

Background

- $E[f(X)] = E_{p(x)}[f(X)] = E_{p(x,y)}[f(X)]$
 - Proof $E_{p(x,y)}[f(X)] = \sum_{x \in \mathcal{X}} \sum_{y \in \mathcal{Y}} p(x, y) f(x) = \sum_{x \in \mathcal{X}} f(x) \sum_{y \in \mathcal{Y}} p(x, y)$

$$= \sum_{x \in \mathcal{X}} f(x) p(x) = E_{p(x)}[f(X)]$$
- $P_{x|x}(x|x) = 1$.
- $p(x, x) = p(x|x)p(x) = 1p(x) = p(x)$
- Convention, based on continuity arguments: $0 \log 0 = 0$, $0 \log \frac{0}{q} = 0$, $0 \log \frac{p}{0} = \infty$.
- Let $\{p(x)\}$ and $\{q(x)\}$ be the pmf for the same alphabet set \mathcal{X} . We say $p = q$ if $\forall x \in \mathcal{X}$ $p(x) = q(x)$.
- **Convexity**
 - Def: A function $f(x)$ is said to be **convex** (convex \cup) over an interval (a, b) if $\forall x_1, \forall x_2 \in (a, b)$ and $0 \leq \mathbf{I} \leq 1$, $f(\mathbf{I}x_1 + (1-\mathbf{I})x_2) \leq \mathbf{I}f(x_1) + (1-\mathbf{I})f(x_2)$.
A function f is said to be **strictly convex** if equality holds only if $\mathbf{I} = 0$ or $\mathbf{I} = 1$.
 - Ex. (strict) $x^2, |x|, e^x, x \log x$ (for $x \geq 0$)
 - Def: A function f is concave (convex \cap) if $-f$ is convex.
 - Ex. (strict) $\log x, \sqrt{x}$ for $x \geq 0$.
 - A function is convex if it always lies below any chord.
A function is concave if it always lies above any chord.
 - If the function f has a second derivative which is non-negative (positive) everywhere, then the function is convex (strictly convex).

Proof. Let $f''(x) > 0 \forall x$. By Taylor's Theorem and Lagrange Remainder Theorem,

$$f(x) = f(x_0) + f'(x_0)(x - x_0) + \frac{f''(x^*)}{2}(x - x_0)^2 \text{ where } x^* \text{ is between } x_0 \text{ and } x.$$

$$\text{x. So, } \frac{f''(x^*)}{2}(x - x_0)^2 \geq 0 \text{ with equality iff } x = x_0.$$

$$\text{Thus, } f(x) \geq f(x_0) + f'(x_0)(x - x_0) \text{ with equality iff } x = x_0.$$

$$\forall x_1 \forall x_2 \neq x_1, \text{ let } x_0 = \mathbf{I}x_1 + (1-\mathbf{I})x_2.$$

$$\begin{aligned} \text{Let } x = x_1. \text{ Then, } f(x_1) &\geq f(x_0) + f'(x_0)(x_1 - x_0) \\ &= f(x_0) + f'(x_0)(1-\mathbf{I})(x_1 - x_2) \end{aligned}$$

with equality iff $\mathbf{I} = 1$.

Thus, $I f(x_1) \geq I f(x_0) + f'(x_0)I(1-I)(x_1 - x_2)$ iff $I = 1$ or $I = 0$

Similarly, let $x = x_2$, then

$(1-I)f(x_2) \geq (1-I)f(x_0) - f'(x_0)(1-I)I(x_1 - x_2)$ with equality iff $I = 0$ or $I = 1$.

So, $I f(x_1) + (1-I)f(x_2) \geq f(x_0)$ with equality iff $I = 1$ or $I = 0$.

- Linear functions $ax+b$ are both convex and concave.
- Jensen's inequality

Let $EX = \sum_{x \in \mathcal{X}} p(x)x$ in discrete case and $EX = \int xf(x)dx$ in the continuous case.

If f is a convex function and X is a random variable, then $E[f(X)] \geq f(EX)$.

Proof by induction.

f is convex \cup ; so, $\mathbf{a}_1 f(x_1) + \mathbf{a}_2 f(x_2)$ where $\mathbf{a}_1 + \mathbf{a}_2 = 1$. So, $E[f(X)] \geq f(EX)$ holds for $|\mathcal{X}| = 2$.

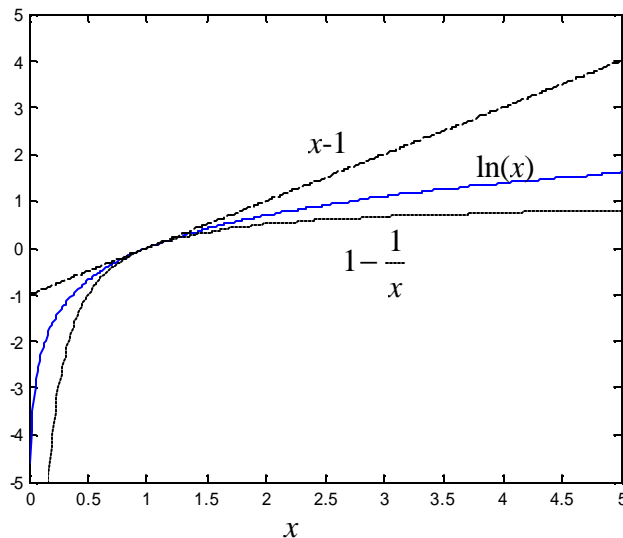
Assume it holds for $|\mathcal{X}| = k-1$, i.e., $\sum_{i=1}^{k-1} \mathbf{a}_i f(x_i) \geq f\left(\sum_{i=1}^{k-1} \mathbf{a}_i x_i\right)$ where $\sum_{i=1}^{k-1} \mathbf{a}_i = 1$.

Then,

If f is strictly convex, then $E[f(X)] = f(EX) \Rightarrow X = EX$ with probability 1, i.e., X is a constant.

- $\log(EX) \geq E[\log(X)]$

- **Fundamental inequality:** $\boxed{1 - \frac{1}{x} \leq \ln(x) \leq x - 1}$ with equality iff $x = 1$



Intro

- Axiomatic Derivation of Information Measure:

Four Postulate

A) Bayesianness: There's a function $f(\mathbf{a}, \mathbf{b})$ such that $i(x, y) = f(\mathbf{a}, \mathbf{b}) \Big|_{\substack{\mathbf{a}=p(x) \\ \mathbf{b}=p(x|y)}}$.

B) Smoothness: $f_1(\mathbf{a}, \mathbf{b}) = \frac{\partial}{\partial \mathbf{a}} f(\mathbf{a}, \mathbf{b})$ and $f_2(\mathbf{a}, \mathbf{b}) = \frac{\partial}{\partial \mathbf{b}} f(\mathbf{a}, \mathbf{b})$ exist.

C) Successive Revelation: $f(\mathbf{a}, \mathbf{g}) = f(\mathbf{a}, \mathbf{b}) + f(\mathbf{b}, \mathbf{g})$, $0 \leq \mathbf{a}, \mathbf{b}, \mathbf{g} \leq 1$.

Justification:

The information you get about X by observing (W, Z) have occurred is that provided by observation that $W = w$ plus that subsequently provided by later learning that $Z = z$.

$$i(x, (w, z)) = i(x, w) + i(x|w, z).$$

$$f\left(\underbrace{p(x)}_a, \underbrace{p(x|w, z)}_g\right) = f\left(\underbrace{p(x)}_a, \underbrace{p(x|w)}_b\right) + f\left(\underbrace{p(x|w)}_b, \underbrace{p(x|w, z)}_g\right).$$

D) Additivity over (independent experiment): $f(\mathbf{ag}, \mathbf{bd}) = f(\mathbf{a}, \mathbf{b}) + f(\mathbf{g}, \mathbf{d})$, $0 \leq \mathbf{a}, \mathbf{b}, \mathbf{g}, \mathbf{d} \leq 1$.

