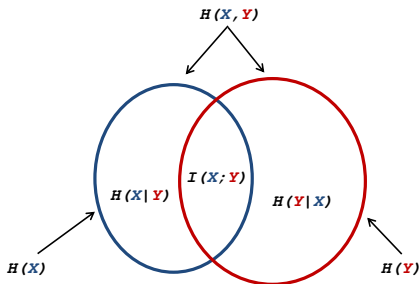


Information Theory

Professor John Daugman

University of Cambridge

Computer Science Tripos, Part II
Michaelmas Term 2015/16



Outline of Lectures

1. **Foundations: probability, uncertainty, information.**
2. **Entropies defined, and why they are measures of information.**
3. **Source coding theorem; prefix, variable-, and fixed-length codes.**
4. **Discrete channel properties, noise, and channel capacity.**
5. **Spectral properties of continuous-time signals and channels.**
6. **Continuous information; density; noisy channel coding theorem.**
7. **Signal coding and transmission schemes using Fourier theorems.**
8. **The quantised degrees-of-freedom in a continuous signal.**
9. **Gabor-Heisenberg-Weyl uncertainty relation. Optimal “Logons”.**
10. **Data compression codes and protocols.**
11. **Kolmogorov complexity. Minimal description length.**
12. **Applications of information theory in other sciences.**

Reference book

(* Cover, T. & Thomas, J.

Elements of Information Theory (second edition).

Wiley-Interscience, 2006

Overview: what is information theory?

Key idea: The movements and transformations of information, just like those of a fluid, are constrained by mathematical and physical laws.

These laws have deep connections with:

- ▶ probability theory, statistics, and combinatorics
- ▶ thermodynamics (statistical physics)
- ▶ spectral analysis, Fourier (and other) transforms
- ▶ sampling theory, prediction, estimation theory
- ▶ electrical engineering (bandwidth; signal-to-noise ratio)
- ▶ complexity theory (minimal description length)
- ▶ signal processing, representation, compressibility

As such, information theory addresses and answers the two fundamental questions which limit all data encoding and communication systems:

1. What is the ultimate data compression?
(*answer: the **entropy of the data**, H , is its compression limit.*)
2. What is the ultimate rate of reliable communication?
(*answer: the **channel capacity**, C , is its transmission rate limit.*)

Information theory studies ways to achieve these two theoretical limits.

Important questions to which information theory offers answers:

- ▶ How should information be measured?
- ▶ How much additional information is gained by some reduction in uncertainty?
- ▶ How do the *a priori* probabilities of possible messages determine the informativeness of receiving them?
- ▶ What is the information content of a random variable?
- ▶ How does the noise level in a communication channel limit its capacity to transmit information?
- ▶ How does the bandwidth (in cycles/second) of a communication channel limit its capacity to transmit information?
- ▶ By what formalism should prior knowledge be combined with incoming data to draw formally justifiable inferences from both?
- ▶ How resistant is a cryptographic key to a brute force attack?
- ▶ How much information is contained in a strand of DNA?
- ▶ How much information is there in the firing pattern of a neurone?

1. Foundations: probability, uncertainty, information

- ▶ *How the concepts of randomness, redundancy, compressibility, noise, bandwidth, and uncertainty are related to information.*
- ▶ *Ensembles, random variables, marginal and conditional probabilities.*
- ▶ *How metrics of information are grounded in the rules of probability.*

Random variables are variables that take on values determined by probability distributions. They may be discrete or continuous, in their domain or in their range.

For example, a stream of ASCII encoded text characters in a transmitted message is a discrete random variable, taking on discrete values that have a known probability distribution for any given natural language.

An analog speech signal represented by a voltage or sound pressure waveform as a function of time (perhaps with added noise), is an example of a continuous random variable described by a continuous probability density function.

Most of Information Theory involves probability distributions of random variables, and conjoint or conditional probabilities defined over ensembles of random variables. Indeed, the information content of a symbol or event is defined in terms of its (im)probability.

Classically, there are two different points of view about what probability actually means:

- ▶ *relative frequency*: sample the random variable a great many times and tally up the fraction of times that each of its different possible values occurs, to arrive at the probability of each.



- ▶ *degree-of-belief*: probability is the plausibility of a proposition or the likelihood that a particular state (or value of a random variable) might occur, even if its outcome can only be decided once (such as the outcome of a particular horse-race).

The first view, the “frequentist” or operationalist view, is the one that predominates in statistics and in information theory. However, it does not capture the full meaning of probability.

For example, the proposition "The moon is made of green cheese" is one which surely has a probability that we should be able to attach to it. We could assess its probability by degree-of-belief calculations which combine our prior knowledge about physics, geology, and dairy products.

Yet it seems the “frequentist” definition of probability could only assign a probability to this proposition by performing (say) a large number of repeated trips to the moon, and tallying up the fraction of trips in which the moon turned out to be a dairy product....

In either case, it seems sensible that the less probable an event is, the more information is gained by noting its occurrence. (Surely discovering that the moon IS made of green cheese would be more “informative” than merely learning that it is just made of earth-like rocks.)

Probability rules

Most of probability theory was laid down by theologians: Blaise PASCAL (1623-1662) who gave it the axiomatization that we accept today; and Thomas BAYES (1702-1761) who expressed one of its most important and widely-applied propositions relating conditional probabilities.

Probability Theory rests upon two rules:

Product Rule:

$$\begin{aligned} p(A, B) &= \text{"joint probability of both } A \text{ and } B\text{"} \\ &= p(A|B)p(B) \end{aligned}$$

$$\begin{aligned} &\text{or equivalently,} \\ &= p(B|A)p(A) \end{aligned}$$

Clearly, in case A and B are *independent* events, they are not conditionalised on each other and so

$$p(A|B) = p(A) \quad \text{and} \quad p(B|A) = p(B),$$

in which case their joint probability is simply $p(A, B) = p(A)p(B)$.

