

1. Coin Flips

$$H(p) = -p \log p - q \log q \quad p = 1-q$$

Suppose a coin with probability of head p is flipped until the first head occurs. X denote the number of flips required

(a) Find the entropy $H(X)$

$$P(X=1) = p \quad P(X=2) = pq \quad P(X=3) = pq^2 \\ \dots \quad P(X=n) = pq^{n-1}$$

$$H(X) = - \sum_{n=1}^{\infty} pq^{n-1} \log (pq^{n-1})$$

$$= -p \log p - pq \log pq - pq^2 \log pq^2 - \dots - pq^{n-1} \log pq^{n-1} - \dots$$

$$= -p \log p - pq [\log p + \log q] - pq^2 [\log p + 2 \log q] -$$

$$pq^3 [\log p + 3 \log q] - \dots - pq^{n-1} [\log p + (n-1) \log q] - \dots$$

$$= -p \log p \sum_{n=0}^{\infty} q^n - p \log q \sum_{n=1}^{\infty} nq^n$$

consider $\sum_{n=0}^{\infty} r^n = \frac{1}{1-r}$ $\sum_{n=1}^{\infty} nq^n = \frac{rq}{(1-q)^2}$

$$\therefore H(X) = -p \log p \frac{1}{1-q} - p \log q \frac{q}{(1-q)^2}$$

$$\begin{aligned}\therefore p+q=1 \quad \therefore 1-q=p \\ H(X) &= -p \log p \cdot \frac{1}{p} - p \log q \cdot \frac{q}{p^2} \\ &= \frac{1}{p} [-p \log p - q \log q] \\ &= \frac{1}{p} H(p) \text{ bits} \\ \therefore H(X) &= \frac{1}{p} H(p) \text{ bits}\end{aligned}$$

(b) what is the value for a fair coin?

$$\begin{aligned}\text{For a fair coin . } p=q=\frac{1}{2} \\ H(X) &= 2 \times H\left(\frac{1}{2}, \frac{1}{2}\right) = 2 \times 1 = 2 \\ H(X) &= 2\end{aligned}$$

(c) A random variable X is drawn according to the distribution with $p = \frac{1}{2}$

$$\rightarrow H(X) = 2$$

$$P(X=1) = \frac{1}{2} \quad P(X=2) = \frac{1}{4} \quad \dots \quad P(X=n) = \frac{1}{2^n}$$

One possible "efficient" sequence yes-no question is :

- (1) Is $X=1$? if not ↘
- (2) Is $X=2$? if not ↘
- (3) Is $X=3$? ↘
- :

So the pattern is if $X=r$, we need to ask r questions

And the expected number of questions should be

$$\frac{1}{2} \cdot 1 + \frac{1}{4} \cdot 2 + \frac{1}{8} \cdot 3 + \dots + \frac{1}{2^n} \cdot n$$

$$= \sum_{n=1}^{\infty} n \cdot \frac{1}{2^n} = 2$$

$\therefore H(X)$ is equal to the expected number of questions required to determine X

2. Minimum Entropy

$$H(p_1, p_2, \dots, p_n) = H(\vec{p})$$

According to the property of entropy, $H(p)$ is a concave function of the distribution.

For $H(\vec{p})$, it achieves its minimum value as $p_i = 1$

for some i and $p_j = 0$ for $j \neq i$

$$\text{e.g. } p = (1, 0, \dots, 0) \quad p = (0, 1, 0, \dots, 0)$$

$$\dots \quad p = (0, 0, \dots, 0, 1)$$

Since $p=0$ or 1 , the variable is not random, and there is no uncertainty. So $H(\vec{p}) = 0$

These points correspond to corners on the simplex

3. Entropy of function of a random variable

$$\begin{aligned}(a) H(x, g(x)) &= H(x) + H(g(x)|x) \\&= -\sum_x p(x) \log p(x) - \sum_x \sum_{g(x)} p(x, g(x)) \log p(g(x)|x) \\&= -\sum_x p(x) \sum_{g(x)} p(g(x)|x) \log p(x) - \sum_x \sum_{g(x)} p(x, g(x)) \log p(g(x)|x) \\&= -\sum_x \sum_{g(x)} p(x, g(x)) \left[\log p(x) + \log p(g(x)|x) \right] \\&= -\sum_x \sum_{g(x)} p(x, g(x)) \log \left(p(x) \cdot p(g(x)|x) \right) \\&= -\sum_x \sum_{g(x)} p(x, g(x)) \log p(x, g(x)) = H(x, g(x))\end{aligned}$$