Justification:

Consider 2 independent experiments:

$$X \rightarrow \square \rightarrow Y$$

$$U \rightarrow \square \rightarrow V$$

$$p(x, u) = p(x)p(u)$$

Then, $p(y, v|x, u) = p(y|x, u)p(v|y, x, u) = p(y|x)p(v|u)$.

$$i((y, v), (x, u)) \text{ should } = i(x, y) + i(u, v).$$

$$\begin{aligned} i((y, v), (x, u)) &= f(p(y, v), p(y, v|x, u)) \\ &= f(p(y)p(v), p(y|x)p(v|u)) \end{aligned}$$

$$f\left(\underbrace{p(y)p(v)}_a, \underbrace{p(y|x)p(v|u)}_g\right) = f(p(y), p(y|x)) + f(p(v), p(v|u))$$

$$A) - D) \Rightarrow i(x, y) = k \log \frac{p(x, y)}{p(x)}$$

Entropy

- **Entropy of a random variable X**

- A measure of the uncertainty of the random variable
- A measure of the amount of information required on the average to describe the random variable.
- Average self information of X .
- Minimum of yes-no questions to get the value of X exactly.

- $$\boxed{\begin{matrix} 0 \\ \text{deterministic} \end{matrix} \leq H(X) = H(\{p(x)\}) = -E[\log p(X)] \leq \log |\mathcal{X}| \begin{matrix} \\ \text{uniform} \end{matrix}}$$

$$\begin{aligned} H(X) &= -\sum_{x \in \mathcal{X}} p(x) \log p(x) = -E_p[\log p(X)] \\ &= E[i(X)] \\ &\geq 0 \text{ with equality iff } \exists x \in \mathcal{X} p(x) = 1 \\ &\leq \log |\mathcal{X}| \text{ with equality iff } \forall x \in \mathcal{X} p(x) = \frac{1}{|\mathcal{X}|} \end{aligned}$$

Proof $H(X) \geq 0$ with equality iff $\exists x \in \mathcal{X} p(x) = 1$.

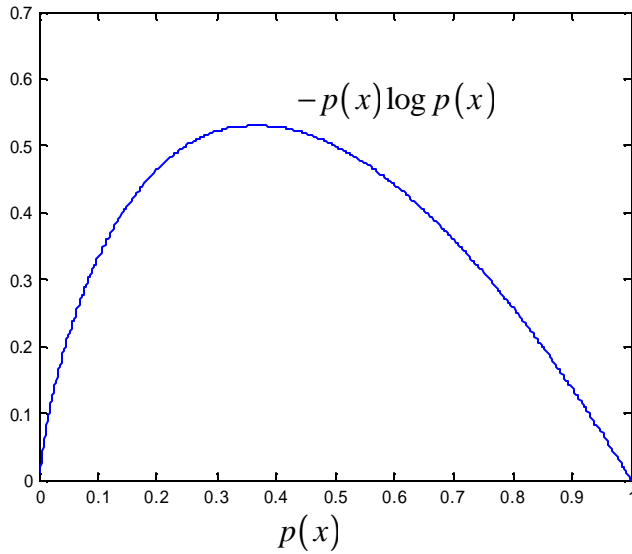
$\forall x p(x)$ and $-\log p(x) \geq 0$. Thus, $\forall x -p(x) \log p(x) \geq 0$.

Hence, $-\sum_{x \in \mathcal{X}} p(x) \log p(x) \geq 0$.

$H(X) = 0 \Leftrightarrow \forall x -p(x) \log p(x) = 0$.

But $p(x) \log p(x) = 0$ if and only if $\forall x p(x) = 0$ or 1 .

$\forall x p(x) = 0$ or 1 iff $\exists x p(x) = 0$.



Proof $H(X) \leq \log|\mathcal{X}|$ with equality iff $\forall x \in \mathcal{X} p(x) = \frac{1}{|\mathcal{X}|}$

$$\begin{aligned} H(X) - \log|\mathcal{X}| &= E[-\log p(X)] - E[\log|\mathcal{X}|] = E\left[\log\frac{1}{|\mathcal{X}|p(X)}\right] \\ &\leq E\left[\frac{1}{|\mathcal{X}|p(X)} - 1\right] = \sum_{x \in \mathcal{X}} p(x) \left(\frac{1}{|\mathcal{X}|p(X)}\right) - 1 \\ &= \sum_{x \in \mathcal{X}} \frac{1}{|\mathcal{X}|} - 1 = \frac{|\mathcal{X}|}{|\mathcal{X}|} - 1 = 0 \end{aligned}$$

$$H(X) = \log|\mathcal{X}| \Leftrightarrow \forall x \frac{1}{|\mathcal{X}|p(x)} = 1.$$

- If the base of the logarithm is b , denote the entropy as $H_b(X)$.
- [bits] if using $\log_2(\cdot)$. [nats] if using $\log_e(\cdot)$.
- A functional of the distribution of X .
- Not depend on the actual value taken by the random variable X .
- $H(X) \geq 0$
- $H_b(X) = (\log_b a) H_a(X)$.
- Ex. entropy of a fair coin toss is $-\frac{1}{2}\log\frac{1}{2} - \frac{1}{2}\log\frac{1}{2} = -\log\frac{1}{2} = 1$.
- $H(X)$ is a function of $\{p_x(x); x \in \mathcal{X}\}$. Hence, should be written as $H(\{p_x(x)\})$.
- $H(\{p(x)\})$ is concave (convex \cap) in $\{p(x)\}$

$\forall I \in [0,1]$ and any two pmf $\{p_1(x), x \in \mathcal{X}\}$ and $\{p_2(x), x \in \mathcal{X}\}$,
 $H(p^*) \geq IH(p_1) + (1-I)H(p_2)$ where $p^*(x) = Ip_1(x) + (1-I)p_2(x) \quad \forall x \in \mathcal{X}$.

Proof

$$\begin{aligned}
 & H(p^*) - IH(p_1) - (1-I)H(p_2) \\
 &= -\sum_{x \in \mathcal{X}} p^*(x) \log p(x) \\
 &\quad + I \sum_{x \in \mathcal{X}} p_1(x) \log p(x) + (1-I) \sum_{x \in \mathcal{X}} p_2(x) \log p(x) \\
 &= -\sum_{x \in \mathcal{X}} (Ip_1(x) + (1-I)p_2(x)) \log p^*(x) \\
 &\quad + I \sum_{x \in \mathcal{X}} p_1(x) \log p_1(x) + (1-I) \sum_{x \in \mathcal{X}} p_2(x) \log p_2(x) \\
 &= I \left(\sum_{x \in \mathcal{X}} p_1(x) \log \frac{p_1(x)}{p^*(x)} \right) + (1-I) \left(\sum_{x \in \mathcal{X}} p_2(x) \log \frac{p_2(x)}{p^*(x)} \right) \\
 &\geq I \left(\sum_{x \in \mathcal{X}} p_1(x) \left(1 - \frac{p^*(x)}{p_1(x)} \right) \right) + (1-I) \left(\sum_{x \in \mathcal{X}} p_2(x) \left(1 - \frac{p^*(x)}{p_1(x)} \right) \right) \\
 &= I \left(\sum_{x \in \mathcal{X}} (p_1(x) - p^*(x)) \right) + (1-I) \left(\sum_{x \in \mathcal{X}} (p_2(x) - p^*(x)) \right) \\
 &= I(1-1) + (1-I)(1-1) = 0
 \end{aligned}$$

- $\boxed{H(g(X)) \leq H(X)}$ with equality iff g is one-to-one.

Proof (1) $H(X, g(X)) = H(X) + H(g(X)|X)$ by chain rule.

But $H(g(X)|X) = 0$; so, $H(X, g(X)) = H(X)$.

(2) Also, by chain rule, $H(X, g(X)) = H(g(X)) + H(X|g(X))$.

Because $H(X|g(X)) \geq 0$ with equality iff g is one-to-one, we have

$H(X, g(X)) \geq H(g(X))$.

Combining part (1) and (2), we have $H(X) \geq H(g(X))$.

- For two random variables X and Y with a joint pmf $p(x, y)$ and marginal pmf $p(x)$ and $p(y)$.

- $\boxed{H(Y|X=x) = -\sum_{y \in \mathcal{Y}} p(y|x) \log p(y|x)}$.