Sum Rule:

If event A is conditionalised on a number of other events B , then the total probability of A is the sum of its joint probabilities with all B :

$$p(A) = \sum_B p(A, B) = \sum_B p(A|B)p(B)$$

From the Product Rule and the symmetry that $p(A, B) = p(B, A)$, it is clear that $p(A|B)p(B) = p(B|A)p(A)$. Bayes' Theorem then follows:

Bayes' Theorem:

$$p(B|A) = \frac{p(A|B)p(B)}{p(A)}$$



The importance of Bayes' Theorem is that it allows us to reverse the conditionalising of events, and to compute $p(B|A)$ from knowledge of $p(A|B)$, $p(A)$, and $p(B)$. Often these are expressed as *prior* and *posterior* probabilities, or as the conditionalising of hypotheses upon data.

Worked Example:

Suppose that a dread disease affects 1/1000th of all people. If you actually have the disease, a test for it is positive 95% of the time, and negative 5% of the time. But if you don't have the disease, the test is positive 5% of the time. We wish to know how to interpret test results.

Suppose you take the test and it is positive. What is the likelihood that you actually have the disease?

We use the above rules, with the following substitutions of “data” D and “hypothesis” H instead of A and B :

D = data: the test is positive

H = hypothesis: you have the disease

\bar{H} = the other hypothesis: you do not have the disease

Before acquiring the data, we know only that the *a priori* probability of having the disease is .001, which gives us $p(H)$. This is called a *prior*. We also need to know $p(D)$.

From the Sum Rule, we can calculate that the *a priori* probability $p(D)$ of testing positive, whatever the truth may actually be, is:

$$p(D) = p(D|H)p(H) + p(D|\bar{H})p(\bar{H}) = (.95)(.001) + (.05)(.999) = .051$$

and from Bayes' Rule, we can conclude that the probability that you actually have the disease given that you tested positive for it, is much smaller than you may have thought:

$$p(H|D) = \frac{p(D|H)p(H)}{p(D)} = \frac{(.95)(.001)}{(.051)} = \boxed{0.019} \quad (\text{less than } 2\%).$$

This quantity is called the *posterior probability* because it is computed after the observation of data; it tells us how likely the hypothesis is, given what we have observed.

(Note: it is an extremely common human fallacy to confound $p(H|D)$ with $p(D|H)$. In the example given, many people would react to the positive test result by concluding that the likelihood that they have the disease is .95, since that is the “hit rate” of the test. They confound $p(D|H) = .95$ with $p(H|D) = .019$, which is what actually matters.)

A nice feature of Bayes' Theorem is that it provides a simple mechanism for repeatedly updating our assessment of the hypothesis as more data continues to arrive. We can apply the rule recursively, using the latest *posterior* as the new *prior* for interpreting the next set of data. In Artificial Intelligence, this feature is important because it allows the systematic and real-time construction of interpretations that can be updated continuously as more data arrive in a time series, such as a stream of images or spoken sounds that we wish to understand.

Information Theory allows us to analyse quantitatively the amount of uncertainty that is reduced, *i.e.* the amount of information that is gained, from an inference using Bayes' Theorem. Now we must develop such metrics that operate on probabilities.

2. Entropies defined, and why they are measures of information.

- ▶ *Marginal entropy, joint entropy, conditional entropies, and the Chain Rule for entropy. The “distance” between random variables.*
- ▶ *Mutual information between random variables. Independence.*

The information measure I of a single event or message is defined as the base-2 logarithm of its probability p of occurring:

$$I = \log_2 p$$

and its *entropy* H is considered the inverse: $H = -I$. Entropy can be regarded intuitively as “uncertainty,” or “disorder.” To gain information is to lose uncertainty by the same amount, so I and H differ only in sign. Entropy and information have units of *bits*.

Note that $I = \log_2 p$ is never positive: it ranges from 0 to $-\infty$, and thus H ranges positively from 0 to $+\infty$, as p varies from 1 to 0.

No information is gained (no uncertainty is reduced) by the appearance of an event or the receipt of a message that was completely certain anyway: $p = 1$, therefore $H = I = 0$. Intuitively, the more improbable an event is, the more significant it is; so the monotonic behaviour seems appropriate.

But why the logarithm?

The logarithmic measure is justified by the desire that information be additive. We want the algebra of our measures to reflect the rules of probability. When independent packets of information arrive, we would like to say that the total information received is the sum of the individual pieces. But the probabilities of independent events multiply to give their combined probabilities, and so we must take logarithms in order for the joint probability of independent events or messages to be combined additively into the total information gained.

This principle can also be understood in terms of the combinatorics of state spaces. Suppose we have two independent problems, one with n possible solutions (or states) each having probability p_n , and the other with m possible solutions (or states) each having probability p_m . Then the number of combined states is mn , and each of these has probability $p_m p_n$. We would like to say that the information gained by specifying the solution to *both* problems is the *sum* of that gained from each one. This desired property is achieved:

$$I_{mn} = \log_2(p_m p_n) = \log_2 p_m + \log_2 p_n = I_m + I_n$$

A note on logarithms:

In information theory we often wish to compute the base-2 logarithms of quantities, but many calculating tools only offer Napierian (base 2.718...) or decimal (base 10) logarithms. So the following conversions are useful:

$$\log_2 X = 1.443 \log_e X = 3.322 \log_{10} X$$

It may be noted in passing that occasionally the “natural” or Napierian (base- e) logarithm is invoked, in which case the information measure is the “nat” or the “nit” (for Napierian bit). Thus 1 bit \approx 0.693 nits.

Henceforward we will omit the subscript; base-2 is always presumed.

You will find it very beneficial to commit to memory now all of the powers of 2 from about -8 to +8 (*i.e.* $\frac{1}{256}$, $\frac{1}{128}$, $\frac{1}{64}$, $\frac{1}{32}$, ..., 1, 2, 4, ..., 128, 256) because we will frequently encounter such numbers, and their base-2 logarithms should be immediately at your fingertips.

