

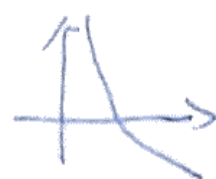
# 概率论：来自 Logistic & Softmax 回归 - 节

② 信息熵 (即概率④)  $H(P) = -\sum_{k=1}^K (P_k \ln P_k)$  是  $n$  个物品装到  $K$  个箱子有多少

种装法, 在  $n \rightarrow +\infty$  时的极限值。熵越高, 表示混乱程度越高,

分法越多, 分到  $K$  个箱子的概率越趋于相同。

③ 信息量:  $h(x) = -\log_2 P(x)$  小概率事件发生携带的信息量大



小概率  $\leftrightarrow$  大信息量

④ 独立事件  $X, Y$  满足  $h(XY) = h(X) + h(Y)$   
 $X, Y$  同时发生。

$$h(XY) = -\log_2 P(XY) = -\log_2 (P(X) \cdot P(Y)) = -\log_2 P(X) - \log_2 P(Y) = h(X) + h(Y)$$

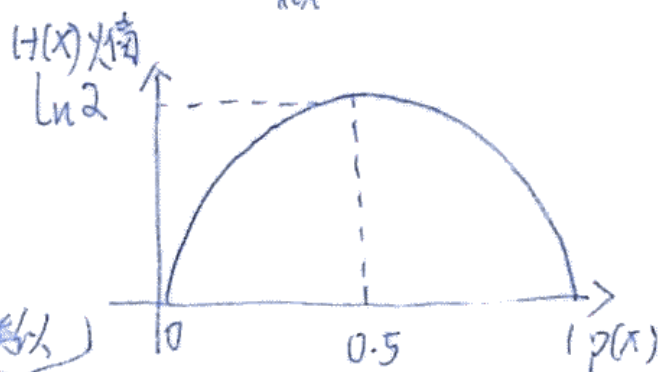
⑤ 用信息量也可以定义信息熵: 即随机事件信息量的数学期望。

$$H(X) = \sum_{x \in X} P(x) h(x) = \sum_{x \in X} P(x) (-\log_2 P(x)) = \boxed{-\sum_{x \in X} P(x) \ln(P(x))}$$

底是  $e$  还是  $2$  不要紧,  
 可以用换底公式替换。

## ③ 概率分布和熵的关系

① 两点分布的熵:  $H(X) = -\sum_{x \in X} P(x) \ln(P(x)) = -p \ln p - (1-p) \ln(1-p)$  根据熵的定义



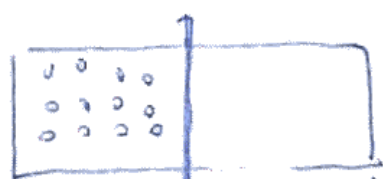
$\Rightarrow$  当  $P(x=0) = P(x=1) = 0.5$  时,  
 信息熵最高, 混乱程度最高。

② 三点分布的熵:  $H(X) = -\sum_{x \in X} P(x) \ln(P(x)) = -p_1 \ln p_1 - p_2 \ln p_2 - p_3 \ln p_3$  其中  $p_1 + p_2 + p_3 = 1$

③ 均匀分布的熵:  $N$  表示某离散分布可取  $N$  个值, 概率都是  $1/N$

离散:  $H(P) = -\sum_{i=1}^N P_i \ln P_i = -\sum_{i=1}^N \frac{1}{N} \ln \frac{1}{N} = \ln N \Rightarrow$  离散分布信息熵位于  $[0, \log N]$

## ③ 从基础问题理解最大熵公式



去掉隔板, 烟自然扩散



熵是随机变量不确定性的度量, 不确定性越大,  
 熵越大 { 随机变量退化成定值, 熵最小, 为 0.  
 { 随机变量为均匀分布, 熵最大。

以上为无条件: 的最大熵分布

若有条件呢? 给定期望和方差前提下, 最大熵的分布形式呢?