- **Joint entropy:** $\boxed{H(X, Y) = -E[\log p(X, Y)]} = -\sum_{x \in \mathcal{X}} \sum_{y \in \mathcal{Y}} p(x, y) \log p(x, y)$

- **Conditional entropy** $\boxed{0 \leq H(Y|X) = -E[\log p(Y|X)] \leq H(Y)}$
 $Y=f(X)$ X, Y independent

$$\begin{aligned} H(Y|X) &= \sum_{x \in \mathcal{X}} p(x) H(Y|X=x) \\ &= - \sum_{x \in \mathcal{X}} \sum_{y \in \mathcal{Y}} p(x, y) \log p(y|x) \\ &= -E_{p(x,y)} [\log p(Y|X)] \\ &\geq 0 \end{aligned}$$

- Conditioning can only decrease entropy: $H(Y|X) \leq H(Y)$

Proof. $I(X; Y) = H(Y) - H(Y|X) \geq 0$.

- $H(X|X) = 0$

Proof $p(X=y|X=x) = \begin{cases} 1 & , y=x \\ 0 & , y \neq x \end{cases}$

$$p(X=y, X=x) = \begin{cases} p(x) & , y=x \\ 0 & , y \neq x \end{cases}$$

$$\begin{aligned} H(X|X) &= - \sum_{x \in \mathcal{X}} \sum_{y \in \mathcal{X}} p(x, y) \log p(y|x) \\ &= - \sum_{x \in \mathcal{X}} \left(p(x, x) \log p(x|x) + \sum_{\substack{y \in \mathcal{X} \\ y \neq x}} p(x, y) \log p(y|x) \right) \\ &= - \sum_{x \in \mathcal{X}} \left(p(x) \log 1 + \sum_{\substack{y \in \mathcal{X} \\ y \neq x}} 0 \log 0 \right) = 0 \end{aligned}$$

- $H(g(X)|X) = 0$

Proof $p(g(X)=y|X=x) = \begin{cases} 1 & , y=g(x) \\ 0 & , y \neq g(x) \end{cases}$

$$p(g(X)=y, X=x) = \begin{cases} p(x) & , y=g(x) \\ 0 & , y \neq g(x) \end{cases}$$

$$\begin{aligned}
H(g(X)|X) &= -\sum_{x \in \mathcal{X}} \sum_{y \in g(\mathcal{X})} p_{X,g(X)}(x,y) \log p_{g(X)|X}(y|x) \\
&= -\sum_{x \in \mathcal{X}} \left(\sum_{\substack{y \in g(\mathcal{X}) \\ y=g(x)}} p_{X,g(X)}(x,y) \log p(y|x) + \right. \\
&\quad \left. \sum_{\substack{y \in g(\mathcal{X}) \\ y \neq g(x)}} p(x,y) \log p_{g(X)|X}(y|x) \right) \\
&= -\sum_{x \in \mathcal{X}} \left(\sum_{\substack{y \in g(\mathcal{X}) \\ y=g(x)}} p(x) \log 1 + \sum_{\substack{y \in g(\mathcal{X}) \\ y \neq g(x)}} 0 \log 0 \right) \\
&= 0
\end{aligned}$$

- $H(g(\bar{X})|\bar{X}) = 0$

Proof $H(g(\bar{X})|\bar{X} = x) = -\sum_{\bar{y}} p_{g(\bar{X})|\bar{X}}(\bar{y}|x) \log p_{g(\bar{X})|\bar{X}}(\bar{y}|x)$

$$\begin{aligned}
&= -\sum_{\bar{y}=g(\bar{x})} p_{g(\bar{X})|\bar{X}}(\bar{y}|x) \log p_{g(\bar{X})|\bar{X}}(\bar{y}|x) \\
&\quad - \sum_{\bar{y} \neq g(\bar{x})} p_{g(\bar{X})|\bar{X}}(\bar{y}|x) \log p_{g(\bar{X})|\bar{X}}(\bar{y}|x) \\
&= -p_{g(\bar{X})|\bar{X}}(g(\bar{x})|\bar{x}) \log p_{g(\bar{X})|\bar{X}}(g(\bar{x})|\bar{x}) - \sum_{\bar{y} \neq g(\bar{x})} 0 \log 0 \\
&= -1 \log 1 + 0 = 0
\end{aligned}$$

$$H(g(\bar{X})|\bar{X}) = \sum_{\bar{x}} p_{\bar{X}}(\bar{x}) H(g(\bar{X})|\bar{X} = x) = \sum_{\bar{x}} p_{\bar{X}}(\bar{x}) 0 = 0.$$

- Chain rule: $H(X,Y) = H(X) + H(Y|X)$.

Proof $p(x,y) = p(x)p(y|x)$

- $H(X,Y|Z) = H(X|Z) + H(Y|X,Z)$.

Proof $p(x,y|z) = p(x|z)p(y|x,z)$.

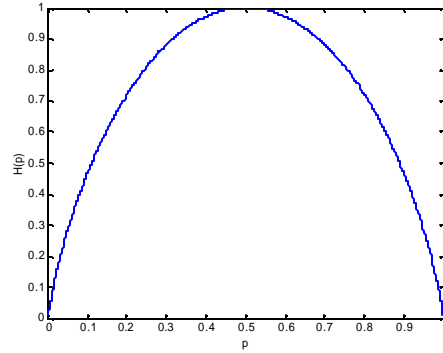
- In general $H(Y|X) \neq H(X|Y)$

- $H(Y) - H(Y|X) = H(X) - H(X|Y)$

Proof $p(x,y) = p(x)p(y|x) = p(y)p(x|y)$

- $H(\{p(x)p(y)\}) = H(\{p(y)\}) + H(\{p(x)\})$.

- Def: $H(p) = -p \log p - (1-p) \log (1-p)$



- Entropy of a collection of random variables.
 - Let X_1^n represents X_1, X_2, \dots, X_n .
 - Let X_1, X_2, \dots, X_n be drawn according to $p(x_1, x_2, \dots, x_n)$.
 - **Joint entropy**: $H(X_1^n) = -E[\log p(X_1^n)]$

$$H(X_1, X_2, \dots, X_n) = -\sum_{x_1 \in \mathcal{X}_1} \sum_{x_2 \in \mathcal{X}_2} \cdots \sum_{x_n \in \mathcal{X}_n} p(x_1, x_2, \dots, x_n) \log p(x_1, x_2, \dots, x_n)$$

$$= -E[\log p(X_1, X_2, \dots, X_n)]$$
 - $H(X_1^n, Y) = H(X_1^n) + H(Y|X_1^n)$

Proof $p(x_1^n, y) = p(x_1^n) p(y|x_1^n)$.
 - **Chain rule for entropy**: $H(X_1, X_2, \dots, X_n) = \sum_{i=1}^n H(X_i|X_{i-1}, \dots, X_1)$.

$$H(X_1^n) = \sum_{i=1}^n H(X_i|X_1^{i-1})$$

Proof $p(x_1^n) = \prod_{i=1}^n p(x_i|x_1^{i-1})$.
- $$H(X_1^n) = -E[\log p(X_1^n)] = \sum_{i=1}^n H(X_i|X_1^{i-1}) \leq \sum_{i=1}^n H(X_i)$$

X_i 's are independent
- $H(X, Y) = H(X) + H(Y|X) = H(Y) + H(X|Y)$
- $H(X_1^n|Y) = \sum_{i=1}^n H(X_i|X_1^{i-1}, Y)$
- $H(X, Y|Z) = H(X|Z) + H(Y|X, Z) = H(Y|Z) + H(X|Y, Z)$
- $H(X) \leq \log |\mathcal{X}|$, where $|\mathcal{X}|$ denotes the number of elements in the range of X with equality if and only if X has a uniform distribution over \mathcal{X} .

- Conditioning reduce entropy: $H(X|Y) \leq H(X)$ with equality iff X and Y are independent.
Proof $I(X;Y) = H(X) - H(X|Y) \geq 0$ with equality iff X and Y are independent.
- Knowing another random variable Y can only reduce the uncertainty in X .
- No general comparison between $H(X|Y = y)$ and $H(X)$.
- Independence bound on entropy: $H(X_1^n) \leq \sum_{i=1}^n H(X_i)$ with equality if and only if the X_i are independent.

Relative Entropy

- **Relative entropy** / Kullback Leibler “distance” between two probability mass functions $p(x)$ and $q(x)$

$$0 \leq D(p||q) = \sum_{x \in \mathcal{X}} p(x) \log \frac{p(x)}{q(x)} = E \left[\log \frac{p(X)}{q(X)} \right] \text{ [bits]}$$

- A measure of the inefficiency of assuming that the distribution is q when the true distribution is p .