Intuitive Example of the Information Measure:

Suppose I select at random one of the 26 letters of the alphabet, and we play the game of “26 questions” in which you try to discover which letter I have chosen. I will only answer ‘yes’ or ‘no,’ always truthfully.

What is the minimum number of such questions that you must ask me, in order to guarantee finding the answer? - Perhaps 25 questions?

What form should such questions take? e.g., “Is it A?” ... “Is it B?” ... or, is there some more intelligent way to solve this problem?

The answer to a Yes/No question having equal probabilities conveys one bit worth of information. In the above example with equiprobable states, you never need to ask more than 5 (well-phrased!) questions to discover the answer, even though there are 26 possibilities.

The information measure tells us that the uncertainty removed as a result of solving this problem corresponds to about 4.7 bits.

Entropy of ensembles

We now move from considering the information content of a single event or message, to that of an *ensemble*. An ensemble is the set of outcomes of one or more random variables. The outcomes have known probabilities attached to them. In general these probabilities are non-uniform, with event i having probability p_i , but they must sum to 1 because all possible outcomes are included; hence they form a discrete probability distribution:

$$\sum_i p_i = 1$$

The *entropy of an ensemble* is simply the average entropy of all of the elements in it. We can compute their average entropy by weighting each of the $\log p_i$ contributions by its probability p_i of occurring:

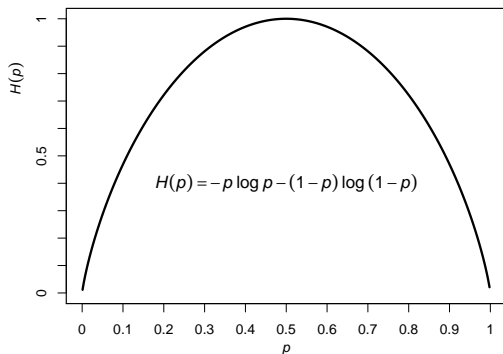
$$H = -I = - \sum_i p_i \log p_i$$

Thus we can speak of the information content or the entropy of a random variable, from knowledge of the probability distribution that it obeys. (*Entropy does not depend upon the actual values taken by the random variable! – Only upon their relative probabilities.*)

Entropy of a binary random variable with outcome probability p

Let us consider a random variable that takes only two values, one with probability p , and the other one with probability $1 - p$. (This is called a Bernoulli process, with parameter p .) How does the entropy of this binary random variable depend on the value of p ?

Plotting $H = -\sum_i p_i \log p_i$ where the index i spans the two possible outcomes, as a function of p shows that the entropy is a symmetric concave function that equals 0 if $p = 0$ or if $p = 1$, and it reaches a maximum value of 1 bit when $p = \frac{1}{2}$:



Entropy of a discrete random variable with non-uniform probabilities

The various letters of the written English language have the following relative frequencies (probabilities), in descending order:

E	T	O	A	N	I	R	S	H	D	L	C	..
.105	.072	.066	.063	.059	.055	.054	.052	.047	.035	.029	.023	..

If all 26 were equiprobable, the entropy of the ensemble would have been $H = -\sum_1^{26} (\frac{1}{26}) \log_2(\frac{1}{26}) = 4.7$ bits. But their non-uniform probabilities reveal that, for example, an E is nearly five times more likely than a C. Surely such prior knowledge must reduce the uncertainty of this random variable, and so the letter-guessing game should be even more efficient.

In fact, the distribution of English letters has an entropy of only 4.0 bits. This means that, on average, only four 'Yes/No' questions are necessary, to discover the secretly chosen letter of the 26 letters in the alphabet.

How can this possibly be true?

That is the subject of Claude Shannon's Source Coding Theorem (so named because it uses the “statistics of the source,” the *a priori* probabilities of the message generator, to construct an optimal code.) Note the important assumption: that the “source statistics” are known!

Shannon's seminal contributions, which essentially created the field of Information Theory, include two other key theorems that we will study: the Channel Coding Theorem (capacity for error-correcting codes); and the Noisy Channel Coding Theorem (channel capacity in Gaussian noise).



Several further measures of entropy first need to be defined, involving the marginal, joint, and conditional probabilities of random variables. Some key relationships will then emerge, that we can apply to the analysis of communication channels and codes.

More concepts and notation

We use capital letters X and Y to name random variables, and we use lower case letters x and y for instances of their respective outcomes.

These are drawn from particular sets \mathcal{A} and \mathcal{B} : $x \in \{a_1, a_2, \dots, a_J\}$, and $y \in \{b_1, b_2, \dots, b_K\}$. The probability of any particular outcome $p(x = a_i)$ is denoted p_i , for $0 \leq p_i \leq 1$ and with $\sum_i p_i = 1$.

An ensemble is just a random variable X . A joint ensemble 'XY' is an ensemble whose outcomes are ordered pairs x, y with $x \in \{a_1, a_2, \dots, a_J\}$ and $y \in \{b_1, b_2, \dots, b_K\}$. The joint ensemble XY defines a probability distribution $p(x, y)$ over all the JK possible joint outcomes x, y .

Marginal probability: From the Sum Rule, we can see that the probability of X taking on any particular value $x = a_i$ equals the sum of the joint probabilities of this outcome for X and all possible outcomes for Y :

$p(x = a_i) = \sum_y p(x = a_i, y)$. We usually simplify this notation for the

marginal probabilities to: $p(x) = \sum_y p(x, y)$ and $p(y) = \sum_x p(x, y)$.

Conditional probability: From the Product Rule, we can easily see that the conditional probability that $x = a_i$, given that $y = b_j$, is:

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

We usually simplify this notation for conditional probability to:

$$p(x|y) = \frac{p(x, y)}{p(y)} \quad \text{and similarly} \quad p(y|x) = \frac{p(x, y)}{p(x)}.$$

It is now possible to define various entropy measures for joint ensembles.

