

Information Theory

Problem Sheet 1

(Most questions are from Cover & Thomas, the corresponding question numbers (as in 1st ed.) are given in brackets at the start of the question)

Notation: \mathcal{X} , \mathbf{x} , \mathbf{X} are scalar, vector and matrix random variables respectively.

The following expressions may be useful: $\sum_{n=1}^{\infty} r^n = \frac{r}{1-r}$ $\sum_{n=1}^{\infty} nr^n = \frac{r}{(1-r)^2}$

1. [2.1] A fair coin is flipped until the first head occurs. Let \mathcal{X} denote the number of flips required.
 - (a) Find the entropy $H(\mathcal{X})$ in bits.
 - (b) A random variable \mathcal{X} is drawn according to this distribution. Find an “efficient” sequence of yes-no questions of the form “Is \mathcal{X} contained in the set S ?”. Compare $H(\mathcal{X})$ to the expected number of questions required to determine \mathcal{X} .

2. [~2.2] \mathcal{X} is a random variable taking integer values. What can you say about the relationship between $H(\mathcal{X})$ and $H(\mathcal{Y})$ if
 - (a) $\mathcal{Y} = \mathcal{X}^2$
 - (b) $\mathcal{Y} = \mathcal{X}^3$

3. [2.3] If \mathbf{p} is an n -dimensional probability vector, what is the maximum and the minimum value of $H(\mathbf{p})$. Find all vectors \mathbf{p} for which $H(\mathbf{p})$ achieves its maximum or minimum value.

4. We write $H(p)$ (with a scalar p) to denote the entropy of the Bernoulli random variable with probability mass vector $\mathbf{p} = [1-p \quad p]$. Prove the following properties of this function:
 - (a) $H'(p) = \log(1-p) - \log p$
 - (b) $H''(p) = \frac{-\log e}{p(1-p)}$
 - (c) $H(p) \geq 2 \min(p, 1-p)$
 - (d) $H(p) \geq 1 - 4(p - 1/2)^2$
 - (e) $H(p) \leq 1 - 2 \log e (p - 1/2)^2$

5. [2.5] Let \mathcal{X} be a discrete random variable and $g(\mathcal{X})$ a deterministic function of it. Show that $H(g(\mathcal{X})) \leq H(\mathcal{X})$ by justifying the following steps:

$$H(\mathcal{X}, g(\mathcal{X})) \stackrel{(a)}{=} H(\mathcal{X}) + H(g(\mathcal{X}) | \mathcal{X}) \stackrel{(b)}{=} H(\mathcal{X})$$

$$H(\mathcal{X}, g(\mathcal{X})) \stackrel{(c)}{=} H(g(\mathcal{X})) + H(\mathcal{X} | g(\mathcal{X})) \stackrel{(d)}{\geq} H(g(\mathcal{X}))$$

6. [2.6] Show that if $H(\mathcal{Y} | \mathcal{X}) = 0$, then \mathcal{Y} is a function of \mathcal{X} , that is for all x with $p(x) > 0$, there is only one possible value of y with $p(x,y) > 0$.

7. [~2.7] x_i is a sequence of i.i.d. Bernoulli random variables with $p(x_i=1) = p$ where p is unknown. We want to find a function f that converts n samples of x into a smaller number, K , of i.i.d. Bernoulli random variables, Z_i , with $p(Z_i=1)=1/2$. Thus $Z_{1:K}=f(x_{1:n})$ where K can depend on the values x_i .

(a) Show that the following mapping for $n=4$ satisfies the requirements and find the expected value of K , $E(K)$.