If we knew the true distribution $\{p(x)\}$ of the random variable, then we could construct a code with average description length $H(p)$. If, instead, we used the code for a distribution q , we would need $H(p) + D(p||q)$ bits on the average to describe the random variable.

$$\text{Proof. } (-p(x) \log q(x)) - (-p(x) \log p(x)) = p(x) \log \frac{p(x)}{q(x)}. \text{ So,}$$

$$E[-\log q(X)] = E[-\log p(X)] + E \left[\log \frac{p(X)}{q(X)} \right].$$

- ≥ 0 , $= 0$ iff $p = q$.

$$\begin{aligned} \text{Proof. } D(p||q) &= \sum_x p(x) \log \frac{p(x)}{q(x)} \geq \sum_x p(x) \left(1 - \frac{q(x)}{p(x)} \right) \\ &\geq \sum_x (p(x) - q(x)) = 1 - 1 = 0 \end{aligned}$$

- Note that this just means if we have two vectors \bar{u}, \bar{v} with the same lengths, each have elements which summed to 1. Then, $\sum_i u_i \log \frac{u_i}{v_i} \geq 0$.

- In general, $D(p||q) \neq D(q||p)$.

- Log sum inequality: for non-negative numbers, a_1, \dots, a_n and b_1, \dots, b_n ,

$$\sum_{i=1}^n a_i \log \frac{a_i}{b_i} \geq \left(\sum_{i=1}^n a_i \right) \log \frac{\sum_{i=1}^n a_i}{\sum_{i=1}^n b_i} \quad \text{with equality iff } \frac{a_i}{b_i} = \text{constant } \forall i.$$

Proof. Define $a'_i = \frac{a_i}{\sum_{i=1}^n a_i} = \frac{a_i}{A}$ and $b'_i = \frac{b_i}{\sum_{i=1}^n b_i} = \frac{b_i}{B}$. Then, from

$$D(p \| q) = \sum_{x \in \mathcal{X}} p(x) \log \frac{p(x)}{q(x)} \geq 0, \text{ we have}$$

$$\begin{aligned} 0 &\leq \sum_{i=1}^n a'_i \log \frac{a'_i}{b'_i} = \sum_{i=1}^n \frac{a_i}{A} \log \frac{\frac{a_i}{A}}{\frac{b_i}{B}} = \sum_{i=1}^n \frac{a_i}{A} \log \frac{a_i B}{b_i A} = \sum_{i=1}^n \frac{a_i}{A} \log \frac{a_i}{b_i} - \sum_{i=1}^n \frac{a_i}{A} \log \frac{A}{B} \\ &= \frac{1}{A} \sum_{i=1}^n a_i \log \frac{a_i}{b_i} - \log \frac{A}{B} = \frac{1}{A} \left(\sum_{i=1}^n a_i \log \frac{a_i}{b_i} - A \log \frac{A}{B} \right) \end{aligned}$$

$$\text{Thus, } \sum_{i=1}^n a_i \log \frac{a_i}{b_i} \geq (A) \log \frac{A}{B}.$$

$$\text{Equality iff } a'_i = b'_i \forall i \Leftrightarrow \frac{a_i}{b_i} = \frac{A}{B} \forall i.$$

- $a \log \frac{a}{0} = \infty$ if $a > 0$, and $0 \log \frac{0}{0} = 0$.
- Not a true distance since symmetry and triangle inequality fail. Nonetheless, it is often useful to think of it as a distance between distributions.
- $D(p \| q)$ is convex \cup in the pair (p, q) .

If (p_1, q_1) and (p_2, q_2) are two pairs of probability mass functions, then

$$D(\mathbf{I} p_1 + (1 - \mathbf{I}) p_2 \| \mathbf{I} q_1 + (1 - \mathbf{I}) q_2) \leq \mathbf{I} D(p_1 \| q_1) + (1 - \mathbf{I}) D(p_2 \| q_2) \quad \forall 0 \leq \mathbf{I} \leq 1.$$

- For fixed p , $D(q \| p)$ is a convex \cup function of q .

$$D(\mathbf{I} q_1 + (1 - \mathbf{I}) q_2 \| p) \leq \mathbf{I} D(q_1 \| p) + (1 - \mathbf{I}) D(q_2 \| p).$$

Proof.

$$\begin{aligned}
p_0(x) \log \frac{p_0(x)}{q_0(x)} &= (\mathbf{I} p_1(x) + (1-\mathbf{I}) p_2(x)) \log \frac{\mathbf{I} p_1(x) + (1-\mathbf{I}) p_2(x)}{\mathbf{I} q_1(x) + (1-\mathbf{I}) q_2(x)} \\
&= A \log \frac{A}{B} \\
&\leq \sum_{i=1}^n a_i \log \frac{a_i}{b_i} \\
&= \mathbf{I} p_1(x) \log \frac{\mathbf{I} p_1(x)}{\mathbf{I} q_1(x)} + (1-\mathbf{I}) p_2(x) \log \frac{(1-\mathbf{I}) p_2(x)}{(1-\mathbf{I}) q_2(x)} \\
&= \mathbf{I} p_1(x) \log \frac{p_1(x)}{q_1(x)} + (1-\mathbf{I}) p_2(x) \log \frac{p_2(x)}{q_2(x)}
\end{aligned}$$

- **Conditional relative entropy** $D(p(y|x) \| q(y|x))$

$$D(p(y|x) \| q(y|x)) = E \left[\log \frac{p(Y|X)}{q(Y|X)} \right] = \sum_x p(x) \sum_y p(y|x) \log \frac{p(y|x)}{q(y|x)}.$$

- $D(p(\bar{x}) \| q(\bar{x})) \geq 0$

Proof. Map $\mathcal{X}^n \xrightarrow{\text{onto}} \{i : i = 1, \dots, |\mathcal{X}|^n\}$, $p(\bar{x}) \rightarrow u_i$, $q(\bar{x}) \rightarrow v_i$. Then, $\sum_i u_i = 1$, and

$$\sum_i v_i = 1. \text{ Use } \sum_i u_i \log \frac{u_i}{v_i} \geq 0.$$

- $D(p(x|z) \| q(x|z)) \geq 0$

Proof.. For any given z , $\sum_x p(x|z) = 1$, and $\sum_x q(x|z) = 1$; thus, $\sum_x p(x|z) \log \frac{p(x|z)}{q(x|z)} \geq$

$$0. \quad D(p(x|z) \| q(x|z)) = \sum_z p(z) \underbrace{\sum_x p(x|z) \log \frac{p(x|z)}{q(x|z)}}_{\geq 0}.$$

- $D(p(\bar{x}|\bar{z}) \| q(\bar{x}|\bar{z})) \geq 0$

- Chain rule for relative entropy:

$$D(p(x, y) \| q(x, y)) = D(p(x) \| q(x)) + D(p(y|x) \| p(y|x))$$

$$\text{Proof } \frac{p(x, y)}{q(x, y)} = \frac{p(x)p(y|x)}{q(x)q(y|x)} = \frac{p(x)}{q(x)} \frac{p(y|x)}{q(y|x)}.$$

Mutual Information

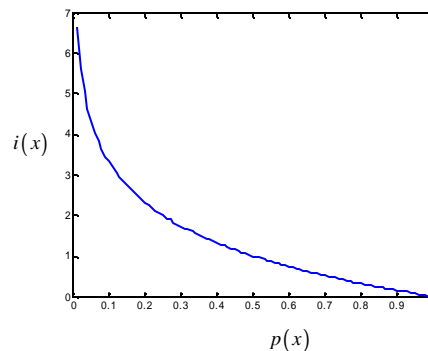
- $i(x, y) = \log \frac{p(x, y)}{p(x)p(y)} = \log \frac{p(x|y)}{p(x)} = \log \frac{p(y|x)}{p(y)}$; can be negative.

$$i(x, y) = i(y, x); \text{ more precisely } i(X = x, Y = y) = i(Y = y, X = x).$$

If $p(x|y) = 1$, the mutual info is equivalent to the self-information of symbol x .

$$i(x) = i(x, y) \Big|_{p(x|y)=1} = \log \frac{p(x|y)}{p(x)} \Big|_{p(x|y)=1} = \log \frac{1}{p(x)} = -\log p(x).$$

$$i(x) = i(x, x) = \log \frac{p(x|x)}{p(x)} = \log \frac{1}{p(x)} = -\log p(x).$$



