

## Jensen's Inequality

$$f \text{ convex} \Rightarrow f(\mathbb{E}x) \leq \mathbb{E}(f(x))$$

In other words, if  $\lambda_1, \dots, \lambda_n$  are such that  $\lambda_i \geq 0$  for  $i=1, \dots, n$  &  $\sum_{i=1}^n \lambda_i = 1$ ,

then

$$f\left(\sum_{i=1}^n \lambda_i x_i\right) \leq \sum_{i=1}^n \lambda_i \cdot f(x_i)$$

Application 1:  $H(x) \leq \log |\Omega|$  for any finite sample space  $\Omega$ , and  $x: \Omega \rightarrow \mathbb{R}$ .

## Application 2

### Theorem (Shannon Inequality)

For any finite sample spaces  $\Omega_1, \Omega_2$ , and  $X: \Omega_1 \rightarrow \mathbb{R}$  and  $Y: \Omega_2 \rightarrow \mathbb{R}$ ,

$$H(X) \geq H(X|Y)$$

(i.e. conditioning never increases entropy.)

Moreover,  $H(X) = H(X|Y)$  iff  $X$  and  $Y$  are independent [HW2]

### Proof

Let

$$f(a) = a \cdot \ln a.$$

Then

$$\begin{aligned} \frac{df}{da} &= a + \frac{1}{a} + \ln a \\ &= 1 + \ln(a) \end{aligned}$$

$$\frac{d^2f}{da^2} = \frac{1}{a} > 0, a > 0.$$

$\Rightarrow f$  is convex.

Now,  $\sum_{i=1}^n \lambda_i f(x_i) \geq f\left(\sum_{i=1}^n \lambda_i x_i\right)$ . for convex weights  $\lambda_i, i=1, \dots, n$ .

Let  $\lambda_i = p(x_i)$ . Let  $x_i = p(y_j | x_i)$  (f free variable).

Then  $\sum_{i=1}^n p(x_i) f(p(y_j | x_i)) \geq f\left(\sum_{i=1}^n p(x_i) \cdot p(y_j | x_i)\right)$ .

i.e.  $\sum_{i=1}^n p(x_i) \cdot p(y_j | x_i) \cdot \log(p(y_j | x_i)) \geq \left(\sum_{i=1}^n p(x_i) p(y_j | x_i)\right) \log\left(\sum_{i=1}^n p(x_i) \cdot p(y_j | x_i)\right)$

i.e.  $\sum_{i=1}^n p(x_i, y_j) \log p(y_j | x_i) \geq \left(\sum_{i=1}^n p(x_i, y_j)\right) \log p(y_j)$   
 $= p(y_j) \log p(y_j)$

Summing over  $j$ ,  
 $\sum_{i=1}^n \sum_{j=1}^m p(x_i, y_j) \log p(y_j | x_i) \geq \sum_{j=1}^m p(y_j) \log p(y_j)$

hence  $-H(Y|X) \geq -H(Y)$

$(\Rightarrow) H(Y) \geq H(Y|X)$  □.

$\longleftarrow$

Kullback-Leibler Divergence.

Let  $P = (p_1, \dots, p_n)$  &  $Q = (q_1, \dots, q_n)$  be two probability distributions. Then the KL divergence of  $Q$  from  $P$  is defined by

$$\begin{aligned}
 D(P||Q) &= \sum_{i=1}^n p_i \log \frac{p_i}{q_i} &= \sum p_i \log p_i + \sum p_i \log \frac{1}{q_i} \\
 & &= \sum p_i \log \left( \frac{1}{q_i} \right) - \underbrace{\left( -\sum p_i \log p_i \right)}_{H(P)} \\
 & &= \sum p_i \log \left( \frac{1}{q_i} \right) - H(P).
 \end{aligned}$$

By Huffman coding, an optimal compressor for  $P$ , the average code length for any message with symbol distribution  $P$  is  $H(P)$ .

Suppose you thought that the actual distribution is  $Q$  (maybe sample error)

The true distribution is  $P$ .

So the average code length =  $\sum p_i \log \frac{1}{q_i}$

$$D(P||Q) = \left( \sum p_i \log \frac{1}{q_i} \right) - \left( \sum p_i \log \frac{1}{p_i} \right).$$

We show  $D(P||Q) \geq 0$ , and  $D(P||Q) = 0$  iff  $P = Q$ .

Lemma.

$$D(P||Q) \geq 0.$$

Proof

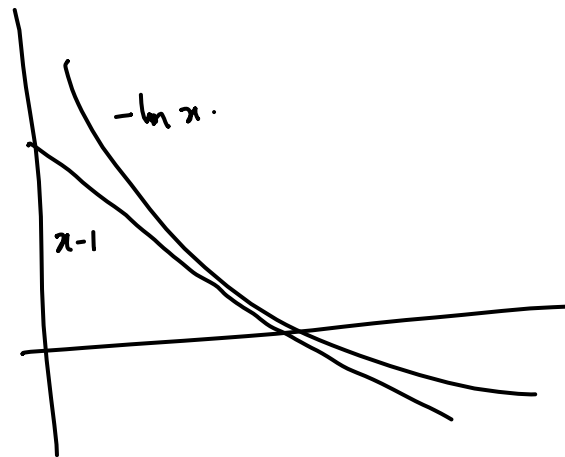
$$D(P||Q) = \sum p_i \log \frac{p_i}{q_i}$$

$$= \sum p_i \left( -\log \frac{q_i}{p_i} \right)$$

$$\geq \sum p_i \left( 1 - \frac{q_i}{p_i} \right)$$

$$= \sum_{i=1}^n (p_i - q_i)$$

$$= \left( \sum_{i=1}^n p_i \right) - \left( \sum_{i=1}^n q_i \right) = 1 - 1 = 0.$$



$$1-x \leq -\ln x$$

[fundamental inequality]