③③ 如果  $\ln(p(x))$  是  $\alpha x^2 + \beta x + \gamma$  形式, 则  $p(x)$  服从正态分布

证明: 利用  $\ln(p(x))$  计算过程的可逆性

$$p(x) = \frac{1}{\sqrt{2\pi}\sigma} e^{-\frac{(x-\mu)^2}{2\sigma^2}} \quad // \text{正态分布概率密度}$$

$$\ln p(x) = \ln \frac{1}{\sqrt{2\pi}\sigma} - \ln 6 - \frac{(x-\mu)^2}{2\sigma^2} = \alpha x^2 + \beta x + \gamma \quad // \text{对数正态分布}$$

该对数正态分布是关于随机变量  $x$  的二次函数, 且计算过程可逆

③④ 如果  $\ln(p(x))$  是  $A \cdot x + B \ln(x) + C$  的形式, 则  $p(x)$  服从 Gamma 分布

证明:

Gamma 函数:  $T(\alpha) = \int_0^{\infty} t^{\alpha-1} e^{-t} dt$

参考: Beta 分布, Gamma 函数  $\rightarrow$  概率 15

$T(n) = (n-1)!$  是阶乘在实数域的推广

Gamma 分布  $f(x; \alpha, \beta) = \frac{\beta^\alpha}{T(\alpha)} \cdot x^{\alpha-1} e^{-\beta \cdot x} \quad x \geq 0 \quad (\text{常数 } \alpha, \beta > 0)$

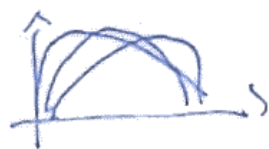
数学期望  $E(x) = \frac{\alpha}{\beta}$

对数 Gamma 分布:

$$\begin{aligned} \ln f(x; \alpha, \beta) &= \alpha \ln \beta + (\alpha-1) \ln x - \beta x - \ln T(\alpha) \\ &= A \cdot x + B \cdot \ln x + C \end{aligned}$$

该计算过程可逆

对比 Beta 分布



抛物线下面积

$$B(\alpha, \beta) = \int_0^1 x^{\alpha-1} (1-x)^{\beta-1} dx = \frac{T(\alpha)T(\beta)}{T(\alpha+\beta)}$$

数学期望是  $\frac{\alpha}{\alpha+\beta}$



⑤ 给定方差的最大熵分布：带约束条件的极值问题，用拉格朗日乘子法

目标函数：  $\arg \max_{p(x)} H(x) = -\sum_x p(x) \ln p(x)$  s.t.  $\begin{cases} E(x) = \mu \\ \text{Var}(x) = \sigma^2 \Leftrightarrow E(x^2) = (E(x))^2 + \text{Var}(x) = \mu^2 + \sigma^2 \end{cases}$

根据约束条件 ⑤, ⑥ 用拉格朗日乘子法建模

$$L(p) = -\sum_x p(x) \cdot \ln p(x) + \lambda_1 (E(x) - \mu) + \lambda_2 (E(x^2) - \mu^2 - \sigma^2)$$

根据方差公式  $\text{Var}(x) = E(x^2) - (E(x))^2$

$$= -\sum_x p(x) \ln p(x) + \lambda_1 \left( \sum_x x \cdot p(x) - \mu \right) + \lambda_2 \left( \sum_x x^2 p(x) - \mu^2 - \sigma^2 \right)$$

要求  $p(x)$  即  $p$  为何值时， $H(x) = -\sum_x p(x) \ln p(x)$  即  $-\sum_x p \ln p$  值最大

对  $p$  求偏导，另偏导为 0，得到驻点

$$\frac{\partial L}{\partial p} = \frac{\partial \left( -\sum_x p \ln p + \lambda_1 \left( \sum_x x \cdot p - \mu \right) + \lambda_2 \left( \sum_x x^2 \cdot p - \mu^2 - \sigma^2 \right) \right)}{\partial p}$$

$$= -\ln p - 1 + \lambda_1 x + \lambda_2 x^2$$

令偏导为 0，得到

$$\ln p(x) = \ln p = \lambda_2 x^2 + \lambda_1 x + 1$$

是关于随机变量  $x$  的二次形，根据引理（概率③③），

$p$  服从正态分布

⇒ 最大熵模型是建立在  $p$  服从正态分布的假设之上

③⑥ 联合熵  $H(X, Y) = -\sum_{x,y} p(x,y) \ln p(x,y)$

随机变量  $X, Y$  联合分布的信息熵

对信息熵定义 (概率④②⑨③⑩)

$$H(Z) = -\sum_z p(z) \ln p(z)$$

其中  $-\ln p(z)$  是信息量

③⑦ 条件熵  $H(Y|X) = H(X, Y) - H(X)$  ①

\* 表示  $X$  发生的前提下,  $X, Y$  同时发生“新”带来的熵。  
(混乱度)

