# Information Theory

**Shanon Entropy**

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

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

**Adaditivity**: For independent events $Z_1$ and $Z_2$

$$\mathbf{Inf}(Z_1 + Z_2) \propto \mathbf{Inf}(Z_1) + \mathbf{Inf}(Z_2)$$

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

Covers two mwntioed properties

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

$$\text{inf}(x) \alpha \frac{1}{P(n)} \quad, \quad \text{inf}(x,y) = \text{inf}(n) + \text{inf}(y) \leftarrow \text{inf}(n) = \log\left(\frac{1}{P(x)}\right) = -\log P(x)$$

$$x,y \text{ indpt} \qquad \qquad x \to P(x) \quad \boxed{-\log P(x)}$$

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

In order to co$\quad$ $\mathbb{H}(X) \triangleq -\sum_{k=1}^{K} p(X=k)\log_2 p(X=k)$ $\quad$ ssage of length 3 bits.

$$0\log_0 = 0 \qquad\qquad E(f(x)) = \sum f(x) P(x)$$

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

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

$$H(x) = E[-\log x] = -\sum P(x)\log P(x) \quad \text{Entropy}$$

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

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

| 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**

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

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

$P(x)$
$q(x)$

$d_{ij}(P(x), q(x))$

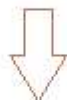$-\int P(x)\ln P(x)\,dx$

# Kullback-Leibler (KL)divergence

$\varepsilon, 7$

$-\int P(x)\ln q(x)\,dx$

$q(x)$

Coding by inexact dist

$$KL(p\|q) = -\int p(\mathbf{x}) \ln q(\mathbf{x})\,d\mathbf{x} - \left(-\int p(\mathbf{x}) \ln p(\mathbf{x})\,d\mathbf{x}\right)$$

$$= -\int p(\mathbf{x}) \ln \left\{\frac{q(\mathbf{x})}{p(\mathbf{x})}\right\}\,d\mathbf{x}.$$

Coding by exact dist

Extra information for decoding

$$KL(p\|q) \neq KL(q\|p).$$

$$KL(p\|q) \geqslant 0$$

p(x):     True unknown distribution

q(x):     An approximation of p(x)

$\boxed{q(x)} \rightarrow \ln q(x)$

$P(x)$

$P$  $x_1$  $x_2$  $x_3$  $x_4$

$\frac{1}{4}$  $\frac{1}{4}$  $\frac{1}{4}$  $\frac{1}{4}$

$Q$  $x_1$  $y_2$  $y_3$  $x_4$

$\frac{1}{4}$  $\frac{3}{8}$  $\frac{1}{8}$  $\frac{1}{4}$

$x_1$  $x_2$  $x_3$  $x_4$

$\frac{1}{2}$  $\frac{1}{2}$  $0$  $\frac{1}{4}$

$$E(f(x)) = \sum f(y) P(x)$$

· $\log P(x)$

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

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

$-H(P)$

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

**cross-entropy**

sym.  $Q$

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

$$f(x)$$

$$E_{x\sim P}\left(-\log P(x)\right)$$

$$argmin_Q H(P,Q) = argmin D_{KL}(P\|Q)$$

$$-\int P(x)\log P(x) \quad \left(-\int P(x)\log q(x)\right)$$

# ML and Cross Entropy Equivalence



$$\boldsymbol{\theta}_{\mathrm{ML}} = \arg\max_{\boldsymbol{\theta}} \sum_{i=1}^{m} \log p_{\mathrm{model}}(\boldsymbol{x}^{(i)}; \boldsymbol{\theta}).$$

$$\boldsymbol{\theta}_{\mathrm{ML}} = \arg\max_{\boldsymbol{\theta}} \mathbb{E}_{x \sim \hat{p}} \log p_{\mathrm{model}}(\boldsymbol{x}; \boldsymbol{\theta}) \quad = \quad argmin \ -\mathbb{E}_{\sim \hat{p}} [\log p_{\mathrm{model}}(\boldsymbol{x})]$$

