MA3466 – Information Theory

David Whyte dawhyte@tcd.ie

Updated April 3, 2012

Contents

1	Proba	bility & statistics	3
	1.1 R	andom variables	3
	1.2 E	xpected value	3
	1.3 M	10ments	4
	1.4 M	fultiple random variables	5
	1.5 B	ayes' rule	6
	1.6 M	Ionty Hall Problem	6
2	Quant	tities in information theory	7
	2.1 E	ntropy	7
	2.2 Jo	Dint entropy	9
	2.3 C	Conditional entropy	9
	2.4 C	hain rule for entropy	9
	2.5 K	ullback–Leibler divergence	10
	2.6 M	Iutual information	10
	2.7 R	elationship between all this stuff	11
	2.8 C	hain rule for entropy	11
	2.9 C	hain rule for mutual information	11
	2.10 C	hain rule for KL divergence	12
3	Jenser	ı's inequality etc.	13
	3.1 C	Cup-like & cap-like functions	13
	3.2 Je	ensen's inequality	13
	3.3 Ir	nformation inequality	14
	3.4 U	Ipper bound on information	15
	3.5 L	og-sum inequality	16
	3.6 Ir	nformation inequality	17
	3.7 C	up-likeness of KL divergence	17
	3.8 C	ap-likeness of entropy	18
	3.9 A	nother theorem	18
4	Marko	ov chains & sufficient statistics	19
	4.1 M	farkov chains	19

	4.2	Data processing inequality	19
	4.3	Sufficient statistics	20
	4.4	Fano inequality	20
5	Asy	mptotic equipartition principle & typical sets	22
	5.1	Convergence	22
	5.2	Law of large numbers	22
	5.3	Asymptotic equipartition principle	22
	5.4	Typical set	23
	5.5	Code word length	24
6	Enc	oding	25
	6.1	Definitions	25
	6.2	Kraft inequality	26
	6.3	Optimal codes	27
	6.4	Huffman coding	28

1 Probability & statistics

1.1 Random variables

When we say *X* is a random variable, we consider a set of outcomes (or "alphabet") $\mathscr{X} = \{x_1, x_2, ..., x_n\}$ where *n* is finite. It is not much more difficult to deal with the infinite discrete case – but the continuous case is more awkward in information theory. In this course we'll deal with the discrete case only.

For the set of outcomes $\mathcal X$, we have a map:

$$p_x \colon \mathscr{X} \to \mathbb{R}_+$$
$$x_i \mapsto p_X(x_i)$$

with

$$\sum_{x_i \in \mathscr{X}} p_X(x_i) = 1.$$

 $p_X(x)$ can also be denoted Pr(x).

For example, if *X* is the roll of a die, then:

$$\mathscr{X} = \{1, 2, 3, 4, 5, 6\}$$

 $p_X(1) = \ldots = p_X(6) = \frac{1}{6}.$

If *X* is the next letter in a text,

$$\mathscr{X} = \{a, \ldots, z\}.$$

Here *X* is not random, but it is unpredictable!

The estimator for the probability is:

$$p_X(x_i) = \lim_{N \to \infty} \frac{\# \text{ times } x_i \text{ happens}}{N}$$

Generally:

$$\Pr(E) = \sum_{x \in \mathscr{X} \mid E \text{ holds}} p_X(x)$$

For example,

$$Pr(3 \text{ or less}) = p_X(1) + p_X(2) + p_X(3) = \frac{1}{6} + \frac{1}{6} + \frac{1}{6} = \frac{1}{2}$$

1.2 Expected value

Say $g : \mathcal{X} \to \mathbb{R}$ (or, more generally, any field). The expected value of g is:

$$E_X g = \langle g(X) \rangle_{p_X} = \sum_{x \in \mathcal{X}} g(x) p_X(x)$$

For example, if the distribution of heights in cm were:

199	200	201
$\frac{1}{4}$	$\frac{1}{2}$	$\frac{1}{4}$
	$\frac{199}{\frac{1}{4}}$	$\begin{array}{ccc} 199 & 200 \\ \frac{1}{4} & \frac{1}{2} \end{array}$

To find the expected value of height, g is the trivial map and:

$$\langle X \rangle = \frac{1}{4} \cdot 199 + \frac{1}{2} \cdot 200 + \frac{1}{4} \cdot 201 = 200.$$

Or to consider a farmyard example,

Item	cow	sheep	turnip
р	$\frac{1}{2}$	$\frac{1}{4}$	$\frac{1}{4}$

Using the trivial map makes no sense since \mathscr{X} has no additive structure. However if we define a map $\notin : \mathscr{X} \to \mathbb{R}$, with $\notin(\text{cow}) = 700$, $\notin(\text{sheep}) = 100$ and $\notin(\text{turnip}) = 0.2$, for example, then:

$$\langle \mathfrak{E}(X) \rangle = \frac{1}{2} \cdot 700 + \frac{1}{4} \cdot 100 + \frac{1}{4} \cdot 0.2$$

= 375.05

In statistics, we have an estimator for $\langle g(x) \rangle$:

$$\langle g(x) \rangle = \lim_{N \to \infty} \frac{1}{N} \sum_{\text{trials}} g(x)$$

If \mathscr{X} does have a field structure, e.g. $\mathscr{X} \subset \mathbb{R}$, $\langle X \rangle$ is called the mean:

$$\langle X \rangle = \sum_{x \in \mathscr{X}} x p_X(x) = \overline{X}.$$

1.3 Moments

The n^{th} central moment is defined as:

$$\mu_n = \langle (X - \overline{X})^n \rangle.$$

Let's calculate μ_1 and μ_2 :

$$\mu_{1} = \langle X - \overline{X} \rangle$$

$$= \sum_{x \in \mathscr{X}} p_{X}(x)x - \sum_{x \in \mathscr{X}} p_{X}(x)\overline{X}$$

$$= \overline{X} - \overline{X} \sum_{x \in \mathscr{X}} p_{X}(x)$$

$$= 0.$$

$$\mu_{2} = \langle (X - \overline{X})^{2} \rangle$$

$$= \langle X^{2} \rangle - 2\overline{X} \langle X \rangle + \langle \overline{X}^{2} \rangle$$

$$= \langle X^{2} \rangle - \overline{X}^{2}.$$

 $\mu_2 = \sigma^2$, the variance. σ is the standard deviation.

 $\gamma_1 = \frac{\mu_3}{\sigma^3}$ is the skewness, and $\gamma_2 = \frac{\mu_4}{\sigma^4} - 3$ is the kurtosis.