- Average Mutual information
 - A measure of the amount of information that one random variable contains about another random variable. ($H(X|Y) = H(X) - I(X;Y)$).
 - The reduction in the uncertainty of one random variable due to the knowledge of the other.
 - A special case relative entropy.
 - Need on average $H(\{p(x, y)\})$ info bits to describe (x, y) . If instead, assume that X and Y are independent, then would need on average $H(\{p(x)p(y)\}) + D(p(x, y) \| p(x)p(y))$ info bits to describe (x, y) .
- **Average mutual information**

$$\boxed{0 \leq I(X;Y) = E \left[\log \frac{P(X,Y)}{p(X)q(Y)} \right] = E \left[\log \frac{P(X|Y)}{p(X)} \right] = E \left[\log \frac{Q(Y|X)}{q(Y)} \right]}.$$

iff independent

$$\begin{aligned}
I(X;Y) &= \sum_{x \in \mathcal{X}} \sum_{y \in \mathcal{Y}} p(x,y) \log \frac{p(x,y)}{p(x)p(y)} \\
&= E_{p(x,y)} \left[\log \frac{p(X,Y)}{p(X)p(Y)} \right] = E[i(X;Y)] \\
&= D(p(x,y) \| p(x)p(y)) \\
&\geq 0 \text{ with equality iff } X \text{ and } Y \text{ are independent}
\end{aligned}$$

Proof

$$\begin{aligned}
I(X;Y) &= \sum_{x \in \mathcal{X}} \sum_{y \in \mathcal{Y}} p(x,y) \log \frac{p(x,y)}{p(x)p(y)} \\
&= \frac{1}{\ln(2)} \sum_{x \in \mathcal{X}} \sum_{y \in \mathcal{Y}} p(x,y) \ln \frac{p(x,y)}{p(x)p(y)} \\
&\geq \frac{1}{\ln(2)} \sum_{x \in \mathcal{X}} \sum_{y \in \mathcal{Y}} p(x,y) \left[1 - \left(\frac{p(x,y)}{p(x)p(y)} \right)^{-1} \right] \\
&= \frac{1}{\ln(2)} \sum_{x \in \mathcal{X}} \sum_{y \in \mathcal{Y}} [p(x,y) - p(x)p(y)] = \frac{1}{\ln(2)} (1-1) = 0
\end{aligned}$$

- | |
|--|
| $ \begin{aligned} I(X;Y) &= H(X) + H(Y) - H(X,Y) \\ &= H(Y) - H(Y X) = H(X) - H(X Y) \end{aligned} $ |
|--|

Proof. $\frac{p(x,y)}{p(x)p(y)} = \frac{p(y|x)}{p(y)} = \frac{p(x|y)}{p(x)}$.

- $I(X;X) = H(X) \Rightarrow$ entropy = self-information.

Proof. $I(X;X) = H(X) - H(X|X) = H(X)$.

- $I(X;Y) = I(Y;X)$

- The X says, on average, as much about Y as Y says, on average, about X.

- Conditional mutual information** of random variables X and Y given Z,

$ \begin{aligned} I(X;Y Z) &= H(X Z) - H(X Y,Z) \\ &= E_{p(x,y,z)} \log \frac{p(X,Y Z)}{P(X Z)p(Y Z)} \\ &\geq 0 \text{ with equality iff } X \text{ and } Y \text{ are conditionally independent given } Z \end{aligned} $

Proof. $\frac{p(x,y|z)}{P(x|z)p(y|z)} = \frac{1}{p(x|z)} \frac{p(y|z)p(x|y,z)}{P(y|z)} = \frac{1}{p(x|z)} p(x|y,z)$.

Proof. $I(X;Y|Z) = D(p(x,y|z) \parallel q(x,y|z)) \geq 0$ where $q(x,y|z) = p(x|z)p(y|z)$.

- $$I(X_1^n; Y) = E \left[\log \frac{p(X_1^n, Y)}{p(X_1^n)p(Y)} \right] = \sum_{x_1 \in \mathcal{X}_1} \sum_{x_2 \in \mathcal{X}_2} \cdots \sum_{x_n \in \mathcal{X}_n} \sum_{y \in \mathcal{Y}} p(x_1^n, y) \log \frac{p(x_1^n, y)}{p(x_1^n)p(y)}$$

$$= H(X_1^n) + H(Y) - H(X_1^n, Y)$$

$$= H(X_1^n) - H(X_1^n | Y) = H(Y) - H(Y | X_1^n)$$

Proof
$$E \left[\log \frac{p(X_1^n, Y)}{p(X_1^n)p(Y)} \right] = E[\log p(X_1^n, Y)] - E[\log p(X_1^n)] - E[\log p(Y)]$$

$$= -H(X_1^n, Y) + H(X_1^n) + H(Y)$$

- Chain rule for information:

$$I(X_1, X_2, \dots, X_n; Y) = \sum_{i=1}^n I(X_i; Y | X_{i-1}, X_{i-1}, \dots, X_1)$$

$$I(X_1^n; Y) = \sum_{i=1}^n I(X_i; Y | X_1^{i-1})$$

Proof
$$I(X_1^n; Y) = H(X_1^n) - H(X_1^n | Y) = \sum_{i=1}^n H(X_i | X_1^{i-1}) - \sum_{i=1}^n H(X_i | X_1^{i-1}, Y)$$

$$= \sum_{i=1}^n (H(X_i | X_1^{i-1}) - H(X_i | X_1^{i-1}, Y)) = \sum_{i=1}^n I(X_i; Y | X_1^{i-1})$$

$$H(X_i | X_1^{i-1}) - H(X_i | X_1^{i-1}, Y) = H(X | Z) - H(X | Z, Y) = I(X; Y | Z)$$

$$= I(X_i; Y | X_1^{i-1})$$

$$I(X_1, X_2; Y) = I(X_1; Y) + I(X_2; Y | X_1)$$

Stationary Information Sources

- Consider stationary source $\{U(k)\}$. Common alphabet \mathcal{U} .

- $$H(U_1^n) = H(U_k^{k+n-1})$$

- Per letter entropy of an L -block:

- $$H_L = \frac{H(U_1^L)}{L} = \frac{H(U_k^{k+L-1})}{L}$$

- $$H_{\text{Volumetric}} = \lim_{L \rightarrow \infty} H_L$$

- Incremental entropy change

- $\boxed{h_L = H(U_L|U_1^{L-1})} = H(U_1^L) - H(U_1^{L-1})$
- $H_{\text{Incremental}} = \lim_{L \rightarrow \infty} h_L$.
- For stationary Markov chain (the initial state of the Markov chain is drawn according to a stationary distribution.):

$$h_{L,\text{markov}} = H(U_L|U_1^{L-1}) = H(U_L|U_{L-1}) = H(U_2|U_1) \quad \forall L \geq 2.$$

$$H = \lim_{L \rightarrow \infty} h_L = H(U_2|U_1) = - \sum_{u_1, u_2} p(u_1) p(u_2|u_1) \log p(u_2|u_1).$$

- $h_1 = H_1 = H(U_1) = H(U_k)$
- Both h_L and H_L are non-increasing \downarrow function of L , converging to same limit H .

$$\boxed{\lim_{L \rightarrow \infty} H_L = \lim_{L \rightarrow \infty} h_L = H}. \text{ Also, } \boxed{h_L \leq H_L} \left(H(X_n|X_1^{n-1}) \leq \frac{H(X_1^n)}{n} \right).$$

Proof. $h_L \leq h_{L-1}$.

$$h_L = H(U_L|U_1^{L-1}) \leq H(U_L|U_2^{L-1}) \underset{\text{stationary}}{=} H(U_{L-1}|U_1^{L-2}) = h_{L-1}.$$

Proof $\boxed{H(U_1^L) = \sum_{k=1}^L h_k}$

$$H(U_1^L) = \sum_{k=1}^L H(U_k|U_1^{k-1}) = \sum_{k=1}^L h_k$$

Proof $h_L \leq H_L$

$$H(U_1^L) = \sum_{k=1}^L h_k \geq \sum_{k=1}^L h_L = Lh_L. \text{ So, } h_L \leq \frac{H(U_1^L)}{L} = H_L.$$

