

# 数基基础 8<sup>th</sup>

135  
246

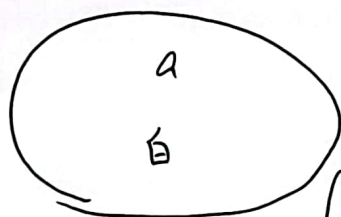
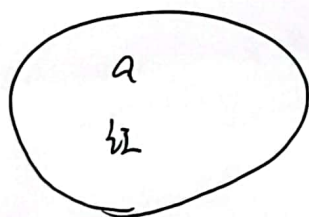
$$1. \quad P(A) = \frac{1}{6} \times \frac{5}{6} \times 2 = \frac{5}{36} \times 2 = \frac{5}{18}$$

$$P(B) = \frac{1}{6} \times \frac{5}{6} \times 2 + \frac{1}{36} = \frac{5}{18} + \frac{1}{36} = \frac{11}{36}$$

$$P(C) = (\text{奇} + \text{奇}), \text{偶} + \text{偶} = \frac{3 \times 3 + 3 \times 3}{6 \times 6} = \frac{18}{36} = \frac{1}{2}$$

$$I_A = -\log\left(\frac{5}{18}\right) \quad I_B = -\log\frac{11}{36} \quad I_C = -\log\frac{1}{2}$$

2.



从 2a 个球中取球

概率空间

$$\begin{pmatrix} X \\ P(X) \end{pmatrix} = \begin{pmatrix} X=x_1 \dots X=x_i \dots X=x_k \\ P(x_1) \dots P(x_i) \dots P(x_k) \end{pmatrix}$$

$x_i = 1/0 \rightarrow$  第  $i$  次取出白球  
↓  
第  $i$  次取出红球

① 有放回情况下, 每个  $P(x_i)$  概率都相同, 因为每次放回都是同样的一堆球, 且  $\sum P_i = 1$

② 无放回情况下, 每次的情况都与前几次抽到什么球有关, 所以  $P_i$  都各不相同, 且  $\sum P_i = 1$

根据熵的极值性: 离散信息源中各消息等概率出现时熵最大

$\therefore$  有放回的熵更大



$$3. D_{KL}(P||Q) = E_{x \sim P} [\log P(x) - \log Q(x)] = E_{x \sim P} \left[ \log \frac{P(x)}{Q(x)} \right]$$

$E_{x \sim P} [f(x)]$   $f(x)$  关于  $P(x)$  的期望

$$\therefore D_{KL}(P||Q) = \left[ \log \frac{P(x)}{Q(x)} \text{ 关于 } P(x) \text{ 的期望} \right]$$

$$H(P, Q) = -E_{x \sim P} \log Q(x)$$

样本  $(x_1, x_2, \dots, x_i, \dots, x_n)$

假设  $x_i$  属于标签  $y$  (真实)

类  $(y_1, y_2, \dots, y_i, \dots, y_k)$

$\Downarrow$   
转换成 one-hot

$$p_i = (0, 0, \dots, \underset{\substack{\downarrow \\ \text{第 } i \text{ 个}}}{1}, \dots, 0)^T$$

对于  $x_i$ , 给定  $p_i$ ,  $p_i$  为第  $i$  个分量为 1 的 one-hot 向量

$q_i$  为预测的标签,  $q_i = t(x_i, \theta)$

$$\therefore D_{KL}(P||Q) = (p_i)^T \left( \log \frac{p_i}{q_i} \right) = \underline{p_i^T \log p_i} - p_i^T \log q_i$$

$$H(p_i, q_i) = -p_i^T \log q_i \quad \text{真实不变的常量}$$

$\therefore D_{KL}(P||Q)$  与  $H(p_i, q_i)$  是等价的

针对  $q_i$  最小化交叉熵等价于最小化 KL 散度, 因为  $q_i$  不参与被省略的一项

$$\therefore \argmin_{\theta} D_{KL}(p_i||q_i) = \argmin_{\theta} H(p_i, q_i)$$

即  $\theta = \hat{\theta}$  时, 可以同时使 KL 散度和交叉熵函数最小



习题 4

$$P(X_1, X_2, \dots, X_n) = P(X_1)P(X_2|X_1) \dots P(X_n|X_{n-1})$$

$$I(X_1; X_2, \dots, X_n) = H(X_1) - H(X_1|X_2, \dots, X_n)$$

$$= H(X_1) - (H(X_1, X_2, \dots, X_n) - H(X_2, \dots, X_n))$$

$$= H(X_1) - (H(X_1) + H(X_2|X_1) + H(X_3|X_1X_2) + \dots + H(X_n|X_1 \dots X_{n-1})) + (H(X_2) + H(X_3|X_2) + H(X_4|X_2X_3) + \dots + H(X_n|X_2X_3 \dots X_{n-1}))$$

$$= H(X_1) - (H(X_1) + \sum_{i=2}^n H(X_i|X_{i-1})) + H(X_2) + \sum_{i=3}^n H(X_i|X_{i-1})$$

$$= -\sum_{i=2}^n H(X_i|X_{i-1}) + H(X_2) + \sum_{i=3}^n H(X_i|X_{i-1})$$

$$= H(X_2) - H(X_2|X_1)$$

$$= I(X_2; X_1) = I(X_1; X_2)$$

① for 马尔科夫

$$P(X_1, X_2) = P(X_1) \cdot P(X_2|X_1) \checkmark$$

$$P(X_1, X_2, X_3) = \underline{P(X_1)} \cdot \underline{P(X_2|X_1)} \cdot P(X_3|X_1X_2) = \underline{P(X_1)} \underline{P(X_2|X_1)} P(X_3|X_2)$$

$$P(X_3|X_1X_2) = P(X_3|X_2)$$

$$P(X_1, X_2, X_3, X_4) = \cancel{P(X_1)} \cdot \cancel{P(X_2|X_1)} \cdot \cancel{P(X_3|X_1X_2)} \cdot \underline{P(X_4|X_1X_2X_3)} = \cancel{P(X_1)} \cancel{P(X_2|X_1)} \cancel{P(X_3|X_2)} \cdot \underline{P(X_4|X_3)}$$

$$P(X_4|X_1X_2X_3) = P(X_4|X_3)$$

$$\therefore \text{for } \forall \quad P(X_i | \prod_{k=1}^{i-1} X_k) = P(X_i | X_{i-1}) \quad (\text{在马尔科夫链中})$$

$$H(X_3|X_1X_2) = E(\log \frac{1}{P(X_3|X_1X_2)}) = E(\log \frac{1}{P(X_3|X_2)}) = H(X_3|X_2)$$

$$H(X_4|X_1X_2X_3) = E(\log \frac{1}{P(X_4|X_1X_2X_3)}) = E(\log \frac{1}{P(X_4|X_3)}) = H(X_4|X_3)$$

$$H(X_i | \prod_{k=1}^{i-1} X_k) = H(X_i | X_{i-1}) \quad (\text{在马尔科夫链中})$$