In passing, the issue of estimators is non-trivial:

$$\overline{X} = \frac{1}{N} \sum_{N \text{ trials}} x$$
$$\sigma^2 = \frac{1}{N-1} \sum_{N \text{ trials}} (x - \overline{X})^2$$

Where the N-1 comes from is very subtle indeed! Notice that as $N \rightarrow \infty$ it becomes unimportant.

1.4 Multiple random variables

Say you have two random variables *X* and *Y*, with sets of outcomes \mathscr{X} and \mathscr{Y} . We can consider $\mathscr{X} \times \mathscr{Y}$ and it is often possible to define the random variable (*X*, *Y*) with probabilities $p_{X,Y}(x, y)$.

For example, if *X* is the first letter of a random word and *Y* is the last letter of the same word,

$$p_X(\mathbf{e}) \approx 0.0186$$

 $p_Y(\mathbf{e}) \approx 0.1916$

(for a certain corpus).

The probability that both the first and last letters are e is:

$$p_{(X,Y)}(e,e) \approx 0.0028$$

Note:

$$p_X(e) p_Y(e) = 0.0035 > p_{(X,Y)}(e).$$

So *X* and *Y* are not independent – a word starting with e is less likely to end in e than a randomly selected word, and vice versa.

 $p_{(X,Y)}(x, y)$ is called the *joint distribution*, and $p_X(x)$ and $p_Y(y)$ are called *marginal distributions*. $p_X(x)$ can be expressed as:

$$p_X(x) = \sum_{y \in \mathscr{Y}} p_{X,Y}(x,y).$$

There are also *conditional distributions* – for example, $p_{X|Y}(x|y)$ is the probability that X = x given that Y = y.

1.5 Bayes' rule

Bayes' rule relates joint, conditional and marginal probabilities. One statement of the rule is:

$$\underbrace{p_{X,Y}(x,y)}_{\text{joint}} = \underbrace{p_{X|Y}(x|y)}_{\text{conditional marginal}} \underbrace{p_{Y}(y)}_{\text{marginal}},$$

i.e. "The probability of getting X = x and Y = y is the probability of getting Y = y multiplied by the probability of getting X = x given that Y = y."

Switching $X \leftrightarrow Y$ gives: $p_{X,Y}(x, y) = p_{Y|X}(y|x)p_X(x)$, so we can combine these to give a more common statement of Bayes' rule.

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

Omitting subscripts (as above) is a common but potentially confusing abuse of notation.

Using the first/last letter example again, $p_{X,Y}(e, e) = 0.0028$ and $p_X(e) = 0.0186$. So:

$$p_{Y|X}(\mathbf{e}, \mathbf{e}) = \frac{p_{X,Y}(\mathbf{e}, \mathbf{e})}{p_X(\mathbf{e})} = 0.1522.$$

Compare this to $p_Y(e) = 0.1916$.

For another example, consider $\mathcal{Y} = \{a, b\}$ and $\mathcal{X} = \{1, 2, 3\}$, with probabilities:

We can calculate the marginal distributions for *X* and *Y* by simple addition:

If we wanted to find $p_{X|Y}(1|b)$, we use Bayes' rule:

$$p_{X|Y}(1|b) = \frac{p_{X,Y}(1,b)}{p_Y(b)} = \frac{\frac{1}{6}}{\frac{7}{12}} = \frac{2}{7}.$$

1.6 Monty Hall Problem

If *C* is the position of the car, and *R* is the door opened by Monty, we can write:

$$\mathscr{C} = \{1, 2, 3\}, \quad \mathscr{R} = \{2, 3\}$$

assuming that the contestant always chooses door 1 for simplicity.

Let's say R = 2 and we want $P_{C|R}(c|2)$. Use Bayes' rule:

$$p_{C|R}(c|r) = \frac{p_{R|C}(r|c)p_C(c)}{p_R(r)}$$

We know that $p_R(2) = \frac{1}{2}$ and $p_C(1) = \frac{1}{3}$. We also know the conditional probabilities $p_{R|C}(2|3) = 1$ and $p_{R|C}(2|1) = \frac{1}{2}$. So:

$$p_{C|R}(1|2) = \frac{p_{R|C}(2|1)p_C(1)}{p_R(2)} = \frac{\frac{1}{2} \cdot \frac{1}{3}}{\frac{1}{2}} = \frac{1}{2},$$
$$p_{C|R}(3|2) = \frac{p_{R|C}(2|3)p_C(3)}{p_R(2)} = \frac{1 \cdot \frac{1}{3}}{\frac{1}{2}} = \frac{2}{3}.$$

So the contestant is better off switching!

Also, $p_{C|R}(2|2) = 0$. Naturally.

2 Quantities in information theory

2.1 Entropy

Shannon's entropy is the central quantity in information theory. It is defined as:

$$H(X) = -\sum_{x \in \mathscr{X}} p_X(x) \log_2 p_X(x)$$
$$= -\langle \log_2 p_X(x) \rangle_p.$$

Note: "log" without a subscript is in this course taken to mean log₂.

This definition can be justified somewhat – no matter what type of thing *x* is, $p_X(x)$ is always a number and so $\langle p_X(x) \rangle$ is well-defined. In addition, the log simplifies the multiplicative structure of probabilities into an additive one.

For example, for this uniform distribution:

_

a
 b
 c
 d
 e
 f
 g
 h

$$\frac{1}{8}$$
 $\frac{1}{8}$
 $\frac{1}{8}$

the entropy is easily calculated:

$$H(x) = -\frac{1}{8}\log\frac{1}{8} - \frac{1}{8}\log\frac{1}{8} - \dots - \frac{1}{8}\log\frac{1}{8}$$
$$= -\log\frac{1}{8}$$
$$= 3 \text{ bits.}$$

(Information entropy is measured in bits).

For a different distribution:

a
 b
 c
 d
 e
 f
 g
 h

$$\frac{1}{2}$$
 $\frac{1}{4}$
 $\frac{1}{8}$
 $\frac{1}{16}$
 $\frac{1}{64}$
 $\frac{1}{64}$
 $\frac{1}{64}$
 $\frac{1}{64}$

the entropy is different:

$$H(x) = -\frac{1}{2}\log\frac{1}{2} - \frac{1}{4}\log\frac{1}{4} - \frac{1}{8}\log\frac{1}{8} - \frac{1}{16}\log\frac{1}{16} - 4 \cdot \frac{1}{64}\log\frac{1}{64}$$
$$= \frac{1}{2} + \frac{1}{2} + \frac{3}{8} + \frac{4}{16} + \frac{24}{64}$$
$$= 2 \text{ bits.}$$

Lower entropy means that we learn less from hearing the 'answer' for this distribution – half the time, the answer is **a** so we are less surprised when it comes up. In the uniform distribution, all answers are equally surprising.

