

Gamma函数

$$\Gamma(x) = \int_0^{\infty} t^{x-1} e^{-t} dt$$

性质 $\Gamma(x+1) = x\Gamma(x)$

于是很容易证出 $\Gamma(n) = (n-1)!$???

即阶乘在实数集上的拓展

推导历史:

1728 哥德巴赫处理阶乘序列 1, 2, 6, 24, 120...

我们可以算 $2!, 3!$ 那 $2.5!$ 怎么算?

哥一写信请教丹尼尔·贝努利 欧拉也在

1729年 欧拉22岁 完美解决了此题

Gamma分布

$$\Gamma(x) = \int_0^{\infty} t^{x-1} e^{-t} dt$$

变形 $\int_0^{\infty} \frac{x^{\alpha-1} e^{-x}}{\Gamma(\alpha)} dx = 1$

取积分中的函数作概率密度

则 $\text{Gamma}(x|\alpha) = \frac{x^{\alpha-1} e^{-x}}{\Gamma(\alpha)}$

令 $x = \beta t$ 得到更一般形式

$$\text{Gamma}(t|\alpha, \beta) = \frac{\beta^{\alpha} t^{\alpha-1} e^{-\beta t}}{\Gamma(\alpha)}$$

α 决定了分布曲线的形状

β 决定了曲线有多陡

Gamma分布作为先验分布很强

指数分布, χ^2 分布都是特殊的 Gamma 分布

指数分布, 泊松分布, 正态分布 都是它的情人

$\beta=1$ 泊松 $\text{Poisson}(X=k|\lambda) = \frac{\lambda^k e^{-\lambda}}{k!}$

离散

即 Gamma 分布中令 $\alpha = k+1$

$$\text{Gamma}(x|\alpha=k+1) = \frac{x^k e^{-x}}{\Gamma(k+1)} = \frac{x^k e^{-x}}{k!}$$

= 二项分布 $B(n, p)$ 在 $np = \lambda, n \rightarrow \infty$ 时的极限分布即是泊松分布

共轭分布. 先验分布和后验分布的形式一样.

二项分布的共轭分布就是 Beta 分布.

= 二项分布. $B(k|n, p) = \binom{n}{k} p^k (1-p)^{n-k}$.

Beta分布 $Beta(p|\alpha, \beta) = \frac{\Gamma(\alpha+\beta)}{\Gamma(\alpha)\Gamma(\beta)} p^{\alpha-1} (1-p)^{\beta-1}$

~~后验分布~~

~~$$\begin{aligned} P(p|n, k, \alpha, \beta) &\propto P(k|n, p) P(p|\alpha, \beta) \\ &= P(k|n, p) P(p|\alpha, \beta) \\ &= B(k|n, p) Beta(p|\alpha, \beta) \end{aligned}$$~~

LDA

文档 \rightarrow 主题 \rightarrow 主题词

文档 \rightarrow 主题 先验分布是 Dirichlet 分布 $\theta_d = D(\vec{\alpha})$ α 是 K 维向量

主题 \rightarrow 主题词 先验分布是 Dirichlet 分布 $\beta_k = D(\vec{\eta})$ η 是 V 维向量 V 是词表
中所有词个数

对于任一文档 d 中第 n 个词, 从主题分布 θ_d 中得到它主题编号 z_{dn} 的分布

$$z_{dn} = \text{multi}(\theta_d)$$

而对于该主题编号, 词 w_{dn} 的概率分布为

$$w_{dn} = \text{multi}(\beta_{z_{dn}})$$

文本建模

文档 对人来说 有序的词序列 $d = w_1, w_2, \dots, w_n$

统计文本建模 对机器 猜测上帝如何抛骰子的

什么样的骰子: 模型参数
每一面的概率
如何抛骰子: 可能有多不同的骰子, 按规则
则投掷

① 最简单无序模型

① 无序模型 \rightarrow 假设词之间, 文档之间独立可交换, 词序不重要

1) 一个骰子, V 面一面对应一个词, 概率不同记 $\vec{p} = (p_1, p_2, \dots, p_V)$

2) 一篇文档有 n 个词, 那么就是独立的抛了 n 次

像一个袋子
也叫词袋模型

这个分布服从多项分布: $w \sim \text{Mult}(w | \vec{p})$

对一篇文档 $d = \vec{w} = (w_1, w_2, \dots, w_n)$

$$p(\vec{w}) = p(w_1, w_2, \dots, w_n) = p(w_1) \cdot p(w_2) \cdot \dots \cdot p(w_n)$$

$$\text{对多篇文档 } p(W) = p(\vec{w}_1) p(\vec{w}_2) \cdot \dots \cdot p(\vec{w}_m)$$

假设总词频为 N , 每个词 v_i 发生次数 n_i

那么 $\vec{n} = (n_1, n_2, \dots, n_V)$ 正好是一个多项分布

$$p(\vec{n}) = \text{Mult}(\vec{n} | \vec{p}, N) = \left(\frac{N}{\vec{n}} \right) \prod_{k=1}^V p_k^{n_k}$$

$$\text{语料概率 } p(W) = p(\vec{w}_1) p(\vec{w}_2) \cdot \dots \cdot p(\vec{w}_m) = \prod_{k=1}^V p_k^{n_k}$$