Proof $H_L \leq H_{L-1}$

$$\begin{aligned} H_L &= \frac{1}{L} H(U_1^L) + \frac{h_L}{L} = \frac{L-1}{L} \frac{H(U_1^{L-1})}{L-1} + \frac{h_L}{L} = \frac{L-1}{L} H_{L-1} + \frac{h_L}{L} \\ &\leq \frac{L-1}{L} H_{L-1} + \frac{H_L}{L} \end{aligned}$$

$$\frac{L-1}{L} H_L \leq \frac{L-1}{L} H_{L-1}$$

$$H_L \leq H_{L-1}$$

Proof $\lim_{L \rightarrow \infty} H_L = \lim_{L \rightarrow \infty} h_L$

$$\text{From } h_L \leq H_L, \lim_{L \rightarrow \infty} H_L \geq \lim_{L \rightarrow \infty} h_L.$$

$$H_{L+M} = \frac{H(U_1^{L+M})}{L+M} = \frac{\sum_{k=L}^{L+M} H(U_k | U_1^{k-1}) + H(U_1^{L-1})}{L+M} = \frac{\sum_{k=L}^{L+M} h_k + H(U_1^{L-1})}{L+M}$$

$$\leq \frac{\sum_{k=L}^{L+M} h_L + H(U_1^{L-1})}{L+M} = \frac{(M+1)h_L + H(U_1^{L-1})}{L+M}$$

Take $M \rightarrow \infty$. $\lim_{M \rightarrow \infty} H_{L+M} \leq h_L$

Take $L \rightarrow \infty$. $\lim_{L \rightarrow \infty} \lim_{M \rightarrow \infty} H_{L+M} = \lim_{L \rightarrow \infty} H_L \leq \lim_{L \rightarrow \infty} h_L$.

- **Entropy rate** of stationary source $\{U_k\}$

$$H(\{U_\ell\}) = H_U = \lim_{L \rightarrow \infty} \frac{H(U_1^L)}{L} = \lim_{L \rightarrow \infty} H(U_L | U_1^{L-1}).$$

- So, for stationary source, the entropy $H(U_1^L)$ grows (asymptotically) linearly with L at a rate H_U .
- For stationary Markov chain of order r , $H_U = H(U_{r+1} | U_1^r) = h_{r+1}$.
- For stationary Markov chain of order 1, $H_U = H(U_2 | U_1) = h_2 < H(U_1) = H(U_2)$.
- Let $\{X_i\}$ be a stationary Markov chain with stationary distribution \bar{u} and transition matrix P . Then, the entropy rate is $H = -\sum_{ij} u_i P_{ij} \log P_{ij}$.

$$P_{ij} = \Pr[\text{Next state is } j | \text{Current state is } i] = \Pr[X_2 = j | X_1 = i].$$

$$u_i = \Pr[X_1 = i].$$

- More than one communicating class: $H_U = \sum_i \Pr[\text{class}_i] H(U_2 | U_1, \text{class}_i)$.
- The best achievable data compression.

Variable-length (VL) lossless source codes

- Stationary discrete memoryless?? source $\{U_k\}$, finite alphabet \mathcal{U} .
- A variable-length D-ary source codes is a mapping $f: \mathcal{U} \rightarrow \{0, \dots, D-1\}^*$
 - binary $D = 2$
 - $D =$ coding alphabet cardinality.
- Def: f is **uniquely decipherable** if $\forall M \forall N$ and any $\underline{U} = (U_1, \dots, U_M)$, $\underline{U}' = (U'_1, \dots, U'_N)$, $f(\underline{U}) = f(\underline{U}') \Rightarrow \underline{U} = \underline{U}'$.

(No two distinct source strings get mapped into same code string.)

- $\ell(u)$ is length of D-ary string $f(u)$.
- Def: $\bar{\ell}$ = the mean code word length = $E[\ell(u)] = \sum_{u \in \mathcal{U}} p(u)\ell(u)$.
- Optimum = $\min \bar{\ell}$, uniquely decipherable.
- **Morse's principle**: To minimize $\bar{\ell}$, if $p(u) = \Pr[U_k = u]$ is small, make $\ell(u)$ large, and conversely.
- **Prefix code**: no short code word is prefix of a longer one. \Rightarrow uniquely decipherable
- **Kraft Inequality** (KI): $\sum_{u \in \mathcal{U}} D^{-\ell(u)} \leq 1$.
- Property of a length set $\{\ell(u), u \in \mathcal{U}\}$

(1) If $\{\ell(u), u \in \mathcal{U}\}$ satisfying KI, then there exists a prefix code (hence, UD) with these lengths.

(2) Every D-ary UD code has word lengths $\{\ell(u), u \in \mathcal{U}\}$ that satisfy KI.

$\boxed{\text{KI} \Rightarrow \exists \text{ prefix (UD)}}$, $\boxed{\text{UD (including prefix)} \Rightarrow \text{KI}}$

We are looking for a UD code with $\min \bar{\ell}$. Suppose we find one. Because it is UD, from (2), it's length set satisfies KI. Then, (1) tells us that there exists a prefix code with exactly the same length set and thus also minimize $\bar{\ell}$. So, (1) and (2) let us restrict search for optimal code to prefix codes.

Proof (2)

Consider L -vector $\underline{u} = (u_1, u_2, \dots, u_L)$.

$f(\underline{u}) = f(u_1)f(u_2) \cdots f(u_L)$. $\ell(\underline{u}) = \ell(u_1) + \ell(u_2) + \cdots + \ell(u_L)$.

$$\begin{aligned} \sum_{\underline{u} \in \mathcal{U}^L} D^{-\ell(\underline{u})} &= \sum_{\underline{u} \in \mathcal{U}^L} D^{-(\ell(u_1) + \ell(u_2) + \cdots + \ell(u_L))} = \sum_{\underline{u} \in \mathcal{U}^L} D^{-\ell(u_1)} D^{-\ell(u_2)} \cdots D^{-\ell(u_L)} \\ &= \sum_{u_1 \in \mathcal{U}} \sum_{u_2 \in \mathcal{U}} \cdots \sum_{u_L \in \mathcal{U}} D^{-\ell(u_1)} D^{-\ell(u_2)} \cdots D^{-\ell(u_L)} \\ &= \left(\sum_{u_1 \in \mathcal{U}} D^{-\ell(u_1)} \right) \left(\sum_{u_2 \in \mathcal{U}} D^{-\ell(u_2)} \right) \cdots \left(\sum_{u_L \in \mathcal{U}} D^{-\ell(u_L)} \right) = \left(\sum_{u \in \mathcal{U}} D^{-\ell(u)} \right)^L \end{aligned}$$

So, we have $\sum_{\underline{u} \in \mathcal{U}^L} D^{-\ell(\underline{u})} = \left(\sum_{u \in \mathcal{U}} D^{-\ell(u)} \right)^L$. (*)

Let $\ell_{\min} = \min_u \ell(u)$, $\ell_{\max} = \max_u \ell(u)$.

A_n = the number of $\underline{u} \in \mathcal{U}^L$ such that $\ell(\underline{u}) = n$.

Note that $L\ell_{\min} \leq n \leq L\ell_{\max}$. The maximum is attained when every u_k in \underline{u} corresponds to ℓ_{\max} . The minimum is attained when every u_k in \underline{u}

corresponds to ℓ_{\min} . Also, $\sum_{\underline{u} \in \mathcal{U}^L} D^{-\ell(\underline{u})} = \sum_{n=L\ell_{\min}}^{L\ell_{\max}} A_n D^{-n}$.

UD implies that $A_n \leq D^n$. (There are only D^n different code sequences of length n . If $A_n > D^n$, then there are at least two \underline{u} which map to the same code sequence.)

$$\begin{aligned} \sum_{\underline{u} \in \mathcal{U}^L} D^{-\ell(\underline{u})} &= \sum_{n=L\ell_{\min}}^{L\ell_{\max}} A_n D^{-n} \leq \sum_{n=L\ell_{\min}}^{L\ell_{\max}} D^n D^{-n} = \sum_{n=L\ell_{\min}}^{L\ell_{\max}} 1 = L\ell_{\max} - L\ell_{\min} + 1 \\ &\leq L\ell_{\max} \end{aligned}$$

So, UD requires $\sum_{\underline{u} \in \mathcal{U}^L} D^{-\ell(\underline{u})} \leq L\ell_{\max}$. (**).