An efficient encoding scheme for the second distribution is:

a	b	С	d	е	f	g	h
0	10	110	1110	111100	111101	111110	111111

No 'comma' is needed – a string of bits can be uniquely resolved into characters.

The average length of one letter is:

$$\langle L \rangle = \frac{1}{2} \cdot 1 + \frac{1}{4} \cdot 2 + \frac{1}{8} \cdot 3 + \frac{1}{16} \cdot 4 + 4 \cdot \frac{1}{64} \cdot 6 = 2$$
 bits.

A naïve code such as:

a	b	С	d	е	f	g	h
000	001	010	011	100	101	110	111

which again needs no comma, would have $\langle L \rangle = 3$. This code is however optimal for the uniform distribution.

The average code length is closely related to the entropy, but they are not always equal as in these examples. Lemma H(x) > 0

Lemma. $H(x) \ge 0$.

Proof. $-\log p_X(x) \ge 0$, since $0 \le p_X(x) \le 1$.

Lemma. $H_b(x) = \log_b a H_a(x)$.

Proof.

$$\log_a B = \frac{\log_c B}{\log_c A}$$

$$\implies \log_b a H_a(x) = -\sum_x p(x) \log_a b \log_a x$$
$$= -\sum_x p(x) \log_b x$$
$$= H_b(x).$$

2.2 Joint entropy

If $p_{X,Y}(x, y)$ is a joint distribution, the *joint entropy* is defined as:

$$H(X, Y) = -\sum p_{X,Y}(x, y) \log p_{X,Y}(x, y),$$

i.e. the normal definition applied to a joint distribution.

2.3 Conditional entropy

The entropy of a conditional distribution $p_{X|Y}(x|y)$ is:

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

The *conditional entropy* is the average of the entropies of the conditional distributions:

$$H(X|Y) = \sum_{y \in \mathscr{Y}} p_Y(y) H(X|Y = y)$$

= $-\sum_{y \in \mathscr{Y}} p_Y(y) \sum_{x \in \mathscr{X}} p_{X|Y}(x|y) \log p_{X|Y} p_{X|Y}(x|y)$
= $-\sum_{y \in \mathscr{Y}} p_{X,Y}(x, y) \log p_{X|Y}(x|y),$

since $p_Y(y)p_{X|Y}(x|y) = p_{X,Y}(x, y)$ by Bayes' rule.

In other words,

$$H(X|Y) = -\langle \log p_{X|Y}(x|y) \rangle_{p_{X,Y}}$$

2.4 Chain rule for entropy

Theorem.

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

Proof.

$$H(X, Y) = -\sum p_{X,Y}(x, y) \log p_{X,Y}(x, y)$$

= $-\sum p_{X,Y} \log [p_Y(y) p_{X|Y}(x|y)]$
= $-\sum p_{X,Y} \log p_Y(y) \underbrace{-\sum p_{X,Y} \log p_{X|Y}(x|y)}_{H(X|Y)}$

As for the first term,

$$\sum_{x,y} p_{X,Y}(x,y) \log p_Y(y) = \sum_{y} \log p_Y(y) \underbrace{\sum_{x} p_{X,Y}(x,y)}_{p_Y(y)}$$
$$= \sum_{y} p_Y(y) \log p_Y(y)$$
$$= H(Y).$$

In words, the chain rule states: "The information in *X* and *Y* is equal to the information in *Y* plus the further information in *X* given that you know *Y*."

A corollary is that:

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

2.5 Kullback–Leibler divergence

Say $\mathscr{X} = \mathscr{Y}$, $p(x) \equiv p_X(x)$, $q(x) \equiv p_Y(x)$. The Kullback–Leibler divergence (or KL divergence, sometimes *relative entropy*) between *p* and *q* is:

$$D(p||q) = \sum_{x \in \mathscr{X}} p(x) \log \frac{p(x)}{q(x)}.$$

In other words,

$$D(p||q) = \left\langle \log \frac{p(x)}{q(x)} \right\rangle_p.$$

The important properties of the KL divergence are:

$$D(p||q) \ge 0 \quad \forall p, q$$

$$D(p||q) = 0 \quad \Longleftrightarrow p(x) = q(x) \forall x \in \mathscr{X}.$$

The KL distribution is somewhat like a 'distance' between two probability distributions. Note that it is *not* a metric – it is not symmetric!

In many applications of the KL divergence, you will have a "real" distribution and a "model" distribution, and want to minimise the divergence between them. It is commonly used in AI, where it is often as practical as the standard L^2 distance:

$$D_2 = \sqrt{\sum \left(p(x) - q(x)\right)^2}$$

One practical problem with the KL divergence is that if $p(x) \neq 0$ and q(x) = 0 for some $x \in \mathcal{X}$, then D(p || q) is undefined.

2.6 Mutual information

The *mutual information* between X and Y is:

$$I(X; Y) = D(p_{X,Y}(x, y) || p_X(x) p_Y(y)),$$

or "how far the distributions are from being independent". Or using the definition of D,

$$I(X;Y) = \sum p_{X,Y}(x,y) \log \frac{p_{X,Y}(x,y)}{p_X(x)p_Y(y)}.$$

2.7 Relationship between all this stuff

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

= $\sum p(x,y) \log \frac{p(x)p(y|x)}{p(x)p(y)}$
= $\sum_{x,y} p(x,y) \log p(y|x) - \sum_{x,y} p(x,y) \log p(y)$
= $-H(Y|X) + H(Y).$

In other words, H(Y) = H(Y|X) + I(X;Y), "the information in *Y* equals the information left in *Y* when you know *X* plus the information knowing *X* gives you about *Y*."

Another form is I(X; Y) = H(X) + H(Y) - H(X, Y), so it is somehow a measure of the *overlap* of *X* and *Y*.

2.8 Chain rule for entropy

Theorem.

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

For example,

$$H(X_1, X_2, X_3) = H(X_1 | X_2, X_3) + H(X_2 | X_3) + H(X_3).$$

Proof.

$$H(X_1, X_2) = H(X_1) + H(X_2|X_1)$$
$$\implies H(X_1, X_2, X_3) = H(X_1) + H(X_2, X_3|X_1)$$
$$= H(X_1) + H(X_3|X_2, X_1) + H(X_2|X_1)$$

and so on.