(b)

$$H(x) + H(g(x)|x)$$

since $g(x)$ is a function of x , when x is fixed
 $g(x)$ is fixed \therefore there is no uncertainty in $p(g(x)|x)$
 $\therefore H(x) + H(g(x)|x) = H(x) + 0 = H(x)$

$$(c) H(x, g(x)) = H(g(x)) + H(x|g(x))$$

(can be proved by chain rule, similar to (a))

$$(d) H(g(x)) + H(x|g(x))$$

We know that $H(x|g(x)) \geq 0$

If $x \rightarrow g(x)$ is one-to-one mapping, we know x
if we already know $g(x)$. If it is not one-to-one,
there is uncertainty in $P(x|g(x))$

$$\therefore H(x|g(x)) \geq 0$$

$$H(g(x)) + H(x|g(x)) \geq H(g(x))$$

$$\therefore H(x, g(x)) = H(x) \geq H(g(x))$$

4. Computing Relative Entropy for 2D p and q

		$p(x,y)$	
		0	1
$x \setminus y$	0	$\frac{1}{6}$	$\frac{7}{12}$
	1	$\frac{1}{6}$	$\frac{1}{12}$

		$q(x,y)$	
		0	1
$x \setminus y$	0	$\frac{1}{4}$	$\frac{1}{2}$
	1	$\frac{1}{12}$	$\frac{1}{6}$

Find $D(p||q)$

$$D(p||q) = \sum_{x,y} p(x,y) \log \frac{p(x,y)}{q(x,y)}$$

$$= -\frac{1}{6} \log \frac{\frac{1}{6}}{\frac{1}{4}} + \frac{7}{12} \log \frac{\frac{7}{12}}{\frac{1}{2}} + \frac{1}{6} \log \frac{\frac{1}{6}}{\frac{1}{12}}$$

$$+ \frac{1}{12} \log \frac{\frac{1}{12}}{\frac{1}{6}}$$

$$= \frac{1}{6} \log \frac{2}{3} + \frac{1}{12} \log \frac{7}{6} + \frac{1}{6} \log 2 + \frac{1}{12} \log \frac{1}{2}$$

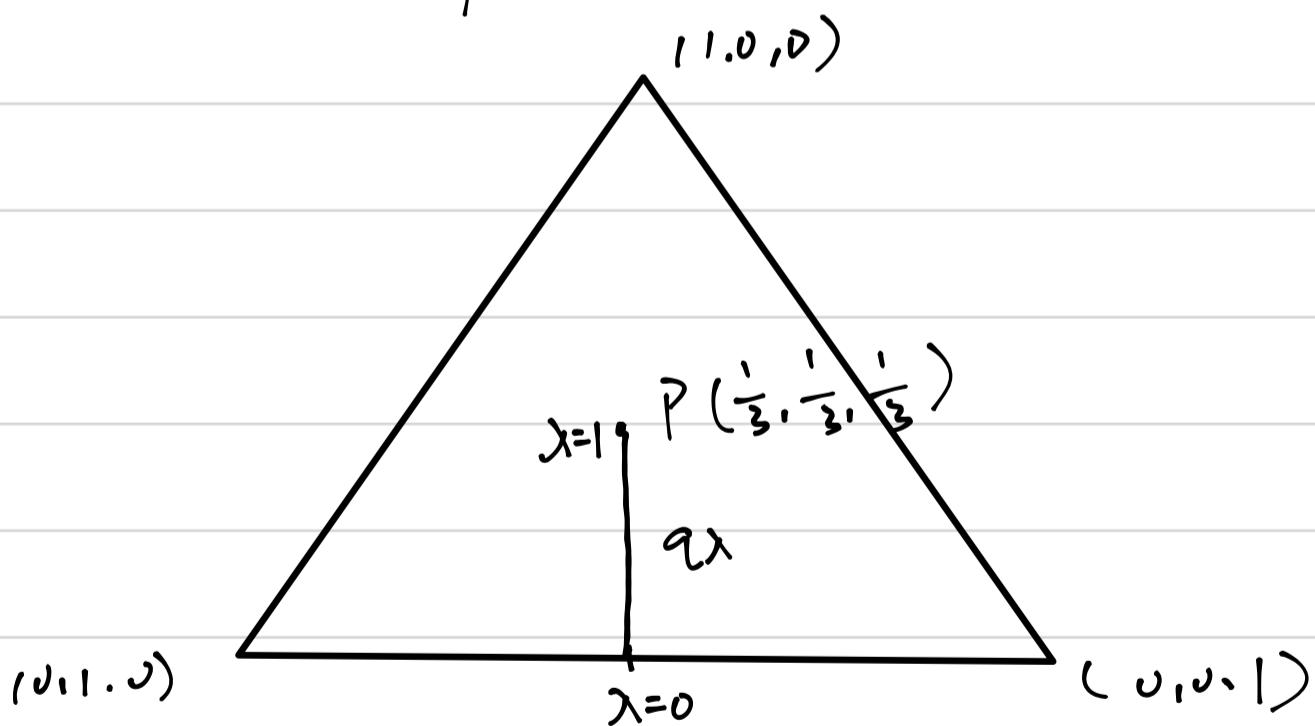
$$= 0.1156 \text{ bit} \quad \therefore D(p||q) = 0.1156 \text{ bit}$$