Combining (*) and (**), we have $\left(\sum_{\underline{u} \in \mathcal{U}^L} D^{-\ell(\underline{u})} \right)^{\frac{1}{L}} \leq L\ell_{\max}$, or equivalently,

$$\sum_{\underline{u} \in \mathcal{U}^L} D^{-\ell(\underline{u})} \leq (L\ell_{\max})^{\frac{1}{L}}. \text{ This has to be true for all } L.$$

Note that $(L\ell_{\max})^{\frac{1}{L}}$ is strictly decreasing as L increase. $\lim_{L \rightarrow \infty} (L\ell_{\max})^{\frac{1}{L}} = 1$. Thus,

$(L\ell_{\max})^{\frac{1}{L}}$ can get arbitrary close to 1 from above. If $\sum_{\underline{u} \in \mathcal{U}^L} D^{-\ell(\underline{u})} > 1$, there will exist

L_0 such that $\sum_{\underline{u} \in \mathcal{U}^L} D^{-\ell(\underline{u})} > (L\ell_{\max})^{\frac{1}{L}}$ for all $L > L_0$. So, to have $\sum_{\underline{u} \in \mathcal{U}^L} D^{-\ell(\underline{u})} \leq (L\ell_{\max})^{\frac{1}{L}}$, need $\sum_{\underline{u} \in \mathcal{U}} D^{-\ell(\underline{u})} \leq 1$.

Proof (1) by induction

We will show that we can embed these KI satisfying word lengths as the terminal nodes in a D -ary branching tree.

Note that putting a terminal node on level ℓ prunes away $D^{L-\ell}$ nodes from level $L \geq \ell$.

Suppose each u such that $\ell(u) \leq \ell - 1$ has been assigned a terminal node on level $\ell(u)$. Now, we want to assign terminal nodes on level ℓ to all u such that $\ell(u) = \ell$. We then need there to be at least $|\{u : \ell(u) = \ell\}|$ nodes on level ℓ not yet pruned away.

Originally, there were D^ℓ nodes on level ℓ . We have pruned away $\sum_{\ell(u) \leq \ell-1} D^{\ell-\ell(u)}$

of them. So, we need $D^\ell - \sum_{\ell(u) \leq \ell-1} D^{\ell-\ell(u)} \geq |\{u : \ell(u) = \ell\}|$.

Trick: $|\{u : \ell(u) = \ell\}| = \sum_{\ell(u)=\ell-1} 1 = \sum_{\ell(u)=\ell-1} D^{\ell-\ell(u)}$.

So, need $1 \geq \sum_{\ell(u) \leq \ell} D^{-\ell(u)}$.

If $\{\ell(u), u \in \mathcal{U}\}$ satisfying KI, then $\sum_{u \in \mathcal{U}} D^{-\ell(u)} \leq 1$, and therefore,

$$\sum_{\ell(u) \leq \ell} D^{-\ell(u)} \leq \sum_{u \in \mathcal{U}} D^{-\ell(u)} \leq 1.$$

- For any UD D -ary code, and any distribution $\{p(u), u \in \mathcal{U}\}$,

$$\bar{\ell} \geq H_D(\{p(u)\}) = -\sum_{u \in \mathcal{U}} p(u) \log_D p(u).$$

Proof. $\bar{\ell} - H_D(\{p(u)\})$

$$= \sum_{u \in \mathcal{U}} p(u) \ell(u) + \sum_{u \in \mathcal{U}} p(u) \log_D p(u) = \sum_{u \in \mathcal{U}} p(u) (\ell(u) + \log_D p(u))$$

$$= \sum_{u \in \mathcal{U}} p(u) (\log_D D^{\ell(u)} + \log_D p(u)) = \sum_{u \in \mathcal{U}} p(u) (\log_D p(u) D^{\ell(u)})$$

$$= \frac{1}{\ln D} \sum_{u \in \mathcal{U}} p(u) \ln p(u) D^{\ell(u)}$$

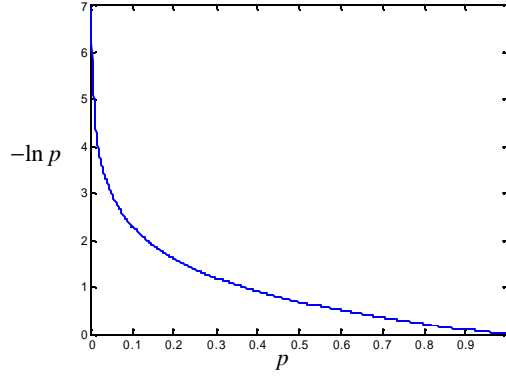
$$\geq \frac{1}{\ln D} \sum_{u \in \mathcal{U}} p(u) \left(1 - \frac{1}{p(u) D^{\ell(u)}}\right) = \frac{1}{\underbrace{\ln D}_{>0}} \left(\sum_{u \in \mathcal{U}} p(u) - \sum_{u \in \mathcal{U}} D^{-\ell(u)}\right)$$

$$\stackrel{(a)}{\geq} \frac{1}{\ln D} (1-1) = 0$$

(a) Code is UD; thus length set satisfies KI.

- $\boxed{KI \Rightarrow \bar{\ell} \geq H_D(\{p_i\})}$
 $\ell_i = \log_D p_i$

- Shannon-Fano codes: $\{\ell(u) = \lceil -\log_D p(u) \rceil = \lceil i(u) \rceil, u \in \mathcal{U}\}$.



- This length assignment is possible because it satisfies KI.

Proof. Because $\lceil -\log_D p(u) \rceil \geq -\log_D p(u) \geq 0$, $-\lceil -\log_D p(u) \rceil \leq \log_D p(u)$, and

$$D^{-\lceil -\log_D p(u) \rceil} \leq D^{\log_D p(u)}.$$

$$\text{Thus, } \sum_{u \in \mathcal{U}} D^{-\ell(u)} = \sum_{u \in \mathcal{U}} D^{-\lceil -\log_D p(u) \rceil} \leq \sum_{u \in \mathcal{U}} D^{\log_D p(u)} = \sum_{u \in \mathcal{U}} p(u) = 1.$$

- $H_D(\{p(u)\}) \leq \bar{\ell}_{SF} < 1 + H_D(\{p(u)\})$

Proof. $1 + -\log_D p(u) \geq \lceil -\log_D p(u) \rceil \geq -\log_D p(u)$

$$\text{So, } 1 + E[-\log_D p(u)] > E[\lceil -\log_D p(u) \rceil] \geq E[-\log_D p(u)].$$

$$\text{Hence, } 1 + H_D(\{p(u)\}) > \bar{\ell}_{SF} \geq H_D(\{p(u)\}).$$

- If $-\log_D p(u)$ is an integer (**D-adic**) for all $u \in \mathcal{U}$, then $\bar{\ell}_{SF} = H_D(\{p(u)\})$.

- Ex. for $D = 2$, $\{p(u)\} = \left\{ \frac{1}{2}, \frac{1}{4}, \frac{1}{8}, \dots, \frac{1}{2^n}, \frac{1}{2^n} \right\}$.

- Ex. for general D , $\{p(u)\} = \left\{ \underbrace{\frac{1}{D}, \dots, \frac{1}{D}}_{D-1 \text{ times}}, \underbrace{\frac{1}{D^2}, \dots, \frac{1}{D^2}}_{D-1 \text{ times}}, \dots, \underbrace{\frac{1}{D^n}, \dots, \frac{1}{D^n}}_{D-1 \text{ times}}, \frac{1}{D^n} \right\}$.

$$\text{Proof } (D-1) \frac{\frac{1}{D} - \frac{1}{D^{n+1}}}{1 - \frac{1}{D}} + \frac{1}{D^n} = \left(1 - \frac{1}{D^n}\right) + \frac{1}{D^n} = 1.$$

- $H_D(\{p(u)\}) \leq \bar{\ell}_{opt} \leq \bar{\ell}_{SF} < 1 + H_D(\{p(u)\})$.

- Block-to-variable length codes

- Instead of $\mathbf{f}: \mathcal{U} \rightarrow \{0, \dots, D-1\}^*$, use $\mathbf{f}: \mathcal{U}^L \rightarrow \{0, \dots, D-1\}^*$ with corresponding $\{p(\underline{u}), \underline{u} \in \mathcal{U}^L\}$.

- Super-letters.