2.9 Chain rule for mutual information

We define *conditional mutual information*:

$$I(X;Y|Z) = \sum p(x, y, z) \log \frac{p(x, y|z)}{p(x|z)p(y|z)},$$

the mutual information of the conditional probabilitied averaged over the condition. By Bayes' rule we can write:

$$I(X;Y|Z) = \sum p(z) \sum p(x,y|z) \log \frac{p(x,y|z)}{p(x|z)p(y|z)}$$
$$= \sum p(z)I(X;Y|Z=z)$$

Theorem.

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

Proof.

$$I(X_1, \dots, X_n; Y) = H(X_1, \dots, X_n) - H(X_1, \dots, X_n | Y)$$

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

2.10 Chain rule for KL divergence

We define *conditional KL divergence*:

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

This is just the average of the KL divergence of the conditional distributions. By Bayes' rule it can also be written:

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

Theorem.

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

Proof.

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

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

3 Jensen's inequality etc.

3.1 Cup-like & cap-like functions

Eschewing the terms "convex" and "concave", we say that a function *f* is *cup-like* if for all $a < b, \lambda \in (0, 1)$:

$$f\bigl((1-\lambda)a+\lambda b\bigr) \leq (1-\lambda)f(a)+\lambda f(b),$$

i.e. the function is below the straight line, as in Fig. 1.

We say *f* is strictly cup-like if the inequality is strict, and *f* is *cap-like* if -f is cup-like.



Figure 1: A cup-like function

3.2 Jensen's inequality

Theorem. If f is cup-like on X,

$$\langle f(X) \rangle \ge f(\langle X \rangle).$$

Furthermore if f is strictly cup-like,

$$\langle f(X) \rangle = f(\langle X \rangle) \Longleftrightarrow X = \langle X \rangle,$$

i.e. X is constant.

Proof. Take the case $\mathcal{X} = \{x_1, x_2\}$ in which case the inequality follows directly from the definition of cup-like, replacing $\lambda \rightarrow p_1$:

$$p_1 f(x_1) + p_2 f(x_2) \ge f(p_1 x_1 + p_2 x_2),$$

Now assume that the inequality holds for $\mathscr{X} = \{x_1, \dots, x_{n-1}\}$. Assuming $p_n \neq 1$, let

$$p_i'=\frac{p_i}{1-p_n},$$

so that $(p'_1, ..., p'_n)$ is a probability distribution on $\mathcal{X}' = \{x_1, ..., x_{n-1}\}$. So:

$$\sum_{i=1}^{n} p_i f(x_i) = p_n f(x_n) + (1 - p_n) \sum_{i=1}^{n-1} p'_i f(x_i)$$

$$\geq p_n f(x_n) + (1 - p_n) f\left(\sum p'_i x_i\right) \qquad \text{(by induce}$$

$$\geq f\left(p_n x_n + (1 - p_n) \sum p'_i x_i\right) \qquad \text{(by define}$$

$$= f\left(\sum p_i x_i\right)$$
i.e. $\langle f(X) \rangle \geq f(\langle X \rangle).$

(by induction assumption) (by definition of cup-like)

If f(x) is strictly cup-like, the 2-point case reads:

$$f(p_1x_1 + p_2x_2) = p_1f(x_1) + p_2f(x_2) \iff p_1 = 1 \text{ or } p_2 = 1,$$

and proceed as above.

3.3 Information inequality

This is also called *Gibbs' inequality*. **Theorem.** Let p(x) and q(x) be probability distributions on \mathscr{X} . Then,

 $D(p \| q) \ge 0,$

with $D(p||q) = 0 \iff p(x) = q(x) \forall x \in \mathcal{X}$.

Proof. Let \mathscr{A} be the support of $p: \mathscr{A} = \{x \in \mathscr{X} : p(x) \neq 0\}$. We can write:

$$-D(p\|q) = -\sum_{x \in \mathscr{A}} p(x) \log \frac{p(x)}{q(x)}$$

Now define a new random variable *Y*:

$$y_i = \frac{q(x_i)}{p(x_i)}$$
, with $p_Y(y_i) = p_X(x_i)$.

We want to apply Jensen's inequality to this new variable.

$$-D(p || q) = -\sum_{x \in \mathscr{A}} p(x) \log \frac{p(x)}{q(x)}$$

= $\langle \log Y \rangle$
 $\leq \log \langle Y \rangle$, (by Jensen) (1)
= $\log \sum_{x \in \mathscr{A}} p(x) \frac{q(x)}{p(x)}$
= $\log \sum_{x \in \mathscr{A}} q(x)$

So,

$$-D(p||q) \leq \log \sum_{x \in \mathcal{A}} q(x)$$

$$\leq \log \sum_{x \in \mathcal{X}} q(x)$$

$$= \log 1 = 0$$

$$\Longrightarrow D(p||q) \geq 0.$$
 (2)

log is strictly cap-like, so if (1) is an equality, then $\frac{p(x)}{q(x)} = c$ with probability 1. This means:

$$\sum_{x \in \mathcal{A}} q(x) = c \sum_{x \in \mathcal{A}} p(x) = c$$
$$= \sum_{x \in \mathcal{X}} q(x) = 1,$$

if (2) is an equality.

Corollary.

$$I(X;Y) \ge 0,$$

with equality $\iff X$, Y independent.

Proof.

$$I(X; Y) = D(p_{X,Y}(x, y) || p_X(x) p_Y(y)) \ge 0,$$

with $D(p_{X,Y}(x, y) || p_X(x) p_Y(y)) = 0 \iff p_X, Y(x, y) = p_X(x) p_Y(y)$, i.e. X and Y are independent.

Corollary.

$$D(p(y|x)||q(y|x)) \ge 0,$$

with equality $\iff q(y|x) = p(y|x) \forall x, y$.

Proof. This follows directly from the information inequality with $p \rightarrow p(y|x)$ and $q \rightarrow q(y|x)$.

Corollary.

$$I(X;Y|Z) \ge 0,$$

with $I(X; Y|Z) = 0 \iff X$, Y conditionally independent, i.e. $p(x, y|z) = p(x|z)p(y|z) \forall x, y, z$.

Proof. This follows directly from the first corollary.

3.4 Upper bound on information

Theorem.

$$H(X) \leq \log |\mathcal{X}|$$

where $|\mathcal{X}|$ is the size of the set of outcomes, with equality $\iff p(x) = 1/|\mathcal{X}|$, i.e. a uniform *distribution*.