S. Computing Relative Entropy for p and q on a line in 3D simplex.

PMF of p is point $(\frac{1}{3}, \frac{1}{3}, \frac{1}{3})$

PMF of q_λ is point $(\frac{\lambda}{3}, \frac{3-\lambda}{6}, \frac{3-\lambda}{6})$

(a) draw a triangle representing the simplex, show the point p and the line segment that shows the trajectory of q as λ varies from 0 and 1

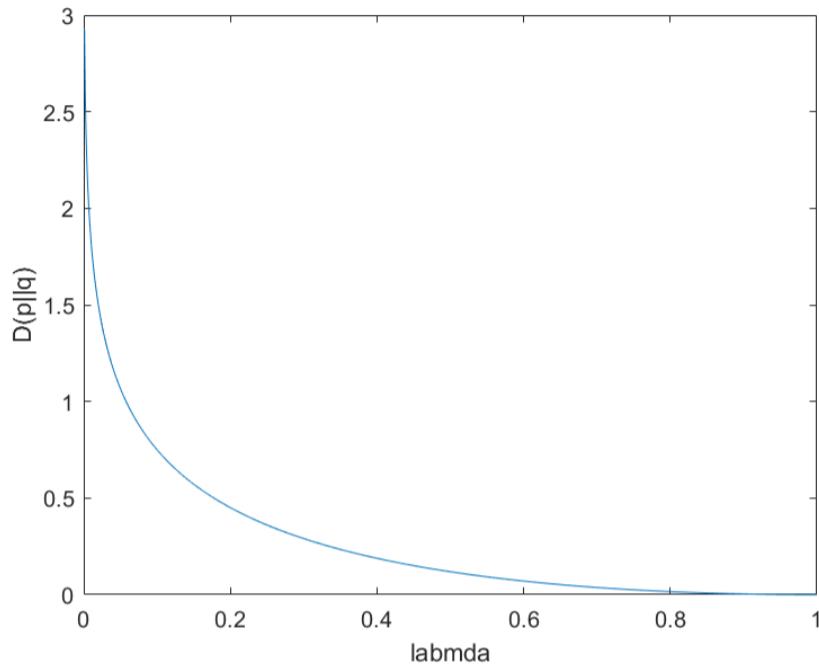


when $\lambda = 0$ $q = (0, \frac{1}{2}, \frac{1}{2})$ when $\lambda = 1$ $q = (\frac{1}{3}, \frac{1}{3}, \frac{1}{3})$

as λ varies from 0 to 1, q varies from the middle of the side to p

(b) Find $D(p||q_\lambda)$ as a function of λ as λ varies between 0 and 1 use MATLAB

$$D(p||q_\lambda)$$



when $\lambda = 1$ p and q represent the same distribution,

$$\therefore D(p||q) = 0$$

when $\lambda = 0$ $q = (0, \frac{1}{2}, \frac{1}{2})$

$$\begin{aligned} D(p||q) &= \sum_x p(x) \log \frac{p(x)}{q(x)} \\ &= \frac{1}{3} \log \frac{\frac{1}{3}}{0} + \frac{1}{3} \log \frac{\frac{1}{3}}{\frac{1}{2}} + \frac{1}{3} \log \frac{\frac{1}{3}}{\frac{1}{2}} \end{aligned}$$