The key thing to notice about all of them is that they are instances of the basic $H = -\sum_i p_i \log p_i$ concept except that the p_i will now be instead a joint probability distribution, or a conditional probability distribution.

For convenience we will often absorb the minus sign into the logarithm by taking the reciprocal inside of the log; and also for convenience sometimes we will replace terms with others by applying the Sum or Product Rules.

Joint entropy of XY

$$H(X, Y) = \sum_{x,y} p(x, y) \log \frac{1}{p(x, y)}$$

Note that we have replaced the usual minus sign in front by taking the reciprocal of $p(x, y)$ inside the logarithm.

From this definition, it follows that joint entropy is additive if X and Y are independent random variables:

$$H(X, Y) = H(X) + H(Y) \quad \text{iff} \quad p(x, y) = p(x)p(y)$$

Otherwise, the joint entropy of the joint ensemble XY is less than the sum of the entropies $H(X)$ and $H(Y)$ of the individual random variables. The amount of that difference (see the Venn diagram on the first slide) will be one of the most important quantities that we encounter.

Conditional entropy of an ensemble X , given that for Y , $y = b_j$

...measures the uncertainty remaining about random variable X after specifying that random variable Y has taken on some particular value $y = b_j$. It is defined naturally as just the entropy of that conditional probability distribution $p(x|y = b_j)$:

$$H(X|y = b_j) = \sum_x p(x|y = b_j) \log \frac{1}{p(x|y = b_j)}$$

If we now consider the above quantity averaged over all the possible outcomes that random variable Y might have, each weighted by its corresponding probability $p(y)$, then we arrive at the...

Conditional entropy of an ensemble X , given an ensemble Y :

$$H(X|Y) = \sum_y p(y) \left[\sum_x p(x|y) \log \frac{1}{p(x|y)} \right]$$

and we know from the Sum Rule that if we move the $p(y)$ term from the outer summation over y , to inside the inner summation over x , the two probability terms combine and become just $p(x, y)$ summed over all x, y .

Hence a simpler expression for this conditional entropy is:

$$H(X|Y) = \sum_{x,y} p(x, y) \log \frac{1}{p(x|y)}$$

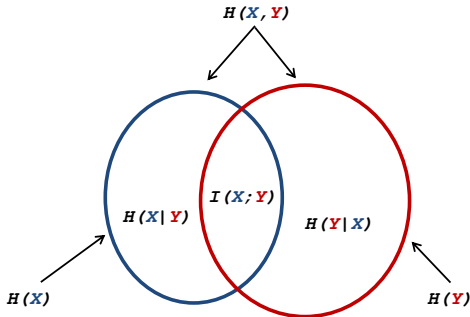
This measures the uncertainty that remains about X , when Y is known, averaged over all possible values of both random variables.

Chain Rule for Entropy

The joint entropy, conditional entropy, and marginal entropy for two random variables X and Y are related by:

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

It should seem natural and intuitive that the joint entropy of a pair of random variables is the entropy of one plus the conditional entropy of the other (the uncertainty that it adds once its dependence on the first one has been discounted by conditioning on it).



“Independence Bound on Entropy”

A consequence of the Chain Rule for Entropy is that if we have many different random variables X_1, X_2, \dots, X_n , then the sum of all their individual entropies must be an upper bound on their joint entropy:

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

Their joint entropy achieves this upper bound only if all of these n random variables are independent.

Another upper bound to note is that for any single random variable X which has N possible values, its entropy $H(X)$ is maximised when all of those values have the same probability $p_i = 1/N$. In that case,

$$H(X) = - \sum_i p_i \log_2 p_i = - \sum_1^N \frac{1}{N} \log_2 \frac{1}{N} = \log_2 N .$$

We shall use this fact when evaluating the efficiency of coding schemes.

Mutual Information between X and Y

The *mutual information* between two random variables measures the amount of information that one conveys about the other. Equivalently, it measures the average reduction in uncertainty about X that results from learning about Y . It is defined:

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

Clearly, X says as much about Y , as Y says about X . Note that in case X and Y are independent random variables, then the numerator inside the logarithm equals the denominator. Then the log term vanishes to 0, and so the mutual information equals zero, as one should expect when random variables are independent.

Mutual information can be related to entropies in three alternative ways, as is apparent from the Venn diagram, and it has important properties. We will see soon that it corresponds to the *relative entropy* between the joint distribution $p(x, y)$ and the product distribution $p(x)p(y)$.

Non-negativity: mutual information is always ≥ 0 . In the event that the random variables are perfectly correlated, then their mutual information is the entropy of either one alone. (Another way to say this is simply: $I(X; X) = H(X)$: the mutual information of a random variable with itself is its entropy. For this reason, the entropy $H(X)$ of a random variable X is sometimes referred to as its *self-information*.)

These properties are reflected in three equivalent definitions for the mutual information between random variables X and Y :

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

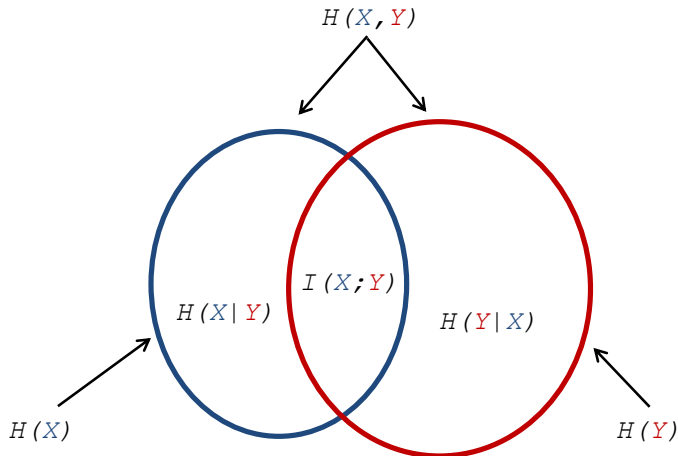
In a sense the mutual information $I(X; Y)$ is the intersection between $H(X)$ and $H(Y)$, since it represents their statistical dependence.