Proof. We have *X* and *p*(*x*). Define another distribution on \mathscr{X} given by $u(x) = \frac{1}{|\mathscr{X}|}$. Then:

$$\begin{split} 0 &\leq D(p \| u) = \sum_{x \in \mathcal{X}} p(x) \log \frac{p(x)}{u(x)} \\ &= -H(x) - \log \frac{1}{|\mathcal{X}|} \\ &= -H(X) + \log |\mathcal{X}|. \end{split}$$

The information inequality also gives that equality $\iff p(x) = u(x) \forall x$. \Box **Theorem.**

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

with equality $\iff X$, Y independent.

Proof.

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

_	_		
L		L	
L		L	

Theorem.

$$H(X_1, X_2, \dots, X_n) \leq \sum_{i=1}^n H(X_i),$$

with equality if X_1, \ldots, X_n all independent.

Proof. For n = 2:

$$\label{eq:eq:constraint} \begin{split} 0 &\leq I(X;Y) = H(X) + H(Y) - H(X,Y) \\ \Rightarrow H(X) + H(Y) \geq H(X,Y). \end{split}$$

Now use the chain rule:

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

$$\leq \sum_{i=1}^{n} H(X_i).$$

3.5 Log-sum inequality

Theorem. For non-negative numbers a_1, \ldots, a_n and b_1, \ldots, b_n ,

$$\sum a_i \log \frac{a_i}{b_i} \ge \left(\sum a_i\right) \log \frac{\sum a_i}{\sum b_i}$$

Proof. Assume without loss of generality that $a_i > 0$, $b_i > 0$. $f(t) = t \log t$ is cup-like (since f''(t) > 0). Define:

$$\alpha_i = \frac{b_i}{\sum_k b_k}, \quad t_i = \frac{a_i}{b_i},$$

and apply Jensen's inequality to $\mathcal{T} = \{t_i\}, p_T(t_i) = \alpha_i, f$:

$$\langle f(T) \rangle \ge f(\langle T \rangle)$$

$$\sum \alpha_i f(t_i) \ge f\left(\sum \alpha_i t_i\right)$$

$$\sum \frac{b_i}{\sum_k b_k} f\left(\frac{a_i}{b_i}\right) \ge f\left(\sum_i \frac{b_i}{\sum_k b_k} \frac{a_i}{b_i}\right)$$

$$\sum b_i f\left(\frac{a_i}{b_i}\right) \ge \sum b_k f\left(\frac{\sum a_i}{\sum b_i}\right)$$

$$\sum a_i \log \frac{a_i}{b_i} \ge (\sum a_i) \log \frac{\sum a_i}{\sum b_i}$$

L
J.

3.6 Information inequality

Theorem.

$$D(p \| q) \ge 0$$

Proof.

$$D(p||q) = \sum_{x \in \mathscr{X}} p(x) \log \frac{p(x)}{q(x)}$$
$$\geq \left(\sum_{x \in \mathscr{X}} p(x)\right) \log \frac{\sum p(x)}{\sum q(x)}$$
$$= 1 \log 1 = 0.$$

Equality $\iff \frac{p(x_i)}{q(x_i)} = \text{const.} \Rightarrow p(x_i) = q(x_i) \text{ for all } i.$

3.7 Cup-likeness of KL divergence

If (p_1, q_1) and (p_2, q_2) are two pairs of probability distributions, then:

$$(p_{\lambda}, q_{\lambda}) = \left(\lambda p_1 + (1 - \lambda)p_2, \lambda q_1 + (1 - \lambda)q_2\right)$$

is also a pair of probability distributions for $0 \le \lambda \le 1$. **Theorem.** For (p_1, q_1) and (p_2, q_2) , $\lambda \in [0, 1]$,

$$D(p_{\lambda}, q_{\lambda}) \leq \lambda D(p_1 || q_1) + (1 - \lambda) D(p_2 || q_2).$$

Proof. Consider some $x \in \mathcal{X}$.

$$\begin{aligned} & \left(\lambda p_1(x) + (1-\lambda)p_2(x)\right)\log\frac{\lambda p_1(x) + (1-\lambda)p_2(x)}{\lambda q_1(x) + (1-\lambda)q_2(x)} \leq \\ & \leq \lambda p_1(x)\log\frac{\lambda p_1(x)}{\lambda q_1(x)} + (1-\lambda)p_2(x)\log\frac{(1-\lambda)p_2(x)}{(1-\lambda)q_2(x)} \end{aligned}$$

by log-sum. Sum both sides over *x*:

$$D(\lambda p_1(x) + (1-\lambda)p_2(x) \|\lambda q_1(x) + (1-\lambda)q_2(x))\lambda D(p_1\|q_1) + (1-\lambda)D(p_2\|q_2)$$

3.8 Cap-likeness of entropy

Theorem.

$$H(\lambda p_1 + (1 - \lambda)p_2) \ge \lambda H(p_1) + (1 - \lambda)H(p_2).$$

Proof. X_1 and X_2 are two random variables on the same set of outcomes \mathscr{X} with distributions p_1 and p_2 ; let θ be a random variable with two outcomes: $p_{\theta}(a) = \lambda$, $p_{\theta}(b) = 1 - \lambda$. Consider the random process X_{θ} : $\theta = a \rightarrow X_1$, $\theta = b \rightarrow X_2$. So,

$$p_{X_{\theta}}(x) = \lambda p_1(x) + (1 - \lambda) p_2(x).$$

$$H(X|\theta) = \lambda H(X_1) + (1 - \lambda)H(X_2)$$
$$= \lambda H(p_1) + (1 - \lambda)H(p_2)$$
$$= -\sum p_1(x)\log p_1(x).$$

We know $H(X) = H(\lambda p_1 + (1 - \lambda)p_2)$, and $H(x) \ge H(X|\theta)$. So:

$$H(\lambda p_1 + (1 - \lambda)p_2) \ge \lambda H(p_1) + (1 - \lambda)H(p_2).$$

3.9 Another theorem

Theorem. Let $(X, Y) \sim p(x, y) = p(x)p(y|x)$. Then the mutual information I(X; Y) can be considered a function of p_X and $p_{Y|X}$, and is cup-like in p(x) for fixed p(y|x) and cap-like in p(y|x) for fixed p(x).

Proof in book.