- Super letter source is still stationary.
- Entropy rate per letter is the same as that of original source (H).

Proof For original source, entropy rate per letter $H = \lim_{L \rightarrow \infty} H_L = \lim_{L \rightarrow \infty} \frac{H(U_1^L)}{L}$.

For the new one, entropy rate per super letter =

$$H_{\text{super}} = \lim_{n \rightarrow \infty} \frac{H(U_1^n)}{n} = \lim_{n \rightarrow \infty} \frac{H(U_1^{Ln})}{n}. \text{ Thus, entropy rate per letter}$$

$$= \lim_{n \rightarrow \infty} \frac{H(U_1^{Ln})}{Ln}. \text{ Note that sequence } \frac{H(U_1^{Ln})}{Ln} \text{ is a subsequence of } \frac{H(U_1^n)}{n}.$$

Because the sequence $\frac{H(U_1^n)}{n}$ converges, the subsequence converge to the same limit.

- $2-L$ delay and extra complexity.
- As $L \rightarrow \infty$, $\bar{\ell} \rightarrow H$, the source's entropy rate.

$$\text{Proof } \bar{\ell} = \frac{E[[-\log p(\underline{U})]]}{L}$$

$$< \frac{E[1 - \log p(\underline{U})]}{L} = \frac{1 + E[-\log p(\underline{U})]}{L}$$

$$= \frac{1 + H(U_1^L)}{L} = \frac{1}{L} + H_L$$

Note: $E[-\log p(\underline{U})] = H(U_1^L)$ because the source is stationary.

$\lim_{L \rightarrow \infty} \bar{\ell} \leq H_L = H$. And we already know that treating a super letter as normal

letter, $H(U_1^L) \leq \bar{\ell}_{\text{super}} < H(U_1^L) + 1$. So, $\frac{H(U_1^L)}{L} \leq \frac{\bar{\ell}_{\text{super}}}{L} < \frac{H(U_1^L)}{L} + \frac{1}{L}$, and

$$\text{thus } \lim_{L \rightarrow \infty} \frac{\bar{\ell}_{\text{super}}}{L} = \lim_{L \rightarrow \infty} \frac{H(U_1^L)}{L} = H.$$

- Huffman code, $D = 2$.

Given $\{p_j, 0 \leq j \leq M-1\}$. Want to build a minimum $\bar{\ell}$ binary prefix code by assigning

lengths $\{\ell_j, 0 \leq j \leq M-1\}$ that minimize $\bar{\ell} = \sum_{j=0}^{M-1} \ell_j p_j$.

$M = |\mathbf{u}|$. Assume $p_0 \geq p_1 \geq \dots \geq p_{M-1}$.

Assume that we have optimal length described by $\{\ell_j, 0 \leq j \leq M-1\}$.

Let ℓ_{\max} be the longest of the optimum ℓ_j 's.

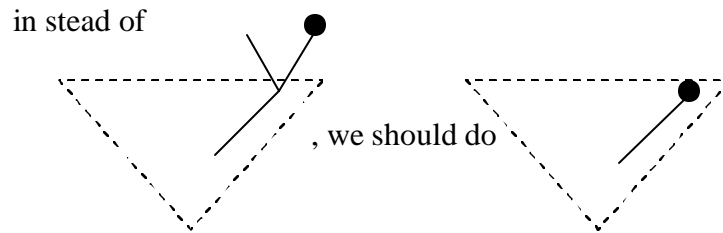
Least likely source symbol has $\ell = \ell_{\max}$. (Morse. If not, switch its assignment with the one that has ℓ_{\max} will give lower $\bar{\ell}$.)

Next-to-least likely source symbol should have $\ell = \ell_{\max}$ also.

Suppose not. Assume the next-to-least likely source symbol has $\ell < \ell_{\max}$.

Then, note that no other symbols can have $\ell = \ell_{\max}$. It has larger probability than the next-to-least likely symbol; so, it should not be assigned larger ℓ .

This means the least likely source symbol is the only one on the level ℓ . This is not optimal because



Without loss of generality, let's have code strings for these two letters identical through level $\ell_{\max}-1$. (Then, one of them ends with 0, the other with 1.) This means they are assigned a common ancestor on level $\ell_{\max}-1$.

Then, define new alphabet set with $|\mathcal{U}'| = M-1$.

$$\begin{aligned}
 p'_i &= p_i \text{ for } 0 \leq i \leq M-3. \quad p'_{M-2} = p_{M-2} + p_{M-1}. \\
 \bar{\ell} &= \sum_{k=0}^{M-1} \ell_k p_k = \sum_{k=0}^{M-3} \ell_k p_k + \ell_{\max} (p_{M-2} + p_{M-1}) \\
 &= \sum_{k=0}^{M-3} \ell'_k p'_k + \ell_{\max} (p'_{M-2}) = \sum_{k=0}^{M-3} \ell'_k p'_k + (\ell'_{M-2} + 1)(p'_{M-2}) \\
 &= \sum_{k=0}^{M-2} \ell'_k p'_k + p'_{M-2}
 \end{aligned}$$

Because p'_{M-2} is constant, we then want to minimize $\bar{\ell}' = \sum_{k=0}^{M-2} \ell'_k p'_k$.

This can be accomplished by recursively applying the above argument.

- D -ary Huffman code, $D \geq 2$.
 - Full D -ary tree: one with D branches out of every internal node.
 - Full tree has $D + k(D-1)$ terminal nodes for some non-negative integer k . So, if $|\mathcal{U}|$ is not in that form, then can't have full tree.
 - First full-fan has D terminal nodes. Growing from this fan, adding one full-fan takes one terminal node off; so, net increase = $D-1$.
 - Optimal code should have full tree except one of the top.
- If any one at lower level is not full, then can move one from the top down and reduce $\bar{\ell}$.

- If there exists k such that $|\mathcal{U}| = D + k(D - 1)$, then, hang D least likely off common ancestor and proceed iteratively as in binary.

If no such k exists, there is a non-full fan of the least likely letters' terminals on ℓ_{\max} . Let $N =$ size of this fan. Then, $2 \leq N \leq D$.

$|\mathcal{U}| = D + k(D - 1) + N - 1$ because we add one fan of N terminal nodes to a full tree. This adds N terminal nodes but takes out 1 terminal node.

$$\begin{aligned} |\mathcal{U}| &= D + k(D - 1) + N - 1 = (k + 1)(D - 1) + 1 + N - 1 \\ &= (k + 1)(D - 1) + N \\ &= (k + 1)(D - 1) + 2 + (N - 2) \end{aligned}$$

Note that $0 \leq N - 2 \leq D - 2$. Therefore, $(N - 2) \bmod (D - 1) = N - 2$.

So, $(|\mathcal{U}| - 2) \bmod (D - 1) = N - 2$. So, $N = 2 + (|\mathcal{U}| - 2) \bmod (D - 1)$.

Grab $N = 2 + (|\mathcal{U}| - 2) \bmod (D - 1)$ least likely at first, then always D at a time.

- For $D = 3$, $N = \begin{cases} 2, & \text{even } |\mathcal{U}| \\ 3, & \text{odd } |\mathcal{U}| \end{cases}$
- Universal lossless coding
 - Borisfitingof

Discrete memoryless stationary source with alphabet \mathcal{U} but unknown distribution $\{p(u), u \in \mathcal{U}\}$.

Encoding: Gather n -block, $\underline{u} = (u_1, \dots, u_n)$. Compute empirical distribution $\tilde{p}(u) = \frac{n(u)}{n}$. $n(u) = |\{k : u_k = u\}|$.

Also need to send $\{\tilde{p}(u)\}$. Note that $0 \leq n(u) \leq n$ for every $u \in \mathcal{U}$. So, $\{\tilde{p}(u)\}$ cannot assume more than $(n + 1)^{|\mathcal{U}|}$ values. So, take no more than $\lceil \log_2 (n + 1)^{|\mathcal{U}|} \rceil$ binary digits to specify $\{\tilde{p}(u), u \in \mathcal{U}\}$.

Thus, per source letter, use fewer than $\frac{\lceil \log_2 (n + 1)^{|\mathcal{U}|} \rceil}{n}$. As $n \rightarrow \infty$, this $\rightarrow 0$.
 - Lempel-Ziv (LZ codes)
 - Arithmetic (Pasco, Rissanen, Langdon)