In the following Venn diagram, the portion of $H(X)$ that does not lie within $I(X; Y)$ is just $H(X|Y)$.

The portion of $H(Y)$ that does not lie within $I(X; Y)$ is just $H(Y|X)$.

Venn diagram summary of concepts and relationships

- ▶ Entropy $H(X)$, $H(Y)$
- ▶ Joint entropy $H(X, Y) = H(X) \cup H(Y)$
- ▶ Conditional entropy $H(X|Y)$, $H(Y|X)$
- ▶ Mutual information $I(X; Y) = H(X) \cap H(Y)$



“Distance” $D(X, Y)$ between two random variables X and Y

The amount by which the joint entropy of two random variables exceeds their mutual information is a measure of the “*distance*” between them:

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

Note that this quantity satisfies the standard axioms for a distance:

- ▶ $D(X, Y) \geq 0$ (distances are non-negative),
- ▶ $D(X, X) = 0$ (distance between something and itself is 0),
- ▶ $D(X, Y) = D(Y, X)$ (symmetry), and
- ▶ $D(X, Z) \leq D(X, Y) + D(Y, Z)$ (triangle inequality for distances).

Relative entropy, or Kullback-Leibler distance

An important measure of the “distance” between two random variables that does *not* satisfy the above axioms for a distance metric, is the

relative entropy or *Kullback-Leibler (KL) distance*. It is also called the *information for discrimination*, and it is used in pattern recognition. If $p(x)$ and $q(x)$ are two different probability distributions defined over the same set of outcomes x , then their relative entropy is:

$$D_{KL}(p\|q) = \sum_x p(x) \log \frac{p(x)}{q(x)}$$

Note that $D_{KL}(p\|q) \geq 0$, and in case $p(x) = q(x)$ then their distance $D_{KL}(p\|q) = 0$, as one might hope. However, this metric is not strictly a “distance,” since in general it lacks symmetry: $D_{KL}(p\|q) \neq D_{KL}(q\|p)$. Note also that major problems arise if there are any outcomes x for which $q(x)$ is vanishingly small relative to $p(x)$, thereby dominating the metric.

The relative entropy $D_{KL}(p\|q)$ is a measure of the “inefficiency” of assuming that a distribution is $q(x)$ when in fact it is $p(x)$.

If we have an optimal code for the distribution $p(x)$ (meaning that we use on average $H(p(x))$ bits, its entropy, to describe it), then the number of additional bits that we would need to use if we instead described $p(x)$ using an optimal code for $q(x)$, would be their relative entropy $D_{KL}(p\|q)$.

Fano's Inequality

We know that conditioning reduces entropy: $H(X|Y) \leq H(X)$. It is clear that if X and Y are perfectly correlated, then their conditional entropies $H(X|Y) = 0$ and $H(Y|X) = 0$. It should also be clear that if X is any deterministic function of Y , then again, there remains no uncertainty about X once Y is known, so their conditional entropy $H(X|Y) = 0$.

Fano's Inequality relates the probability of error P_e in guessing X from knowledge of Y to their conditional entropy $H(X|Y)$, when the number of possible outcomes is $|\mathcal{A}|$ (e.g. the length of a symbol alphabet):

$$P_e \geq \frac{H(X|Y) - 1}{\log |\mathcal{A}|}$$

The lower bound on P_e is a linearly increasing function of $H(X|Y)$.

The “Data Processing Inequality”

If random variables X , Y , and Z form a Markov chain (which means that the conditional distribution of Z depends only on Y and is independent of X), normally denoted as $X \rightarrow Y \rightarrow Z$, then the mutual information must be monotonically decreasing over steps along the chain:

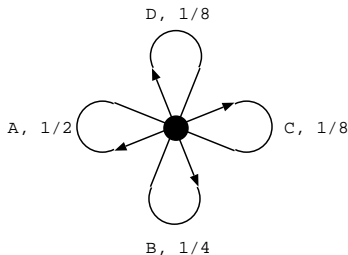
$$I(X; Y) \geq I(X; Z)$$

We turn next to applying these measures and relationships to the study of communication channels, symbol sources, codes, error correction, and channel capacity under various conditions.

3. Source coding theorem; variable-length and prefix codes

- ▶ Discrete symbol sources as Markov processes, with one or many states.
- ▶ Entropy of sources. Code rates and compression. Fixed-length codes.
- ▶ Capacity of a noiseless channel. Huffman codes and the prefix property.

We model a source of symbols as a **Markov process**, in which letters are emitted with known probabilities. Initially we consider just a one-state Markov process. (In a two-state Markov process, after the emission of certain symbols, the state may change to a different one having different emission probabilities for the symbols.)



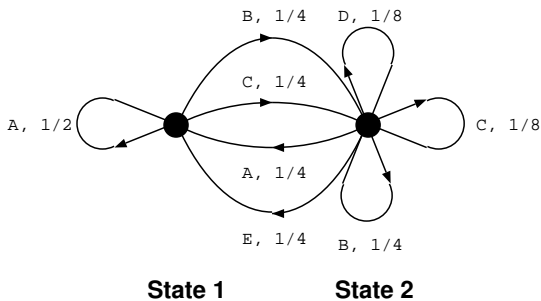
Such a Markov process (having any number of states) has an **entropy**. For the one-state Markov process above this is easily calculated from the probability distribution for letter emissions, and it is:

$$H = - \sum_i p_i \log p_i = \left(\frac{1}{2}\right)(1) + \left(\frac{1}{4}\right)(2) + \left(\frac{1}{8}\right)(3) + \left(\frac{1}{8}\right)(3) = 1.75 \text{ bits.}$$

Note that this entropy is based only on the symbol emission probabilities, regardless of the pace at which the source emits symbols, and so it is expressed as **bits per symbol**. If the symbols are emitted at a known rate, we may also characterise this symbol source in units of **bits per second**: (bits/symbol) \times (symbols/second).