4 Markov chains & sufficient statistics

4.1 Markov chains

Random variables *X*, *Y*, *Z* are said to form a *Markov chain* denoted $X \rightarrow Y \rightarrow Z$ if p(x, y, z) = p(x)p(y|x)p(z|y). Equivalently, *X* and *Z* are conditionally independent:

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

Note also that the definition is symmetric: $X \to Y \to Z \iff Z \to Y \to X$. Also, if Z = f(Y) then $X \to Y \to Z$.

4.2 Data processing inequality

Theorem. If $X \to Y \to Z$, then $I(X; Y) \ge I(X; Z)$.

Proof.

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

Since p(x, z|y) = p(x|y)p(z|y), then:

$$I(X; Z|Y) = \sum p(y) \sum \log \frac{p(x, z|y)}{\underbrace{p(x|y)p(z|y)}_{=1}} = 0$$

So $I(X; Y) = I(X; Z) + I(X; Y|Z) \ge I(X; Z)$, with equality if $I(X; Y|Z) = 0 \iff X \rightarrow Z \rightarrow Y$.

Corollary. *If* Z = g(Y) *then* $I(X; Y) \ge I(X; g(Y))$.

Proof. $X \rightarrow Y \rightarrow g(Y)$.

Corollary. *If* $X \to Y \to Z$, *then* $I(X; Y|Z) \leq I(X; Y)$.

Proof. Recall I(X; Y) + I(X; Z|Y) = I(X; Z) + I(X; Y|Z). Noting I(X; Z|Y) = 0 and $I(X; Z) \ge 0$ gives result.

4.3 Sufficient statistics

 $\{p_{\theta}(X)\}\$ is a set of probability distributions indexed by some underlying variable θ . T(X) is a statistic, e.g. mean, variance, mean and variance, ...

Since T(X) is a function of X,

$$\theta \to X \to T(X),$$

so $I(\theta; X) \ge I(\theta; T(X))$.

T(X) is a sufficient statistic if $I(\theta; X) = I(\theta; T(X))$. In other words, if:

$$\theta \to T(X) \to X.$$

For example, let $X_1, ..., X_n$ with $X_i \in \{0, 1\}$ be an independent and identically distributed (iid) set of random variables. Let $\theta = p_{X_i}(1)$. We want to estimate θ .

Let:

$$T(X) = \sum_{i=1}^{n} X_i.$$

Claim: T(X) is a sufficient statistic. We show this by showing that X is independent of θ given T(X).

$$\Pr\left((X_1,\ldots,X_n)=(x_1,\ldots,x_n)\Big|\sum X_i=k\right) = \begin{cases} \frac{1}{N} & \text{if } \sum x_i=k\\ 0 & \text{otherwise,} \end{cases}$$

where $N = {n \choose k}$. This is clearly independent of θ . This means $\theta \to T(X) \to X$, so T(X) is a sufficient statistic.

4.4 Fano inequality

X is a random variable. *Y* is a measurement. \hat{X} is a guess of *X*. We are interested in the probability of error.

If $H(X|Y) \sim H(X)$, then reconstructing *X* from *Y* is just a guess – i.e. the chance of error is not different from guessing without knowing *Y*. If H(X|Y) = 0, then there is no uncertainty left in *X* when given *Y*, so *Y* determines *X* and there will be no error. The Fano inequality quantifies this.

X is a random variable; *Y* is another variable which is related to it. \hat{X} is the reconstruction of *X* from *Y*; we usually imagine \hat{X} to be a function of *Y*. Note that in general $\mathscr{X} \neq \hat{\mathscr{X}}$. Denote the probability of error $p_e = \Pr(X \neq \hat{X})$. **Theorem.** If \hat{X} is an estimator: $X \to Y \to \hat{X}$, and $p_e = \Pr(X \neq \hat{X})$, then:

$$\begin{split} H(p_e) + p_e \log |\mathcal{X}| &\geq H(X|\hat{X}) \\ &\geq H(X|Y). \end{split}$$

This inequality is often weakened to:

$$1 + p_e \log |\mathcal{X}| \ge H(X|Y),$$

(since $H(p_e) \ge 1$) or:

$$p_e \ge \frac{H(X|Y) - 1}{\log|\mathcal{X}|}$$

Proof.

$$H(X|\hat{X}) \ge H(X|Y)$$

since $I(X; \hat{X}) \leq I(X; Y)$. So:

$$H(X) - H(X|\hat{X}) \le H(X) - H(X|Y).$$

Define:

$$E = \begin{cases} 1 & \text{if } \hat{X} \neq X \\ 0 & \text{if } \hat{X} = X \end{cases}$$

so that $p_E(1) = p_e$.

$$H(E, X|\hat{X}) = H(X|\hat{X}) + H(E|X, \hat{X})$$
$$= H(X|\hat{X}),$$

since *E* is known given *X* and \hat{X} .

$$H(E, X|\hat{X}) = H(E|\hat{X}) + H(X|E, \hat{X}).$$

We know:

$$H(E|\hat{X}) \leq H(E) = H(p_e)$$

and

$$\begin{split} H(X|E,\hat{X}) &= p_E(0)H(X|\hat{X},E=0) + p_E(1)H(X|\hat{X},E=1) \\ &\leq 0 + \log |\mathcal{X}|p_e, \end{split}$$

since $X = \hat{X}$ when E = 0, Putting this together gives the Fano inequality:

$$H(p_e) + p_e \log |\mathcal{X}| \ge H(X|Y)$$

Corollary. If $\hat{\mathscr{X}} \subseteq \mathscr{X}$ then the inequality can be tightened to:

$$H(p_e) + p_e \log(|\mathcal{X}| - 1) \ge H(X|\hat{X}).$$

Proof. Before, we had $H(X|\hat{X}, E = 1) \le \log |\mathscr{X}|$. However if $\hat{X} \in \mathscr{X}$ and we know E = 1, then *X* can only take $|\mathscr{X}| - 1$ different values.

Corollary. For any two random variables X and Y, let $p = Pr(X \neq Y)$. Then $H(p) + p\log|\mathcal{X}| \ge H(X|Y)$.

Proof. Simply apply the Fano inequality to the Markov chain $X \to Y \to Y$.

5 Asymptotic equipartition principle & typical sets

5.1 Convergence

Convergence: given a sequence of random variables $Y_1, Y_2, ...$ we say the sequence converges to Y:

- "in probability" if $\forall \varepsilon > 0$, $\Pr(|Y_n Y| > \varepsilon) \rightarrow 0$
- "in mean square" if $\langle (Y Y_n)^2 \rangle \rightarrow 0$.
- "with probability 1" or "almost surely" if:

$$\Pr\left(\lim_{n \to \infty} Y_n = Y\right) = 1$$

5.2 Law of large numbers

If X_1, X_2, \dots is a sequence of iid random variables with $X_i \sim X$, $\mu = \langle X \rangle$, $\overline{X}_n = \frac{1}{n} \sum_{i=1}^n X_i$,

- Weak law: $\overline{X}_n \rightarrow \mu$ in probability.
- **Strong law:** $\overline{X}_n \rightarrow \mu$ almost surely.