$$= \quad argmin \ H(P, Q)$$

$$D_{\mathrm{KL}}(\hat{p}_{\mathrm{data}} \| p_{\mathrm{model}}) = \mathbb{E}_{\sim \hat{p}} [\log \hat{p}_{\mathrm{data}}(\boldsymbol{x}) - \log p_{\mathrm{model}}(\boldsymbol{x})]$$
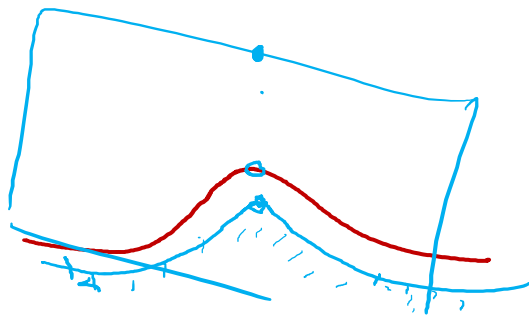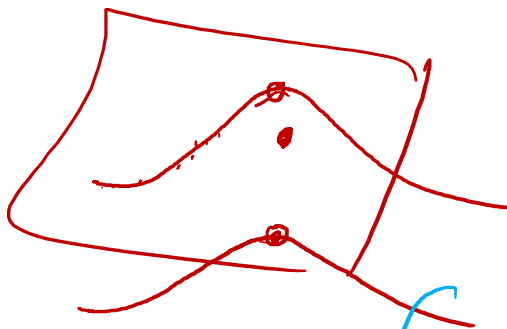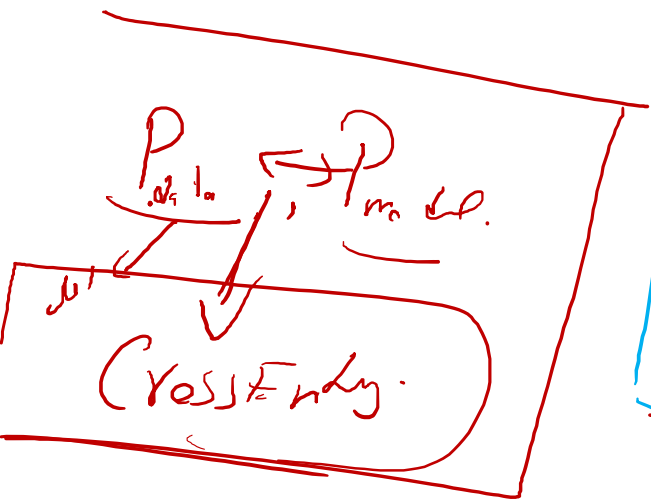
One way to interpret maximum likelihood estimation is to view it as minimizing the dissimilarity between the empirical distribution $\hat{p}_{\mathrm{data}}$ defined by the training set and the model distribution, with the degree of dissimilarity between the two measured by the KL divergence.

$P , q \qquad P , Q \qquad$ <span style="color:blue">$P, Q$</span>

$P_{N_q l_q} \longrightarrow P_{m} \, \text{d.l.}$

$M$

CrossEntropy.

<span style="color:blue">$-\int P(x) \log P(x)$</span>

<span style="color:blue">$\int P(x) \log Q(x)$</span>

'61

# Mutual Information

$$KL\left(p(x,y) \,\|\, p(x)p(y)\right)$$

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

$$I[x,y] \equiv KL(p(x,y)\|p(x)p(y))$$

$$= -\iint p(x,y)\ln\left(\frac{p(x)p(y)}{p(x,y)}\right)dx\,dy$$

$$-\iint p(x,y)\ln p(x,y) + \iint \boxed{p(x)p(y)} \ln p(x)p(y)$$

$$\frac{p(x,y)}{}$$

$$= \iint p(x,y)\ln\frac{p(x)p(y)}{p(x,y)}\,dx\,dy$$

$$I[x,y] = H[x] - H[x|y] = H[y] - H[y|x]$$

# Conditional log-Likelihood

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

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

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

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