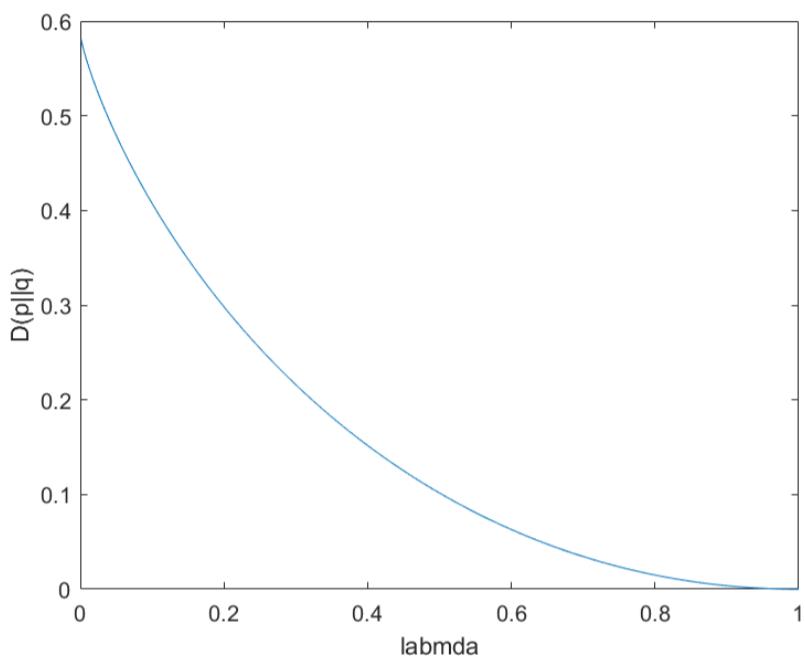
As defined in the book, we consider $p \log \frac{p}{q} = \infty$

$$\begin{aligned} \therefore D(p||q) &= +\infty + \frac{1}{3} \log \frac{2}{3} + \frac{1}{3} \log \frac{2}{3} \\ &= +\infty \end{aligned}$$

\therefore when $\lambda = 0$ $D(p||q)$ reaches positive definite

(c) Find $D(q_\lambda || p)$ as a function of λ as λ varies between 0 and 1 use MATLAB

$$D(q_\lambda || p)$$



Again when $\lambda=1$, p, q represent the same distribution

$$\therefore D(q||p)=0$$

when $\lambda=0$

$$D(q||p) = 0 \log \frac{0}{\frac{1}{3}} + \frac{1}{2} \log \frac{\frac{1}{2}}{\frac{1}{3}} + \frac{1}{2} \log \frac{\frac{1}{2}}{\frac{1}{3}}$$

as defined in the book, we consider $0 \log \frac{0}{q} = 0$

$$\therefore D(q||p) = (\frac{1}{2} \log \frac{3}{2}) \times 2 = \log \frac{3}{2}$$

(As $\log_2 0$ is not defined, we cannot get the point when $\lambda=0$ in MATLAB)

(d) Discuss the differences between $D(p||q_\lambda)$ and $D(q||p)$

$D(p||q)$ is the penalty when p is the true distribution, q is the wrong distribution

On the contrary, $D(q||p)$ is the penalty when q is the true distribution, p is the wrong distribution

We note that the "distance" between p and q is not symmetric as $D(p||q) \neq D(q||p)$

when $\lambda = 0$

$$p = \left(\frac{1}{3}, \frac{1}{3}, \frac{1}{3} \right) \quad q = (0, \frac{1}{2}, \frac{1}{2})$$

For $D(p||q)$,

The $p = (\frac{1}{3}, \frac{1}{3}, \frac{1}{3})$ is true distribution, but we wrongly set one variable to 0 possibility, this leads to the infinite penalty

For $D(q||p)$

The $q = (0, \frac{1}{2}, \frac{1}{2})$ is true distribution, but we wrongly choose $(\frac{1}{3}, \frac{1}{3}, \frac{1}{3})$. the penalty for this situation is smaller than the previous situation.

6. Mutual Information?

Can the relative entropy computed in problem 4 be expressed as a mutual information?

Yes.