The asymptotic equipartition principle is basically the law of large numbers for information theory.

5.3 Asymptotic equipartition principle

The AEP states that:

$$p(X_1, X_2, \dots, X_n) \sim 2^{-nH(X)}$$

for iid random variables X_1, \ldots, X_n with $X_i \sim X$.

For example, if $\mathscr{X} = \{0, 1\}$ with probabilities p and q respectively, one possible sequence is (1, 1, ..., 1), but it is not the most likely: it has probability q^n .

Usually the sequence will consist of ~ np 0s and ~ nq 1s. Such a sequence has probability $p^{np}q^{nq}$, and:

$$p^{np}q^{nq} = 2^{\log p^{np}} 2^{\log q^{nq}}$$
$$= 2^{\log p^{np} + \log q^{nq}}$$
$$= 2^{-nH}.$$

So 2^{-nH} is, roughly speaking, the probability of a typical outcome. **Theorem.** *If* $X_1, X_2, ..., X_n$ *are iid with* $X_i \sim X$ *, then:*

$$-\frac{1}{n}\log p(X_1, X_2, \dots, X_n) \to H(X) \text{ in probability.}$$

Proof. Functions of independent variables are independent: *Y* independent of $Z \Rightarrow f(Y)$ independent of f(Z). So $\{\log p(X_i)\}$ are independent variables. By the weak law of large numbers,

 $\frac{1}{N} \sum \log p(X_i) \to \langle \log p(X) \rangle$

The theorem follows from noting that $\langle \log p(x) \rangle = -H(X)$, and:

$$\sum \log p(X_i) = \log \prod p(X_i)$$
$$= \log p(X_1, \dots, X_n).$$

-		_
-	-	_

5.4 Typical set

The typical set $A_{\varepsilon}^{(n)}$ is the set of elements $(x_1, x_2, \dots, x_n) \in \mathcal{X}^n$ such that:

$$2^{-n(H(X)+\varepsilon)} \le p(x_1, \dots, x_n) \le 2^{-n(H(X)-\varepsilon)}$$

We saw that $p(X_1,...,X_n) \rightarrow 2^{-nH(X)}$, so $A_{\varepsilon}^{(n)}$ consists of the $(x_1,...,x_n)$ which are close to the limit.

Theorem. 1. If $(x_1, \ldots, x_n) \in A_{\varepsilon}^{(n)}$, then:

$$H(X) - \varepsilon \leq -\frac{1}{n} \log p(x_1, \dots, x_n) \leq H(X) + \varepsilon$$

- 2. $\Pr(A_{\varepsilon}^{(n)}) > 1 \varepsilon$ for sufficiently large n.
- 3. $|A_{\varepsilon}^{(n)}| \leq 2^{n(H(X)+\varepsilon)}$, and:
- 4. $|A_{\varepsilon}^{(n)}| \ge 2^{n(H(X)-\varepsilon)}$.

Proof. 1. Follows from taking the log of the definition of the typical set.

2. We know:

 $\frac{1}{n}\log p(X_1,\ldots,X_n) \to H(X)$ in probability.

In other words:

$$\Pr\left(\left|-\frac{1}{n}\log p(X_1,\ldots,X_n)\right| < \varepsilon\right) \to 1.$$

or in terms of δ ,

$$\Pr\left(\left|-\frac{1}{n}\log p(X_1,\ldots,X_n)\right| < \varepsilon\right) > 1 - \delta$$

and let $\delta = \epsilon$ for result.

3.

$$1 = \sum_{\mathbf{x}\in\mathscr{X}^n} p(\mathbf{x}) \ge \sum_{\mathbf{x}\in A_{\varepsilon}^{(n)}} p(\mathbf{x}) \ge \sum_{\mathbf{x}\in A_{\varepsilon}^{(n)}} 2^{-n(H(X)+\varepsilon)}.$$

So:

$$\begin{split} 1 &\geq 2^{-n(H(X)+\varepsilon)} \sum_{\mathbf{x} \in A_{\varepsilon}^{(n)}} 1 \\ &= 2^{-n(H(X)+\varepsilon)} \left| A_{\varepsilon}^{(n)} \right|. \end{split}$$

$$\begin{split} \mathbf{l} - \varepsilon &< \Pr(A_{\varepsilon}^{(n)}) \\ &= \sum_{\mathbf{x} \in A_{\varepsilon}^{(n)}} p(\mathbf{x}) \\ &\leq \sum_{\mathbf{x} \in A_{\varepsilon}^{(n)}} 2^{-n(H(X) - \varepsilon)} \\ &= \left| A_{\varepsilon}^{(n)} \right| 2^{-n(H(X) - \varepsilon)}. \end{split}$$

	-	-	-	

5.5 Code word length

Let $X_1, X_2,...$ be iid with probability distribution p(x), and we wish to find a short description of a length-n output. We can construct a code word by letting the first bit be 0 if $\mathbf{x} \in A_{\varepsilon}^{(n)}$ or 1 otherwise. Then assign some order to elements of both $A_{\varepsilon}^{(n)}$ and $\mathscr{X}^n \setminus A_{\varepsilon}^{(n)}$, convert the number of the element to binary and concatenate.

For example, the 18th element of the typical set would have a code word 0 10010 by this scheme.

The length of this binary number must be less than $\log|\text{set}| + 1$. Hence the elements of the typical set have code words of length $l < \log 2^{n(H+\varepsilon)} + 1 + 1 = n(H+\varepsilon) + 2$. Similarly words not in the typical set have length $l < \log|\mathcal{X}^n \setminus A_{\varepsilon}^{(n)}| + 2$, $\Rightarrow l < \log|\mathcal{X}^n| + 2$.