$$H(Y|X) = -\sum_{x,y} p(x,y) \ln p(y|x) \quad \text{⑥}$$

$$= -\sum_{x,y} [p(y|x) \cdot p(x)] \ln p(y|x)$$

\* 也可以理解为条件概率信息量的数学期望 (概率密度)  
\* 公式① 等价于公式⑥

$$H(Y|X) = H(X, Y) - H(X) \quad \text{|| 条件熵定义①}$$

$$= -\sum_{x,y} p(x,y) \ln p(x,y) + \sum_x p(x) \ln p(x) \quad \text{|| 信息熵, 联合熵定义①}$$

$$* H(Y|X) = H(X, Y) - H(X)$$

$$= H(Y) - I(X, Y)$$

互信息: 见

概率③⑨④⑩

$$= -\sum_{x,y} p(x,y) \ln p(x,y) + \sum_x \left( \sum_y p(x,y) \right) \ln p(x)$$

$$= -\sum_{x,y} p(x,y) \ln p(x,y) + \sum_{x,y} p(x,y) \ln p(x)$$

$$= -\sum_{x,y} p(x,y) \ln \frac{p(x,y)}{p(x)}$$

$$= -\sum_{x,y} p(x,y) \ln p(y|x) \quad \text{|| 公式⑥}$$

设用  $p(y|x)$

因为要表示信息增量

$$= -\sum_x p(x) \cdot p(y|x) \ln p(y|x)$$

$$= \sum_x p(x) H(Y|X=x) \quad \text{③}$$

\* 等价于公式③: 给定  $X=x$  时,  $Y$  的熵的数学期望

$p(x, y=0)$

$p(x, y=1)$



$p(y=0)$   $p(y=1)$

$p(x,y)$  的分母是整个概率

$p(x|y)$  的分母是  $p(y)$ , 不

能表示增量



③ 相对熵 (也叫互熵, 交叉熵, 鉴别信息, Kullback 熵, Kullback-Leibler 散度, 简称 K-L 距离), 可以理解为对  $\log \frac{p(x)}{q(x)}$  求数学期望。

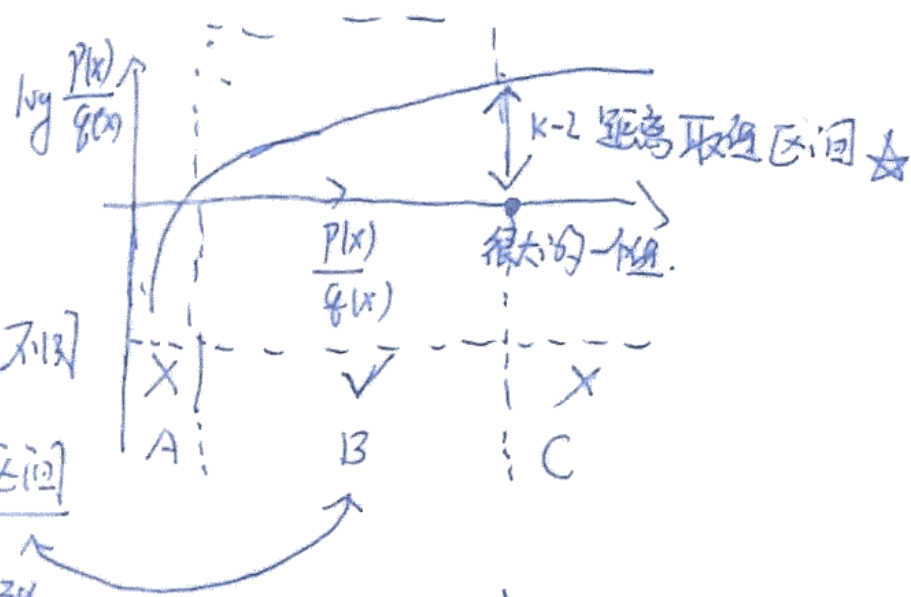
$$D(p||q) = \sum_x p(x) \log \frac{p(x)}{q(x)}$$

难点:  $D(p||q)$  与  $D(q||p)$  该用哪个?

