

09/07/2022 Wed

Entropy of coin shows heads

$$H(x) = - \sum_{k=1}^k p \log_2 p$$

- $P(x) = 0.5$ $\log_2 \frac{1}{2}$

$$\begin{aligned} H(x) &= - 0.5 \log_2 0.5 - 0.5 \log_2 0.5 \\ &= - 0.5 \cdot (-1) - 0.5 \cdot (-1) \\ &= 1 \end{aligned}$$

- $P(x) = 1$ getting heads

$$\begin{aligned} H(x) &= - 1 \log_2 1 - 0 \log_2 0 \\ &= 0 \end{aligned}$$

KL divergence

$$KL(p \parallel q) \triangleq \sum_{k=1}^k p \log \frac{p}{q}$$

$$\log \frac{p}{q} = \log p - \log q$$

$$= \sum_k p \log p - \sum_k p \log q$$

$$= -H(p) + H(p, q)$$

entropy

Cross-entropy

$$KL(p \parallel q) \neq KL(q \parallel p)$$

$$D = \{1, 0, 1, 1\}$$

$$X \sim \text{Ber}(\theta), \quad 0 \leq \theta \leq 1$$

$$\text{Likelihood of } D \quad P(D|\theta) = \prod_{n=1}^N P(x_n|\theta)$$

$$P(D|\theta) = \theta \cdot (1-\theta) \cdot \theta \cdot \theta$$

$$= \theta^3 \cdot (1-\theta)$$

$$D = \{N_1 \text{ positive samples} ; N_2 \text{ negative samples}\}$$

$$P(D|\theta) = \theta^{N_1} (1-\theta)^{N_2}$$

MLE : maximum Likelihood Estimate

$$\hat{\theta}_{MLE} = \arg \max_{\theta} P(D|\theta)$$

$$= \arg \max_{\theta} \theta^{N_1} (1-\theta)^{N_2}$$

$$\frac{\partial P(D|\theta)}{\partial \theta} = 0$$

$$\frac{\partial \theta^{N_1} (1-\theta)^{N_2}}{\partial \theta}$$

$$= N_1 \theta^{N_1-1} (1-\theta)^{N_2} - N_2 \theta^{N_1} (1-\theta)^{N_2-1} = 0$$

$$N_1 (1-\theta) = N_2 \theta$$

$$N_1 - N_1 \theta = N_2 \theta$$

$$\hat{\theta}_{MLE} = \frac{N_1}{N_1 + N_2}$$

Posterior \propto likelihood \times prior

$$P(\theta|D) \propto P(D|\theta) \cdot P(\theta)$$

$$\propto \theta^{N_1} (1-\theta)^{N_2} \theta^{a-1} (1-\theta)^{b-1}$$

$$\propto \theta^{N_1+a-1} (1-\theta)^{N_2+b-1}$$

$$\text{Beta}(\theta|a,b) \propto \theta^{a-1} (1-\theta)^{b-1}$$

MAP : maximum a posteriori estimate

$$\hat{\theta}_{\text{MAP}} = \underset{\theta}{\text{argmax}} P(\theta|D)$$

$$= \frac{N_1 + a - 1}{N_1 + N_2 + a + b - 2}$$

$$E(\theta|D) = \frac{N_1 + a}{N_1 + N_2 + a + b}$$

Bayesian Averaging

$$P(x^* = 1|D) = \int_0^1 P(x^* = 1|\theta) P(\theta|D) d\theta$$

$$= \int_0^1 \theta \cdot \text{Beta}(\theta|N_1+a, N_2+b) d\theta$$

$$= \frac{N_1 + a}{N_1 + N_2 + a + b}$$

$$KL(P \parallel Q) = \sum_y P(y) \log \frac{P(y)}{Q(y)} = \underbrace{\sum_y P \log P}_{\text{entropy } P} - \underbrace{\sum_y P \log Q}_{\text{cross-entropy}}$$

$$P(y) = P_D(y) \quad \text{empirical distribution} \triangleq \frac{1}{N} \sum_{n=1}^N \delta(y - y_n)$$

$$Q(y) = P(y|\theta) \quad \text{learned model}$$

$$\underline{KL(P \parallel Q)} = \sum_y P_D(y) \log P_D(y) - \sum_y \underline{P_D(y)} \log P(y|\theta)$$

$$= -H(P_D) - \frac{1}{N} \sum_{n=1}^N \log(y_n|\theta)$$

$$= \text{const} + NLL(\theta)$$

$$\underline{NLL(\theta)} = -\log P(D|\theta)$$

$$= -\log \prod_y P(y|\theta)$$

$$= -\sum_y \log(y|\theta)$$

$$\hat{h}_{MAP} = \underset{h}{\operatorname{argmax}} \underbrace{P(D|h)}_{\text{likelihood}} \underbrace{p(h)}_{\text{prior}}$$

$$= \underset{h}{\operatorname{argmax}} (\log P(D|h) + \log p(h))$$