What is the average length of a code word?

$$\langle l(\mathbf{x}) \rangle = \sum_{x \in A_{\varepsilon}^{(n)}} p(\mathbf{x}) l(\mathbf{x})$$

$$= \sum_{x \in A_{\varepsilon}^{(n)}} p(\mathbf{x}) l(\mathbf{x}) + \sum_{x \notin A_{\varepsilon}^{(n)}} p(\mathbf{x}) l(\mathbf{x})$$

$$\leq \sum_{x \in A_{\varepsilon}^{(n)}} p(\mathbf{x}) [n(H+\varepsilon)+2] + \sum_{x \notin A_{\varepsilon}^{(n)}} p(\mathbf{x}) [\log |\mathscr{X}^{n}|+2]$$

$$= \Pr(A_{\varepsilon}^{(n)}) [n(H+\varepsilon)+2] + \Pr(\mathscr{X}^{n} \setminus A_{\varepsilon}^{(n)})$$

$$\leq (1-\varepsilon) [n(H+\varepsilon)+2] + \varepsilon [\log |\mathscr{X}^{n}|+2]$$

$$= n(H+\varepsilon'),$$

where $\varepsilon' \equiv \varepsilon + \varepsilon \log |\mathcal{X}| + \frac{2}{n}$. So: $\frac{\langle l(\mathbf{x}) \rangle}{n} \leq H + \varepsilon'$.

 ε' can be made arbitrarily small for sufficiently large *n* and suitable choice of ε . This means that for iid X_1, \ldots, X_n with $X_i \sim X$, the code length per outcome can be made arbitrarily close to *H*.

6 Encoding

6.1 Definitions

Ideally we want common outcomes to have short code words, and rare outcomes to have long code words. For example, in Morse code, e is \cdot , t is –, a is \cdot –, etc.

We define a *source code C* for a variable *X* as a mapping:

$$C: \mathscr{X} \to \mathscr{D}^*, \quad x \mapsto C(x)$$

where \mathcal{D}^* is the set of finite-length strings of elements from the *D*-symbol alphabet \mathcal{D} . The length of C(x) is denoted l(x).

The *expected length* of a source code, *L*(*C*) is the average length:

$$L(C) = \langle l(x) \rangle$$
$$= \sum_{x \in \mathscr{X}} p(x) l(x).$$

We say a code is *non-singular* if every $x \in \mathcal{X}$ maps to a different string:

$$x \neq x' \Rightarrow C(x) \neq C(x').$$

i.e. *c* is injective.

The *extension* C^* of a code *C* is the induced map:

$$C^*:\mathscr{X}^*\to\mathscr{D}^*$$

given by concatenation:

$$C(x_1 x_2 \dots x_n) = C(x_1) C(x_2) \dots C(x_n).$$

A code is *instantaneous* (or a *prefix code*) if no code word is the prefix of another, e.g.:

A prefix code (e.g. 01101010111) can be uniquely split up:

$$\begin{array}{c} 0 1101010111 \\ a & c & b & b & d \end{array}$$

However, codes can have this property without being prefix codes.

6.2 Kraft inequality

Theorem. For any prefix code over an alphabet of size D with code word lengths $l_1 = l(x_1), l_2, ..., l_n$,

$$\sum_{i=1}^n D^{-l_i} \leq 1.$$

Conversely, given $\{l_i\}$ satisfying this inequality it is possible to construct a prefix code.

Note that this inequality has nothing to do with efficiency - it is about prefix-ness.

Proof. Construct a tree where code words give a 'route map'. Each node in the tree has $\leq D$ outgoing connections labelled $0, \dots, D-1$. For example, here outcomes a - k are assigned code words from a 3-nary alphabet:



This example is a prefix code, since all of the outcomes correspond to *leaves*: nodes with no outgoing connections. Leaves have no *descendants*: nodes which a given node prefixes.

Call the length of the longest code word l_* for a given prefix code. If all leaves at this level were in the code, we would have D^{l_*} leaves. If there is a code word of length $l_i < l_*$, then it would have had $D_{l_*-l_i}$ descendants at length l_* . For example, here a code word of length 1 would have had $3^{3-1} = 9$ descendants of length 3.



Because this is a prefix code, no two code words share any descendants. So the number of children at level l_* must be less than or equal to the maximum number of children at

level l_* :

$$\sum_{i=1}^n D^{l_*-l_i} \leq D^{l_*}$$

Dividing by D^{l_*} gives the Kraft inequality.

Conversely if l_1, \ldots, l_n satisfy the Kraft inequality, we can choose disjoint family subsets of size $D^{l_*-l_i}$ at level l_* and cut them off at their parent to form a prefix code.

Note that this is for $|\mathscr{X}|$ finite. There does exist an extension for countably infinite \mathscr{X} .

6.3 Optimal codes

The idea is to find a lower bound on L(C). We can do this by minimising $L(C) = \sum p_i l_i$ subject to $\sum D^{-l_i} \le 1$. Use Lagrange multipliers!

$$J = \sum p_i l_i + \lambda \left(\sum D^{-l_i} - 1 \right).$$

Now minimise *J*:

$$\frac{\partial J}{\partial l_i} = p_i - \lambda D^{-l} \ln D_i = 0$$
$$\Rightarrow D^{-l_i} = \frac{p_i}{\lambda \ln D}$$

Sum both sides over *i*:

$$1 = \frac{1}{\lambda \ln D}$$

and substitute back in:

$$D^{-l_i} = p_i$$

So the optimal length, l_i^* , is given by $-\log_D p_i$. The average length of a code word is given by:

$$L^* = \sum p_i l_i^*$$

= $-\sum p_i \log_D p_i$
= $H_D(X)$.

Note that $-\log_D p_i$ is in general not an integer. We say that the best code has $l_i = \lceil l_i^* \rceil$, i.e. l_i^* rounded up to the nearest integer.

The Shannon-Fano theorem states that for an optimal prefix code:

$$H_D \le L \le H_D + 1.$$

In other words, the average code word length of an optimal prefix code is within 1 digit of the minimum possible average length.

6.4 Huffman coding

An elegant scheme for constructing an optimal prefix code.

Begin with the outcomes as leaves labelled by their probabilities. Then repeatedly join the *D* nodes with the lowest probabilities (sometimes this can be done in more than one way) and label the new node with the sum of these probabilities. Here this is performed for D = 2:



Labelling routes yields code words for each outcome:



This gives a prefix code where the most probable outcomes have the shortest code words.

For $D \ge 3$ there may not be the correct number of symbols. Since each joining replaces D symbols with 1, and we want 1 left at the end, this alogrithm only works for numbers of the form k(D - 1) + 1. To solve this, introduce a sufficient number of 0-probability nodes. For example, for D = 3:



This means that there will be codes that do not correspond to any outcomes (for example, 222 above). This means the code is in some sense not as efficient as it 'could be'. However the code will still have $L \le H_D + 1$.