因为 K-L 距离是非对称的, 上面 2 个值不同

选择依据: 尽量让  $\frac{p(x)}{q(x)}$  或  $\frac{q(x)}{p(x)}$  落在 B 区间

这样 K-L 距离取值区间不会太夸张



两个 K-L 散度的区别: PPT

概率 ③⑧

③⑨ 互信息: 随机变量  $X, Y$  的 ④联合分布 与 ⑥独立分布 的乘积 的 相对熵

$$I(X, Y) = D(p(x, y) || p(x)p(y))$$

表示没有  $X$  时  $Y$  的不确定性

$$= \sum_{x,y} p(x, y) \log \frac{p(x, y)}{p(x)p(y)}$$

给定  $X$  时  $Y$  的不确定性  
的差值

④ 条件熵与互信息关系

$$H(Y|X) = H(X, Y) - H(X)$$

定义式:  $X$  发生前提下,  $X, Y$  同时发生新带来的混乱度

$$H(Y|X) = H(Y) - I(X, Y)$$

不依赖  $X$  时  $Y$  的不确定性 差值  
表示 给定  $X$  时,  $Y$  的不确定性。即。

推导:  $H(Y) - I(X, Y)$

$$= -\sum_y p(y) \log p(y) - \sum_{x,y} p(x, y) \log \frac{p(x, y)}{p(x)p(y)}$$

$$= -\sum_{x,y} p(x, y) \log(p(y)) - \sum_{x,y} p(x, y) \log \frac{p(x, y)}{p(x)p(y)}$$

$$= -\sum_{x,y} p(x, y) \log \frac{p(x, y)}{p(x)}$$

$$= -\sum_{x,y} p(x, y) \log p(y|x)$$

$$= H(Y|X)$$

#### ④ 互信息与信息熵，联合熵之间的关系

$$I(X, Y) = H(X) + H(Y) - H(X, Y)$$

证明:  $H(X) + H(Y) - H(X, Y)$

$$= -\sum_x p(x) \log p(x) - \sum_y p(y) \log p(y) + \sum_{x,y} p(x,y) \log p(x,y)$$

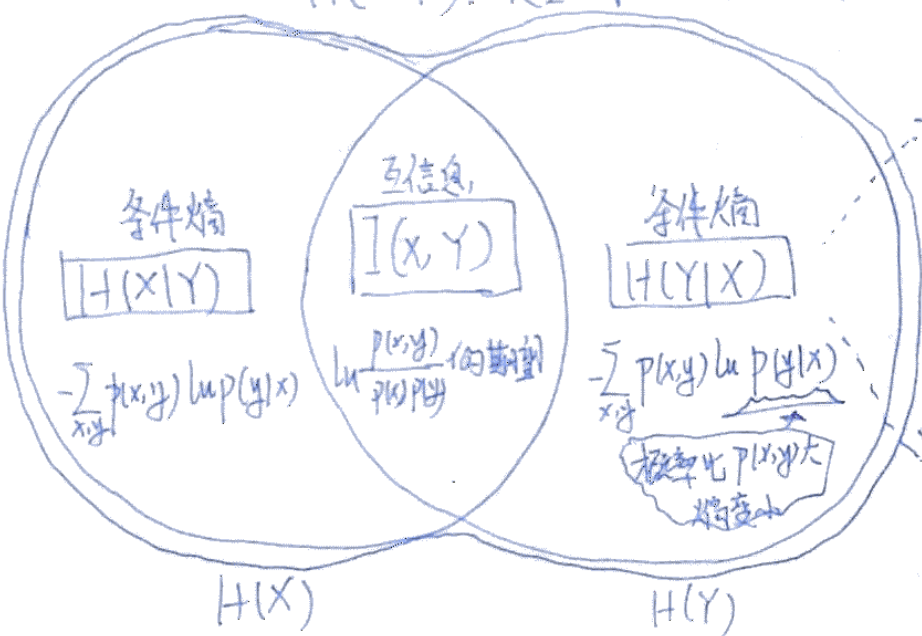
$$= -\sum_{x,y} p(x,y) \log p(x) - \sum_{x,y} p(x,y) \log p(y) + \sum_{x,y} p(x,y) \log p(x,y)$$

$$= -\sum_{x,y} p(x,y) \log \frac{p(x,y)}{p(x)p(y)} = I(X, Y)$$

#### ④ Venn图: 信息熵，联合熵，条件熵，互信息之间的关系。

$-\sum_{x,y} p(x,y) \ln p(x,y)$  概率越小，混乱程度越高，熵越大

$H(X, Y)$ : 联合熵,  $X, Y$  联合分布的信息熵



$X$  的信息熵

$$-\sum_x p(x) \ln p(x)$$

$Y$  的信息熵

$$-\sum_y p(y) \ln p(y)$$

$$H(Y|X) = H(X, Y) - H(X)$$

$\downarrow$   
X发生前提下  
X, Y同时发生(概率小, 熵大)  
新带来的不确定性

$$H(Y|X) = H(Y) - I(X, Y)$$

$$I(X, Y) = H(Y) - H(Y|X)$$

$\uparrow$   
不考虑X取值,  
Y自身的不确定性  
 $\downarrow$   
X发生前提下  
Y的不确定性  
指  $p(y|x)$   
(不是  $p(x, y)$ )  
概率大熵小