The definition of $I(X;Y)$ is

$$I(X;Y) = \sum_{x,y} p(x,y) \log \frac{p(x,y)}{p(x)p(y)}$$

$$\text{For problem 4 } D(p||q) = \sum_{x,y} p(x,y) \log \frac{p(x,y)}{q(x,y)}$$

		0	1
0	0	$\frac{1}{6}$	$\frac{7}{12}$
	1	$\frac{1}{6}$	$\frac{1}{12}$

$$\therefore p(x=0) = \frac{3}{4} \quad p(x=1) = \frac{1}{4}$$

$$p(Y=0) = \frac{1}{3} \quad p(Y=1) = \frac{2}{3}$$

		0	1
0	0	$\frac{1}{4}$	$\frac{1}{2}$
	1	$\frac{1}{2}$	$\frac{1}{6}$

$$q(0,0) = p(x=0) \cdot p(Y=0) = \frac{1}{4}$$

$$q(0,1) = p(x=0) \cdot p(Y=1) = \frac{1}{2}$$

$$q(1,0) = p(x=1) \cdot p(Y=0) = \frac{1}{12}$$

$$q(1,1) = p(x=1) \cdot p(Y=1) = \frac{1}{6}$$

$$\therefore q(x,y) = p(x) \cdot p(y)$$

$$\therefore D(p||q) = \sum_{x,y} p(x,y) \log \frac{p(x,y)}{q(x,y)}$$

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

∴ It could be expressed as a mutual information

7. Example of joint entropy

		Y\X	
		0	1
X	0	$\frac{1}{3}$	$\frac{1}{3}$
	1	0	$\frac{1}{3}$

$$(a) P(X=0) = \frac{2}{3} \quad P(X=1) = \frac{1}{3}$$

$$H(X) = -\frac{2}{3} \log \frac{2}{3} - \frac{1}{3} \log \frac{1}{3} = 0.918 \text{ bit}$$

$$P(Y=0) = \frac{1}{3} \quad P(Y=1) = \frac{2}{3}$$

$$\therefore H(Y) = -\frac{1}{3} \log \frac{1}{3} - \frac{2}{3} \log \frac{2}{3} = 0.918 \text{ bit}$$

$$H(X) = H(Y) = 0.918 \text{ bit}$$

$$(b) H(X|Y) = \sum_y P(y) H(X|Y=y)$$

$$= \frac{1}{3} H(X|Y=0) + \frac{2}{3} H(X|Y=1)$$

$$= \frac{1}{3} H(1,0) + \frac{2}{3} H(\frac{1}{2}, \frac{1}{2}) = \frac{2}{3} = 0.667 \text{ bit}$$

$$H(Y|X) = \sum_x P(x) H(Y|X=x)$$

$$= \frac{2}{3} H(Y|X=0) + \frac{1}{3} H(Y|X=1)$$

$$= \frac{2}{3} H(\frac{1}{2}, \frac{1}{2}) + \frac{1}{3} H(0,1) = \frac{2}{3} = 0.667 \text{ bit}$$

$$(c) H(X,Y) = - \sum_{x,y} P(x,y) \log P(x,y)$$

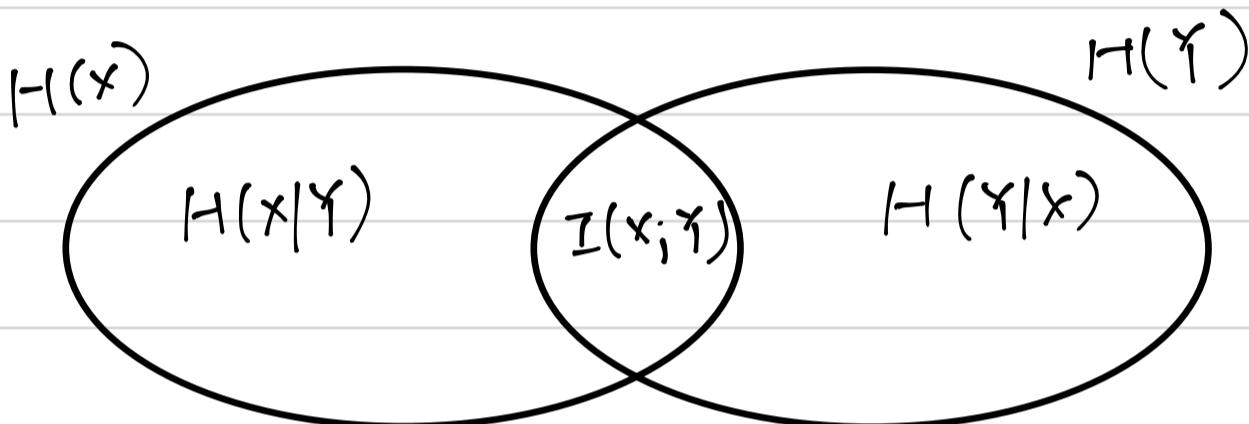
$$= -3 \times \frac{1}{3} \log \frac{1}{3} = \log 3 = 1.585 \text{ bits}$$