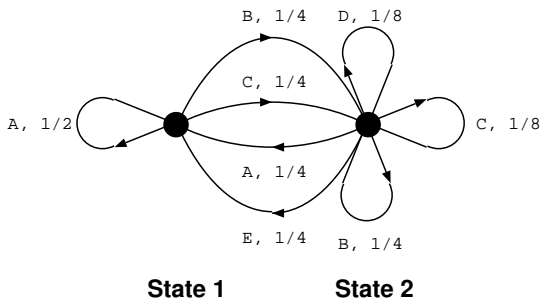
Now suppose that the Markov process has multiple states, and that transitions occur between states when certain symbols are emitted. For example, in English the letter 'q' is almost always followed by the letter 'u', but in other states (e.g. after 'z') the letter 'u' is far less probable.

First we can calculate the entropy associated with each state, as above; but then we also want to characterise the entropy of the entire process by taking into account the occupancy probabilities of these various states.



In the above two-state Markov process, which has “memory,” there are two different sets of probabilities for the emission of symbols.

If a Markov process has several states $\{S_1, S_2, \dots, S_n\}$, with associated emission probabilities $p_i(j)$ being the probability of emitting symbol j when in State S_i , then we first define the entropy of each of these states H_i in the normal manner, and then take the weighted average of those, using the occupancy probabilities P_i for the various states to arrive at an overall entropy $H = \sum_i P_i H_i$ for the multi-state Markov process.



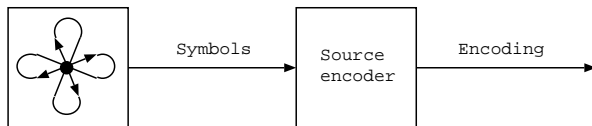
$$H = \sum_i P_i H_i = - \sum_i P_i \sum_j p_i(j) \log p_i(j)$$

In State 1, this system emits letters (B or C) with probability 0.5 that generate a transition to State 2. Similarly, with probability 0.5, State 2 emits letters (A or E) generating a transition back to State 1. Thus the state occupancies are equiprobable, $P_1 = P_2 = 0.5$, and it is easy to see that the overall entropy of this two-state Markov process is:

$$H = (0.5)(1.5) + (0.5)(2.25) = 1\frac{7}{8} = 1.875 \text{ bits.}$$

Fixed-length codes

We now consider various schemes for encoding symbols into codewords, as binary digits, with focus on the **average codeword length per symbol** when symbol probabilities are taken into account. We are interested in **data compression**, and the efficient use of **channel capacity**. We will also study the complexity of symbol decoding, and (later) the potential these schemes offer for **error correction**.



First consider encoding the set of N symbols $\{s_i\}$ having a probability distribution with entropy H , as a fixed-length (R) block of binary digits.

To ensure that the symbols can be decoded, we need a block of length $R = \log_2(N)$ if N is a power of 2, or otherwise $R = \lfloor \log_2(N) \rfloor + 1$ where $\lfloor X \rfloor$ is the largest integer less than X .

The **code rate** then is R bits per symbol, and as we noted earlier that entropy has an upper bound $H \leq \log_2(N)$, it follows that $H \leq R$.

The **efficiency** η of the coding is given by: $\eta = \frac{H}{R}$.

Note that fixed-length codes are inefficient for a number of reasons:

- ▶ If N is not a power of two, then much of the “address space” of the codewords is wasted. For example, if we have $N = 16$ symbols then 4 binary digits generate exactly the required number of codewords; but if $N = 17$ symbols then we need 5 bits, generating an address space of 32 codewords, of which 15 are wasted.
- ▶ If the N symbols occur with non-uniform probabilities, then even if N is a power of 2 the fixed-length code is still inefficient because we have $H < R$.
- ▶ In general, both of these sources of inefficiency for fixed-length codes will exist.

Variable-length codes, and those with the prefix property

In general symbols are not equiprobable, and we would hope to achieve some more compressed form of encoding by using variable-length codes, just as telegraphers used Morse code with short encodings for the more common letters and longer encodings for the less frequently used letters.

Consider a four-symbol alphabet and three possible variable-length codes:

x	$p(x)$	Code 1	Code 2	Code 3
A	$1/2$	1	0	0
B	$1/4$	00	10	01
C	$1/8$	01	110	011
D	$1/8$	10	111	111

The entropy of this alphabet is $H = (\frac{1}{2})(1) + (\frac{1}{4})(2) + (\frac{1}{8})(3) + (\frac{1}{8})(3) = 1.75$ bits. Note too that the **average codeword length in bits/symbol** (weighting each codeword length by its symbol's probability) for Code 2 and for Code 3 is also $R = 1.75$ bits; indeed the arithmetic is identical to that for entropy H . But for Code 1, the average codeword length is instead $R = (\frac{1}{2})(1) + (\frac{1}{4})(2) + (\frac{1}{8})(2) + (\frac{1}{8})(2) = 1.5$ bits/symbol.

Now let us examine some properties of each code in more detail.

x	$p(x)$	Code 1	Code 2	Code 3
A	1/2	1	0	0
B	1/4	00	10	01
C	1/8	01	110	011
D	1/8	10	111	111

Code 1 is not **uniquely decodable**: a bit sequence such as 1001 could be decoded either as ABA, or as DC (unless punctuation is added).