0000,1111 → ignore; 1010 → 0; 0101 → 1; 0001,0011,0111 → 00;
 0010,0110,1110 → 01; 0100,1100,1101 → 10; 1000,1001,1011 → 11

(b) Justify the steps in the following bound on $E(K)$

$$\begin{aligned} nH(p) &\stackrel{(a)}{=} H(x_{1:n}) \stackrel{(b)}{\geq} H(Z_{1:K}, K) \stackrel{(c)}{=} H(K) + H(Z_{1:K} | K) \\ &\stackrel{(d)}{=} H(K) + EK \stackrel{(e)}{\geq} EK \end{aligned}$$

8. [2.10] Give examples of joint random variables x , y and Z such that:

(a) $I(x; y | Z) < I(x; y)$

(b) $I(x; y | Z) > I(x; y)$

9. [2.12] We can define the “mutual information” between three variables as

$$I(x; y; z) = I(x; y) - I(x; y | z)$$

(a) Prove that

$$\begin{aligned} I(x; y; z) &= H(x, y, z) - H(x, y) - H(y, z) - H(z, x) \\ &\quad + H(x) + H(y) + H(z) \end{aligned}$$

(b) Give an example where $I(x; y; z)$ is negative. This lack of positivity means that it does not have the intuitive properties of an “information” measure which is why I put “mutual information” in quotes above.

10. [2.17] Show that $\log_e(x) \geq 1-x^{-1}$ for $x > 0$.

11. [~2.16] x and y are correlated binary random variables with $p(x=y=0)=0$ and all other joint probabilities equal to $1/3$. Calculate $H(x)$, $H(y)$, $H(x|y)$, $H(y|x)$, $H(x,y)$, $I(x,y)$.

12. [~2.22] If $x \rightarrow y \rightarrow z$ form a Markov chain, and for y , the alphabet size $|\mathcal{Y}| = k$, show that $I(x; z) \leq \log k$. What does this tell you if $k = 1$?

13. [2.29] Prove the following and find the conditions for equality:

(a) $H(x, y | z) \geq H(x | z)$

(b) $I(x, y; z) \geq I(x; z)$

(c) $H(x, y, z) - H(x, y) \leq H(x, z) - H(x)$

(d) $I(x, z | y) \geq I(z, y | x) - I(z, y) + I(x; z)$

Information Theory

Solution Sheet 1

1. (a) $X = n$ means that Tail occurs for the first $n - 1$ flips, while the last flip is Head. Thus, X has distribution $P(X = n) = 2^{-(n-1)} 2^{-1} = 2^{-n}$. Thus

$$\begin{aligned} H(X) &= \sum_{n=1}^{\infty} 2^{-n} \log 2^n = \sum_{n=1}^{\infty} 2^{-n} n \log 2 \\ &= \sum_{n=1}^{\infty} n 2^{-n} = \frac{1/2}{(1-1/2)^2} = 2 \end{aligned}$$

- (b) Ask if $x = 1, 2, 3, \dots$ in turn, i.e., ask the following questions:

Is $X = 1$?
 If not, is $X = 2$?
 If not, is $X = 3$?
 ...

Expected number of questions is $\sum_{n=1}^{\infty} n 2^{-n} = 2$.

2. $H(X, Y) = H(X) + H(Y|X) = H(Y) + H(X|Y)$, but $H(Y|X) = 0$ since Y is a function of X so $H(Y) = H(X) - H(X|Y) \leq H(X)$ with equality iff $H(X|Y) = 0$ which is true only if X is a function of Y , i.e. if Y is a one-to-one function of X for every value of x with $p(x) > 0$. Hence

(a) $H(Y) \leq H(X)$ because, for example $1^2 = (-1)^2$

(b) $H(Y) = H(X)$

3. Maximum is $\log n$ iff all elements of \mathbf{p} are equal. Minimum is 0 iff only one element of \mathbf{p} is non-zero; there are n possible elements that this could be.

4. (a) and (b) are straightforward calculus: easiest to convert logs to base e first. For the others, assume $1/2 < p < 1$ for convenience (other half follows by symmetry). Since $H''(p) < 0$, $H(p)$ is concave and so lies above the straight line $2 - 2p$ defined in (c).