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

$$= 0.918 \text{ bit} - 0.667 \text{ bit} = 0.251 \text{ bit}$$

$$(e) I(X;Y) = H(Y) - H(Y|X) = 0.251 \text{ bit}$$

(f) Draw a Venn diagram for the quantities above



$$I(X;Y) = H(X) + H(Y) - H(X,Y)$$

8. Mutual Information and the Weather

Let random variable W denote the weather in El Nino University

according to the question

$$P(W=\text{rain}) = 0.9 \quad P(W=\text{sun}) = 0.1$$

(a) Stormy is a very lazy forecaster, it achieves his 90% correct forecast by simply always predict rain

\therefore Let random variable P_S denote the predicted weather for stormy forecaster

$$\therefore P(P_S=\text{sun}) = 0 \quad P(P_S=\text{rain}) = 1$$

Therefore, we can derive the joint distribution

$P_S \backslash W$	sun	rain
sun	0	0
rain	0.1	0.9

$$I(P_S; W) = H(W) - H(W|P_S)$$

$$H(W) = H(0.1, 0.9)$$

$$H(W|P_S) = \sum_{P_S} P(P_S) H(W|P_S = P_S)$$

$$= 0 + 1 \cdot H(0.1, 0.9) = H(0.1, 0.9)$$

$$\therefore I(P_S; W) = H(W) - H(W|P_S) = 0 \text{ bit}$$

(b) Wendy is a hard-working forecaster.

Let random variable P_W denote the predicted weather for Wendy forecaster

$$\text{According to the (b)} \quad P(P_W = \text{sun}) = 0.25$$

$$P(P_W = \text{rain}) = 0.75$$

Whenever she predicts rain she is correct

$$\therefore P(W = \text{rain} | P_W = \text{rain}) = 1$$

$$\begin{aligned} P(W = \text{rain}, P_W = \text{rain}) &= P(W = \text{rain} | P_W = \text{rain}) \cdot P(P_W = \text{rain}) \\ &= 1 \cdot 0.75 = 0.75 \end{aligned}$$

Whenever she predicts sun she is correct 40% of the time

$$\therefore P(W=\text{sun} | P_W=\text{sun}) = 0.4$$

$$P(W=\text{sun}, P_W=\text{sun}) = P(W=\text{sun} | P_W=\text{sun}) \cdot P(P_W=\text{sun}) \\ = 0.4 \cdot 0.25 = 0.1$$

\therefore we can draw the joint distribution here

$P_W \backslash W$	Sun	Rain
Sun	0.1	0.15
Rain	0	0.75

$$I(P_W; W) = H(W) - H(W|P_W)$$

$$H(W) = -0.1 \log 0.1 - 0.9 \log 0.9 = 0.469 \text{ bit}$$

$$H(W|P_W) = \sum_{P_W} P(P_W) H(W|P_W=P_W) \\ = 0.25 \cdot H(0.6, 0.4) + 0.75 H(0.1)$$

$$\textcircled{1} \quad H(0.6, 0.4) = -0.6 \log 0.6 - 0.4 \log 0.4 \\ = 0.971$$

$$\textcircled{2} \quad H(0.1) = 0$$

$$\therefore H(W|P_W) = 0.2427 \text{ bit}$$

$$\therefore I(P_W; W) = 0.469 - 0.2427 = 0.2263 \text{ bit}$$

$$(c) \quad I(P_S; W) = 0 \text{ bit}, \quad I(P_W; W) = 0.2263 \text{ bit}$$

\therefore wendy forecaster provides the most information about the weather

(d) Suppose you want to plant your tulip bulbs on a day when you knew it was going to rain. which forecaster would be helpful in achieving this?

