(x) dx$$
$$A''(\eta) = E_{x \sim p(x|\eta)} [\phi^2(x)] - (E_{x \sim p(x|\eta)} [\phi(x)])^2$$

$$A''(\eta) = \text{Var}_{x \sim p(x|\eta)} [\phi(x)]$$

$$\sim | \quad A'(\eta) = E_{x \sim p(x|\eta)} [\phi(x)]$$

$$A''(\eta) = \text{Var}_{x \sim p(x|\eta)} [\phi(x)]$$

极大似然估计 (MLE)

上一节的推导是直接根据标准形式确定的。本节将使用样本来估计 η 。

标准指数族分布的形式：

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

数据 $D = \{x_1, x_2, \dots, x_N\}$

$$\eta_{MLE} = \arg\max_{\eta} \log P(D|\eta)$$

$$= \arg\max_{\eta} \log \prod_{i=1}^N p(x_i|\eta)$$

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

由于 $h(x_i)$ 与 η 无关，故舍去：

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

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

$$\Rightarrow A'(\eta_{MLE}) = \frac{1}{N} \sum_{i=1}^N \phi(x_i), \quad \text{记 } A^{(-1)'(\cdot)} \text{ 为 } A'(\cdot) \text{ 的反函数.}$$

$$\text{则 } \eta_{MLE} = A^{(-1)'(\cdot)}(\eta_{MLE})$$

最大熵角度

(Max Entropy Perspective)

对未知概率分布二猜测

最大熵 \Leftrightarrow 等可能.

信息量: $-\log p$.

熵: $E_{x \sim p(x)} [-\log p] = \int -p(x) \cdot \log p(x) dx$ 连续

记为 HLP

$$= -\sum_x p(x) \log p(x)$$

离散.

假设 x 是离散的. 分布为:

x	1, 2, ..., K
p	p_1, p_2, \dots, p_k

$\sum_{i=1}^K p_i = 1$

2. 最大熵问题就是一个优化问题

$$\begin{cases} \max_p \text{HLP} = \max - \sum_{i=1}^K p_i \log p_i \\ \text{s.t. } \sum_{i=1}^K p_i = 1 \end{cases}$$

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

拉格朗日函数: $L(p, \lambda) = \sum_{i=1}^K p_i \log p_i + \lambda (1 - \sum_{i=1}^K p_i)$

$$\frac{\partial L}{\partial p_i} = \log p_i + 1 - \lambda = 0 \Rightarrow p_i = \exp\{\lambda - 1\}.$$

$$\text{即 } p_1 = p_2 = \dots = p_k.$$

$$\text{又因 } \sum_{i=1}^K p_i = 1$$

$$\text{2.1 } p_1 = p_2 = \dots = p_k = \frac{1}{k} \quad \text{均匀分布}$$

在没有任何已知情况下,
均匀分布的熵最大

最大熵原理

满足已知事实的最大熵，数据服从指数族分布。

已知事实: $\text{Data} = \{x_1, x_2, \dots, x_N\}$

通过经验分布量化该已知事实

$$\text{经验分布: } \hat{p}(X=x) = \hat{p}(x) = \frac{\sum_i \mathbb{1}(x_i=x)}{N} \quad \mathbb{1}(\text{condition}) = \begin{cases} 1, & \text{condition} = \text{true} \\ 0, & \text{condition} = \text{false} \end{cases}$$

empirical distribution

则可根据该分布求得 $E[X]$, $\text{Var}[X]$ 等数字特征。

则定义 $E_p[f(x)] = \Delta \rightarrow \text{已知}$ 。

$$\text{其中: } f(x) = \begin{pmatrix} f_1 \\ f_2 \\ \vdots \\ f_q \end{pmatrix}, \quad \Delta = \begin{pmatrix} \Delta_1 \\ \Delta_2 \\ \vdots \\ \Delta_q \end{pmatrix}$$

$$\begin{array}{c|c} & x \\ \hline & p_1 \ p_2 \ p_3 \end{array}$$

则根据最大熵理论，得到优化问题是

$$\begin{cases} \min \sum_x p(x) \log p(x) \\ \text{s.t. } \sum_x p(x) = 1 \\ E_p[f(x)] = E_{\hat{p}}[f(x)] = \Delta \end{cases}$$

$$\text{Lagrange 乘数: } \mathcal{L}(p, \lambda_0, \lambda) = \sum_x p(x) \log p(x) + \lambda_0 (1 - \sum_x p(x)) + \lambda^T (\Delta - E_p[f(x)])$$

对 $p(x)$ 求偏导

$$\begin{aligned}\frac{\partial \mathcal{L}}{\partial p(x)} &= \sum_x (\log p(x) + 1) - \sum_x \lambda_0 - \sum_x \lambda^T f(x) \quad \text{这一步需要仔细思考一下.} \\ &= \sum_x [\log p(x) + 1 - \lambda_0 - \lambda^T f(x)] \\ &\triangleq 0\end{aligned}$$

这里注意，假设没将 $p(x)$ 看作是离散，则对于 x 的每个取值对应的 p_i ，其拉格朗日函数的偏导数都为 0（见该小节第一页），则

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

蓝字这里还需要具体验证一下，
但不知道哪里有理论证明。

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

$$\text{令 } \eta = \begin{pmatrix} \lambda_0 \\ \lambda \end{pmatrix}, \quad \phi(x) = \begin{pmatrix} 1 \\ f(x) \end{pmatrix}, \quad A(\eta) = (1, 0) \cdot \eta + 1, \quad h(x) =$$

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

\therefore 满足已知事实的最大熵所得到的数据分布就是指数族分布。

可以看一下最大熵模型

同时 PRML 中，还用指数族分布引出了 sigmoid 和 softmax，也可以看看