At $p = 1/2$ the bound in (e) has the same value and first two derivatives as $H(p)$. For $1/2 < p < 1$ its second derivative is greater than $H''(p)$ and so the bound follows.

For (d) we consider $D(p) = H(p) - 1 + 4(p - 1/2)^2$. $D'(p) = 0$ is a quadratic in p and has only two solutions $p = 1/2 \pm \sqrt{(2 - \log e)/8} = 0.5 \pm 0.26$. Therefore $D'(p)$ increases from 0 at $p = 0.5$ to reach a maximum at $p = 0.76$ and decreases thereafter. This implies that $D'(p) = 0$ has only one solution for $p > 1/2$ and therefore that $D(p)$ has a single maximum. Since $D(1/2) = D(1) = 0$ we must have $D(p) > 0$ for $1/2 < p < 1$.

5. (a) chain rule, (b) $g(X|X)$ has only one possible value and hence zero entropy, (c) chain rule, (d) entropy is positive. We have equality at (d) iff $g(X)$ is a one-to-one function for every x with $p(x) > 0$.

6. $H(Y|X) = \sum_x p(x) H(Y|X=x)$

All terms are non-negative so the sum is zero only if all terms are zero. For any given term this is true either if $p(x) = 0$ or if $H(Y|X=x)$ is zero. The second case arises only if $H(Y|X=x)$ has only one value, i.e. Y is a function of X . The first case is why we needed the qualification about $p(x) > 0$ in answers 2 and 4 above.

7. (a) The probability of any given value of $X_{1:4}$ depends on the number of 1's and 0's. We create four subsets with equal probabilities to generate a pair of bits and two other subsets to generate one bit only. The expected number of bits generated is

$$EK = 8p(1-p)^3 + 10p^2(1-p)^2 + 8p^3(1-p)$$

- (b) (a) i.i.d entropies add, (b) functions reduce entropy, (c) chain rule, (d) Z_i are i.i.d. with entropy of 1 bit, (e) entropy is positive.

8. (a) This is true for any Markov chain $X \rightarrow Y \rightarrow Z$. One possibility is $X=Y=Z$ all fair Bernoulli variables.

- (b) An example of this was given in lectures. A slightly different example is if X and Y are fair binary variables and $Z = XY$. Knowing Z , entangles X and Y .

9. (a) $I(X, Y, Z) = \{H(X) - H(X|Y)\} - \{H(X|Z) - H(X|Y, Z)\} = H(X) - \{H(X, Y) - H(Y)\} - \{H(X, Z) - H(Z)\} + \{H(X, Y, Z) - H(Y, Z)\}$

- (b) Use the example from 8(b) above.

10. Define $f(x) = \ln(x) + x^{-1} - 1$. This is continuous and differentiable in $(0, \infty)$. Differentiate twice to show that the only extremum occurs at $x=1$ and that it is a minimum. Hence $f(x) \geq f(1) = 0$.
11. $H(x) = H(y) = 0.918$; $H(x|y) = H(y|x) = 0.667$; $H(x, y) = 1.58$; $I(x, y) = 0.252$.
12. The data processing inequality says that $I(x, z) \leq I(x, y) = H(y) - H(y|x) \leq H(y) \leq \log k$ where the last inequality is the uniform bound on entropy. If $k=1$ then $\log k = 0$ and so x and z must be independent.
13. (a) $H(x, y|z) = H(x|z) + H(y|x, z) \geq H(x|z)$ with equality if y is a function of x and z .
- (b) $I(x, y, z) = I(x, z) + I(y, z|x) \geq I(x, z)$ with equality if y and z are conditionally independent given x .
- (c) $H(x, y, z) - H(x, y) = H(z|x, y) = H(z|x) - I(y, z|x) \leq H(z|x) = H(x, z) - H(x)$ with equality if y and z are conditionally independent given x .
- (d) $I(x, y, z) = I(y, z) + I(x, z|y) = I(x, z) + I(y, z|x)$. Rearrange this to give the inequality which is in fact always an equality (trick question).