I will select Wendy forecast, as the Storm forecaster always predict rain. it provides no information to me, it is helpless.

For forecaster Wendy, when she predicts rain. it will rain, this will help me in this situation

9. Concavity of entropy

(a) Show that $\log x$ is concave in x for positive x

$$\log x = \ln x / \ln 2$$

$\ln 2$ is a constant, if the base for \log is b ,

$$\log_b x = \ln x / \ln b. \text{ it should be a constant}$$

$$\frac{d}{dx} \log x = \frac{d}{dx} \frac{1}{\ln b} \ln x = \frac{1}{\ln b} \frac{1}{x}$$

$$\frac{d^2}{dx^2} \log x = \frac{1}{\ln b} \cdot \left(-\frac{1}{x^2}\right) < 0 \text{ for } x > 0$$

∴ for positive x , $\log x$ is concave

(b) show that $x \log x$ is convex in x for positive x .

Suppose the base of \log is b

$$\therefore \log_b x = \ln x / \ln b$$

$$\therefore \frac{d}{dx} x \log x = \frac{d}{dx} \times \frac{\ln x}{\ln b} = \frac{1}{\ln b} [\ln x + 1]$$

$$\frac{d^2}{dx^2} x \log x = \frac{1}{\ln b} \cdot \frac{1}{x} < 0 \quad \text{for positive } x$$

$\therefore x \log x$ is convex

(c) Show that $H(p) = -p \log p - (1-p) \log(1-p)$ is concave in p for $0 \leq p \leq 1$

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

$$= -\frac{1}{\ln b} [p \ln p - (1-p) \ln(1-p)]$$

$$\frac{d H(p)}{dp} = -\frac{1}{\ln b} [\ln p + 1 - \ln(1-p) - 1]$$

$$\frac{d^2}{dp^2} H(p) = -\frac{1}{\ln b} \left[\frac{1}{p} + \frac{1}{1-p} \right]$$

$$\therefore 0 \leq p \leq 1 \quad \therefore \frac{1}{p}, \frac{1}{1-p} \geq 0 \quad \therefore \frac{d^2}{dp^2} H(p) \leq 0$$

$\therefore H(p)$ is concave in p

for $0 \leq p \leq 1$

10. Maximum entropy

what is the maximum value of $H(p_1, \dots, p_n) = -I(p)$

let u be the uniform distribution on \mathcal{X}

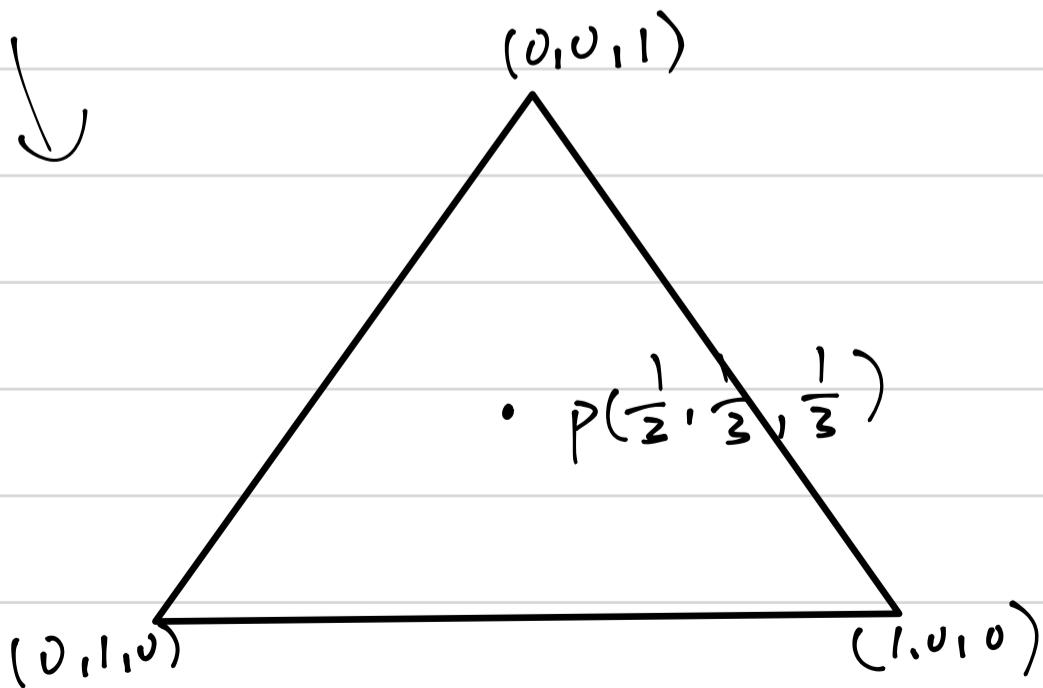
$$E_p[-\log u(x)] = E_p[-\log \frac{1}{n}] = \log n$$