Code 2 is uniquely decodable, and **instantaneously**: once we have the bits for an encoded symbol we can decode immediately, without backtracking or waiting for more. This is called the **prefix** property: no codeword is the prefix of a longer codeword. For example, the bit string 0110111100110 is instantaneously decodable as ACDBAC.

Code 3 lacks the prefix property and it is not an instantaneous code: the codeword for B is also the start of the codeword for C; and the codeword for A cannot be decoded until more bits are received to exclude B or C.

Clearly, **self-punctuating** Code 2 is the most desirable of the three codes.

Shannon's Source-Coding Theorem

A remarkably far-reaching result, which was illustrated already in Code 2, is that: *it is possible to compress a stream of data whose entropy is H into a code whose rate R approaches H in the limit, but it is impossible to achieve a code rate $R < H$ without loss of information.*

Thus **the Source-Coding Theorem establishes limits to data compression**, and it provides operational meaning to entropy.

The usual statement of the theorem is that for a discrete source with entropy H , for any $\epsilon > 0$, it is possible to encode the symbols at an average rate R such that

$$R = H + \epsilon$$

as an asymptotic limit ($\epsilon \rightarrow 0$ as the number of symbols gets large).

Proof is given in Shannon and Weaver (1949). This is sometimes called the **Noiseless Coding Theorem** as it does not consider noise processes, such as bit corruption in a communication channel. Note the assumption that the **source statistics are known**, so H can be calculated.

Huffman codes

An optimal prefix code (having the shortest possible average codeword length in bits/symbol) for any given probability distribution of symbols can be constructed using an algorithm discovered by Huffman. This is a constructive illustration of Shannon's source-coding theorem.

The idea is to assign the bits in a reverse sequence corresponding to increasing symbol probability, so that the more probable symbols are encoded with shorter codewords. Thus we start with the two least frequent symbols and assign the "least significant bit" to them. More probable symbols will lack "less significant bits" in their codewords, which are therefore shorter, given the prefix property.

A **binary tree** is constructed in reverse, which specifies a hierarchical partitioning of the alphabet using a **priority queue** based (inversely) on probability.

Constructing a Huffman tree (example on next slide)

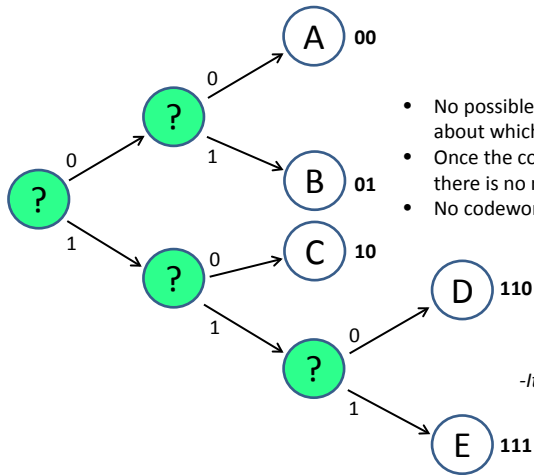
1. Find the two symbols having the lowest probabilities and assign a bit to distinguish them. This defines a branch in the binary tree.
2. Combine those two into a virtual “symbol” node whose probability is the sum of their two probabilities.
3. From this new shorter list of symbol nodes, repeat Step 1.
4. Repeat Step 2.
5. Continue this process until there is just one symbol node. That is the root node of the Huffman tree.

In the event that the symbol probabilities are powers of $1/2$, then a Huffman code achieves perfect efficiency ($R = H$). This is because each extra bit in a codeword removes half of the remaining uncertainty about the symbol, given that the bits so far have not specified it.

Note that there is not a unique Huffman code for any symbol alphabet. (The codewords having any given length could always be interchanged.) But a Huffman code is as efficient (compressive) as possible.

Example of a uniquely decodable, instantaneous, prefix code over 5 letters {A,B,C,D,E}

$p(A)=1/4$
 $p(B)=1/4$
 $p(C)=1/4$
 $p(D)=1/8$
 $p(E)=1/8$



- No possible received string of bits is ambiguous about which symbols were encoded.
- Once the codeword for any symbol is received, there is no need to wait for more bits to resolve it.
- No codeword is a prefix of another codeword.

How efficient is this code?
-It achieves optimal Shannon efficiency:

Note the entropy of this alphabet is:
 $3 \cdot (1/4) \cdot 2 + 2 \cdot (1/8) \cdot 3 = \underline{2.25 \text{ bits}}$

Note the average codeword length is also:
 $3 \cdot (1/4) \cdot 2 + 2 \cdot (1/8) \cdot 3 = \underline{2.25 \text{ bits/codeword}}$

Kraft-McMillan inequality

Any instantaneous code (one with the **prefix property**) must satisfy the following condition on the codeword lengths: if the N codewords have lengths $c_1 \leq c_2 \leq \dots \leq c_N$, then

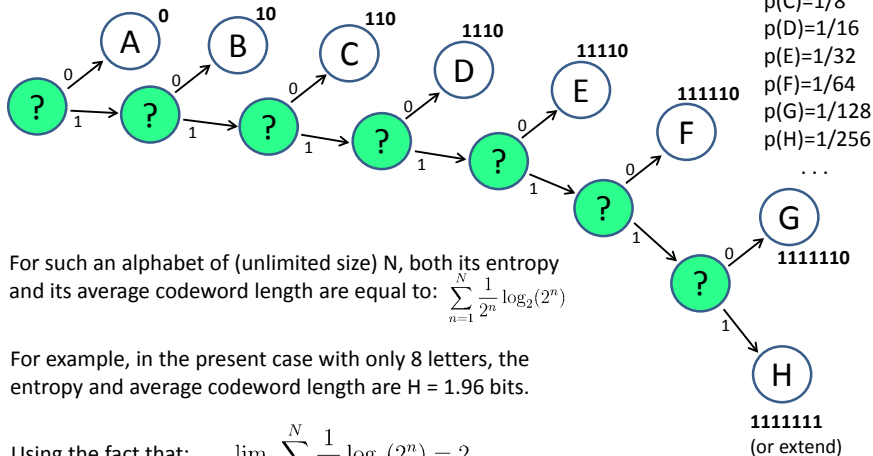
$$\sum_{i=1}^N \frac{1}{2^{c_i}} \leq 1.$$

Although this is a necessary condition, it is not a sufficient condition for a code to be an instantaneous code. For example, in our earlier study of three codes, we saw that Code 3 was not an instantaneous code but its codeword lengths were the same as those of Code 2, and both of them do satisfy the Kraft-McMillan inequality. (Note that the summation above equals 1.0 for both Codes 2 and 3, but it equals 1.25 for Code 1.)