用最大似然估计最大化 $p(W)$ $\hat{p}_i = \frac{n_i}{N}$

② 贝叶斯无参模型 ②

1) 无穷个骰子 每个骰子有 V 面

2) 上帝抽了一个骰子出来, 用这一个骰子不断抛, 产生了所有词。

有些类型的骰子多, 有些少, 骰子 \vec{p} 服从分布 $p(\vec{p})$

每个骰子都可能被使用, 由先验分布 $p(\vec{p})$ 决定。这个分布称为参数 \vec{p} 的先验分布。

对一个具体的骰子 \vec{p} , 产生数据概率是 $p(W|\vec{p})$ 。

则 最终 $p(W) = \int p(W|\vec{p}) p(\vec{p}) d\vec{p}$ 积分累加求和。

而注意到, $p(\vec{n}) = \text{Mult}(\vec{n}|\vec{p}, N)$ 实际上是计算一个多项分布的概率。

所以先验分布的一个比较好的选择就是多项分布的对应分布。

取 Dirichlet 分布, $\text{Dir}(\vec{p}|\vec{\alpha}) = \frac{1}{\Delta(\vec{\alpha})} \prod_{k=1}^V p_k^{\alpha_k - 1}$ $\vec{\alpha} = (\alpha_1, \dots, \alpha_V)$

$\Delta(\vec{\alpha})$ 就是归一化因子 $\text{Dir}(\vec{\alpha})$

$$\text{即 } \Delta(\vec{\alpha}) = \int \prod_{k=1}^V p_k^{\alpha_k - 1} d\vec{p}$$

Dirichlet 先验 + 多项分布数据 \rightarrow 后验分布也是 Dirichlet 分布

$$\text{Dir}(\vec{p}|\vec{\alpha}) + \text{Mult}(\vec{n}) = \text{Dir}(\vec{p}|\vec{\alpha} + \vec{n})$$

即在给定参数 \vec{p} 的先验分布为 $\text{Dir}(\vec{p}|\vec{\alpha})$ 时,

各个词出现的频次是 $\vec{n} \sim \text{Mult}(\vec{n}|\vec{p}, N)$ 为多项分布。

则后验分布是 $\vec{p}(\vec{p}, N, \vec{\alpha}) = \text{Dir}(\vec{p}|\vec{n} + \vec{\alpha})$

$$= \frac{1}{\Delta(\vec{n} + \vec{\alpha})} \prod_{k=1}^V p_k^{n_k + \alpha_k - 1} d\vec{p}$$

\vec{p} 如何估计呢? $E(\vec{p}) = (\frac{n_1 + \alpha_1}{\sum_{i=1}^V (n_i + \alpha_i)}, \frac{n_2 + \alpha_2}{\sum_{i=1}^V (n_i + \alpha_i)}, \dots, \frac{n_V + \alpha_V}{\sum_{i=1}^V (n_i + \alpha_i)})$

$$\hat{p}_i = \frac{n_i + \alpha_i}{\sum_{i=1}^V (n_i + \alpha_i)} \Leftrightarrow \frac{\text{先验估计数} + \text{数据中的计数}}{\text{总估计数}}$$

文本语料产生的概率

$$P(W|\vec{\alpha}) = \int P(W|\vec{p}) P(\vec{p}|\vec{\alpha}) d\vec{p}$$

$$= \int \prod_{k=1}^V P_k^{n_k} \text{Dir}(\vec{p}|\vec{\alpha}) d\vec{p}$$

$$= \int \prod_{k=1}^V P_k^{n_k} \frac{1}{\Delta(\vec{\alpha})} \prod_{k=1}^V P_k^{\alpha_k-1} d\vec{p}$$

$$= \frac{1}{\Delta(\vec{\alpha})} \int \prod_{k=1}^V P_k^{n_k + \alpha_k - 1} d\vec{p}$$

$$= \frac{\Delta(\vec{n} + \vec{\alpha})}{\Delta(\vec{\alpha})}$$

③
PLSA Topic 模型

Hoffman 1999年 PLSA (Probabilistic Latent Semantic Analysis)

1) 有两种骰子: doc-topic K 面 每面一个 topic $1 \sim K$
topic-word V 面 每面一个 word, 有 K 个, 编号 $1 \sim K$

2) 先有一个特定的 doc-topic 骰子.

投掷得到 topic 编号 z

投掷编号 z 的骰子得到单词.

即文档之间独立可交换; 一个文档内的词独立可交换.

有 K 个 topic-word 骰子. 记 $\vec{\psi}_1 \dots \vec{\psi}_K$

语料 C 有 M 篇文档. 记 d_1, d_2, \dots, d_M 都对应一个 doc-topic 骰子 $\vec{\theta}_1, \dots, \vec{\theta}_M$

则第 m 篇文档中 d_m 里每个词生成的概率.

$$P(\vec{w} | d_m) = \prod_{i=1}^n \sum_{z=1}^K P(w_i | z) P(z | d_m) = \prod_{i=1}^n \sum_{z=1}^K \psi_{z w_i} \theta_{z m}$$

求解可用 EM 算法. 参考 Hoffman 原文.

④ LDA

Latent Dirichlet Allocation

doc-topic 的骰子和 topic-word 的骰子

都是参数 要有先验分布!!

而 ϕ_k 和 θ_m 对应多项分布 先验是 Dirichlet 分布