$$\begin{aligned} D(p||u) &= E_p[\log \frac{p(x)}{u(x)}] \\ &= E_p[\log p(x) - \log u(x)] \\ &= E_p[-\log u(x)] - E_p[-\log p(x)] \\ &= \log n - H(x) \geq 0 \\ \therefore H(x) &\leq \log n \end{aligned}$$

\therefore The maximum value of $H(\vec{p})$ is $\log n$, when p is uniform distribution

That is $p_1 = p_2 = \dots = p_n = \frac{1}{n}$

For $n=3$ $p_1 = p_2 = p_3 = \frac{1}{3}$



11. Drawing with and without replacement

which has higher entropy. drawing $k \geq 2$ balls from the urn with replacement or without replacement

Let X_i denote the situation for i th draw

$$\therefore P(X_i = \text{red}) = r/(r+w+b)$$

$$P(X_i = \text{white}) = w/(r+w+b)$$

$$P(X_i = \text{blue}) = b/(r+w+b)$$

The entropy we compare is $H(X_1, X_2, \dots, X_k)$

Using the chain rule:

$$H(X_1, X_2, \dots, X_k) = \sum_{i=1}^k H(X_i | X_{i-1}, X_{i-2}, \dots, X_1)$$

(1) with replacement:

$H(X_i | X_{i-1}, X_{i-2}, \dots, X_1) = H(X_i)$ (since with replacement, the probability for each draw is same)

$$\begin{aligned} H(X_i) &= -r/(r+w+b) \log(r/(r+w+b)) - w/(r+w+b) \\ &\quad (\log(w/(r+w+b))) - b/(r+w+b) \log(b/(r+w+b)) \\ &= \log(r+w+b) - \sqrt{(r+w+b)} \log r - w/(r+w+b) \log w - b/(r+w+b) \log b \end{aligned}$$

(2) without replacement:

The conditional entropy $H(X_i | X_{i-1}, \dots, X_1) \leq H(X_i)$

for each i (since each draw is not independent)

\therefore the entropy of drawing without replacement is smaller

12. Conditional Mutual Information

(a) Show that $X \rightarrow Y \rightarrow Z$ forms a Markov chain

$$I(X; Y|Z) \leq I(X; Y)$$

$$\begin{aligned} \text{Consider } I(X; Y, Z) &= I(X; Z) + I(X; Y|Z) \\ &= I(X; Y) + I(X; Z|Y) \end{aligned}$$

$$\therefore I(X; Z|Y) = 0$$

$$\therefore I(X; Z) + I(X; Y|Z) = I(X; Y)$$

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

$$\therefore I(X; Y) \geq 0$$

$$\therefore I(X; Y|Z) \leq I(X; Y)$$

(b) Not always true

As mentioned in the lecture

Consider X and Y , two independent binary random variables and set $Z = X \oplus Y$

$$I(X; Y) = 0 \quad I(X; Y|Z) = 1$$

\therefore There exists examples that conditions increase mutual information.

13. Find the gap

For $X \rightarrow Y \rightarrow Z \quad I(X; Z) \leq I(Y; Z)$, find the exact gap between these mutual information

$X \rightarrow Y \rightarrow Z$ also implies $Z \rightarrow Y \rightarrow X$

Consider $I(Z; X, Y) = I(Z; X) + I(Z; Y|X)$
 $= I(Z; Y) + I(Z; X|Y)$

$X \rightarrow Y \rightarrow Z$ is a markov chain

$$\therefore I(Z; X|Y) = 0$$

$$\therefore I(Z; Y) = I(Z; X) + I(Z; Y|X)$$

$$\therefore I(Y; Z) - I(X; Z) = I(Y; Z|X)$$

The gap is a conditional mutual information
 $I(Y; Z|X)$