

Information Theory

Machine learning, 2021

Mansoor Rezghi

Department of Computer science, TMU

Shanon Entropy

Information $\propto \frac{1}{Possibility}$

Information $\propto \frac{1}{P(x)}$

Adaditivity: For independent events Z_1 and Z_2

Inf($Z_1 + Z_2$) \propto **Inf**(Z_1) + **Inf**(Z_2)

$I(X) = -\log P(x)$

Covers two mwntioed properties

Shannon entropy $H(x) = \mathbb{E}_{x \sim P}[I(x)] = -\mathbb{E}_{x \sim P}[\log P(x)]$.

Entropy

The **entropy** of a random variable X with distribution p , denoted by $\mathbb{H}(X)$ or sometimes $\mathbb{H}(p)$, is a measure of its uncertainty. In particular, for a discrete variable with K states, it is defined by

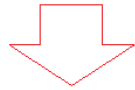
$$\mathbb{H}(X) \triangleq - \sum_{k=1}^K p(X = k) \log_2 p(X = k)$$

In order to communicate the value of x to a receiver, we would need to transmit a message of length 3 bits.

$$\mathbb{H}[x] = -8 \times \frac{1}{8} \log_2 \frac{1}{8} = 3 \text{ bits.}$$

Event $\{a, b, c, d, e, f, g, h\}$

Probability $\frac{1}{2}, \frac{1}{4}, \frac{1}{8}, \frac{1}{16}, \frac{1}{64}, \frac{1}{64}, \frac{1}{64}, \frac{1}{64}$



Coding

0, 10, 110, 1110, 111100, 111101, 111110, 111111



$$\mathbb{H}[x] = -\frac{1}{2} \log_2 \frac{1}{2} - \frac{1}{4} \log_2 \frac{1}{4} - \frac{1}{8} \log_2 \frac{1}{8} - \frac{1}{16} \log_2 \frac{1}{16} - \frac{4}{64} \log_2 \frac{1}{64} = 2 \text{ bits.}$$



$$\text{average code length} = \frac{1}{2} \times 1 + \frac{1}{4} \times 2 + \frac{1}{8} \times 3 + \frac{1}{16} \times 4 + 4 \times \frac{1}{64} \times 6 = 2 \text{ bits}$$

Kullback-Leibler (KL)divergence

Coding by inexact dist



$$\begin{aligned}\text{KL}(p\|q) &= - \int p(\mathbf{x}) \ln q(\mathbf{x}) \, d\mathbf{x} - \left(- \int p(\mathbf{x}) \ln p(\mathbf{x}) \, d\mathbf{x} \right) \\ &= - \int p(\mathbf{x}) \ln \left\{ \frac{q(\mathbf{x})}{p(\mathbf{x})} \right\} \, d\mathbf{x}.\end{aligned}$$



Coding by exact dist



Extra information for decoding

$$\text{KL}(p\|q) \neq \text{KL}(q\|p).$$

$$\text{KL}(p\|q) \geq 0$$

$p(\mathbf{x})$: True unknown distribution

$q(\mathbf{x})$: An approximation of $p(\mathbf{x})$

Kullback-Leibler (KL) divergence: The similarity of two distributions

$$D_{\text{KL}}(P\|Q) = \mathbb{E}_{x \sim P} \left[\log \frac{P(x)}{Q(x)} \right] = \mathbb{E}_{x \sim P} [\log P(x) - \log Q(x)]$$

$$D_{\text{KL}}(P\|Q) \neq D_{\text{KL}}(Q\|P)$$

cross-entropy

$$H(P, Q) = H(P) + D_{\text{KL}}(P\|Q) = -\mathbb{E}_{x \sim P} \log Q(x)$$

$$\operatorname{argmin}_Q H(P, Q) = \operatorname{argmin}_Q D_{\text{KL}}(P\|Q)$$

ML and Cross Entropy Equivalence

$$\theta_{\text{ML}} = \arg \max_{\theta} \sum_{i=1}^m \log p_{\text{model}}(\mathbf{x}^{(i)}; \theta). \quad \Rightarrow \quad \theta_{\text{ML}} = \arg \max_{\theta} \mathbb{E}_{\mathbf{x} \sim \hat{p}} \log p_{\text{model}}(\mathbf{x}; \theta) = \operatorname{argmin}_{\theta} - \mathbb{E}_{\mathbf{x} \sim \hat{p}} [\log p_{\text{model}}(\mathbf{x})]$$

$= \operatorname{argmin}_{\theta} H(\hat{p}, p_{\theta})$

$$D_{\text{KL}}(\hat{p}_{\text{data}} \| p_{\text{model}}) = \mathbb{E}_{\mathbf{x} \sim \hat{p}} [\log \hat{p}_{\text{data}}(\mathbf{x}) - \log p_{\text{model}}(\mathbf{x})]$$

One way to interpret maximum likelihood estimation is to view it as minimizing the dissimilarity between the empirical distribution \hat{p}_{data} defined by the training set and the model distribution, with the degree of dissimilarity between the two measured by the KL divergence.

Mutual Information

$$\begin{aligned} I[\mathbf{x}, \mathbf{y}] &\equiv \text{KL}(p(\mathbf{x}, \mathbf{y}) \| p(\mathbf{x})p(\mathbf{y})) \\ &= - \iint p(\mathbf{x}, \mathbf{y}) \ln \left(\frac{p(\mathbf{x})p(\mathbf{y})}{p(\mathbf{x}, \mathbf{y})} \right) d\mathbf{x} d\mathbf{y} \end{aligned}$$

$$I[\mathbf{x}, \mathbf{y}] = H[\mathbf{x}] - H[\mathbf{x}|\mathbf{y}] = H[\mathbf{y}] - H[\mathbf{y}|\mathbf{x}]$$

Conditional log-Likelihood

If \mathbf{X} represents all our inputs and \mathbf{Y} all our observed targets, then the conditional maximum likelihood estimator is

$$\boldsymbol{\theta}_{\text{ML}} = \arg \max_{\boldsymbol{\theta}} P(\mathbf{Y} \mid \mathbf{X}; \boldsymbol{\theta}).$$



If the examples are assumed to be i.i.d., then this can be decomposed into

$$\boldsymbol{\theta}_{\text{ML}} = \arg \max_{\boldsymbol{\theta}} \sum_{i=1}^m \log P(\mathbf{y}^{(i)} \mid \mathbf{x}^{(i)}; \boldsymbol{\theta}).$$