Finally, we note an amazing consequence of the material covered in this section: It is possible to encode an alphabet of INFINITE length, provided its symbols have a particular probability distribution, with a prefix code whose average codeword length is no more than 2 bits per codeword!

Efficiency of prefix codes

Example of an alphabet of unlimited size, with a special probability distribution, that can be uniquely encoded with average codeword length $< \underline{2 \text{ bits/codeword}}!$



For such an alphabet of (unlimited size) N , both its entropy and its average codeword length are equal to: $\sum_{n=1}^N \frac{1}{2^n} \log_2(2^n)$

For example, in the present case with only 8 letters, the entropy and average codeword length are $H = 1.96$ bits.

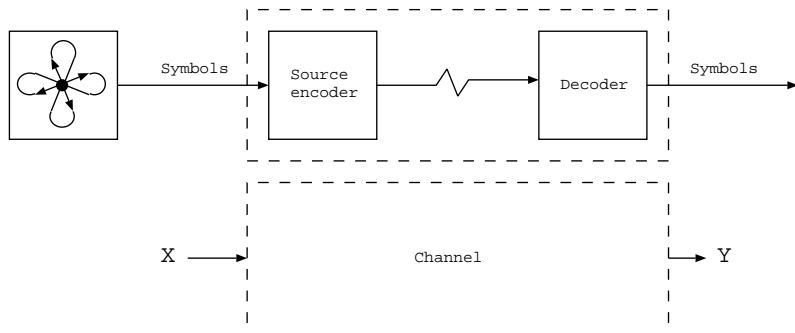
Using the fact that: $\lim_{N \rightarrow \infty} \sum_{n=1}^N \frac{1}{2^n} \log_2(2^n) = 2$

we see that even if the size of this alphabet grows indefinitely, it can still be uniquely encoded with an average codeword length just below 2 bits/codeword.

4. Discrete channel capacity, noise, and error correction

- ▶ *Channel matrix. Mutual information between input and output.*
- ▶ *Binary symmetric channel with error probability. Channel coding.*
- ▶ *Capacity of a noisy discrete channel. Error-correcting codes.*

We have considered discrete symbol sources, and encodings for them, and their code rates and compression limits. Now we consider channels through which such encodings pass, relating the input random variable X to the (perhaps randomly corrupted) output symbol, random variable Y .



Channel matrix

We shall apply all the tools and metrics introduced in the first part of this course. An input alphabet is random variable $X = \{x_1, \dots, x_J\}$, and the output symbol is drawn from random variable $Y = \{y_1, \dots, y_K\}$.

Note that J and K need not be the same. For example, for binary input $X = \{0, 1\}$ we could have an output alphabet $Y = \{0, 1, \perp\}$ where \perp means the decoder has detected some error.

A discrete memoryless channel can then be represented as a set of transition probabilities $p(y_k|x_j)$: that if symbol x_j is injected into the channel, symbol y_k is emitted. These conditional probabilities form the **channel matrix**:

$$\begin{pmatrix} p(y_1|x_1) & p(y_2|x_1) & \dots & p(y_K|x_1) \\ p(y_1|x_2) & p(y_2|x_2) & \dots & p(y_K|x_2) \\ \vdots & \vdots & \ddots & \vdots \\ p(y_1|x_J) & p(y_2|x_J) & \dots & p(y_K|x_J) \end{pmatrix}$$

For every input symbol we will get something out, so $\sum_{k=1}^K p(y_k|x_j) = 1$.

Average probability of symbol error, or correct reception

Using the channel matrix as well as some known probability distribution $\{p(x_j), j = 1, 2, \dots, J\}$ for a memoryless symbol source X , we can now apply the product and sum rules of probability to compute certain useful quantities. The joint probability distribution for the random variables X (input) and Y (output symbols) is: $p(x_j, y_k) = p(y_k|x_j)p(x_j)$, and the marginal probability distribution for output symbol y_k appearing is:

$$p(y_k) = \sum_{j=1}^J p(x_j, y_k) = \sum_{j=1}^J p(y_k|x_j)p(x_j)$$

Finally we can define the **average probability of symbol error** P_e as the sum of all elements in the channel matrix in which a different symbol was emitted than the one injected, which we will signify as $k \neq j$, weighted by the probability distribution $p(x_j)$ for the input symbols:

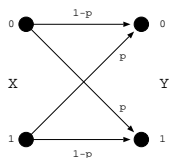
$$P_e = \sum_{j=1}^J \sum_{(k=1, k \neq j)}^K p(y_k|x_j)p(x_j)$$

and so the **average probability of correct reception** is: $1 - P_e$.

Binary symmetric channel

A binary symmetric channel has two input and two output symbols (denoted $\{0, 1\}$ for simplicity), and a common probability p of incorrect decoding of the input at the output. Thus its channel matrix is:

$$\begin{pmatrix} 1-p & p \\ p & 1-p \end{pmatrix} \text{ which we can graph as:}$$



Using the concepts and tools introduced at the beginning, we now wish to characterise the channel in terms of the **mutual information** between the input and output, and the **conditional entropy** $H(X|Y)$: how much uncertainty remains about the input X , given receipt of the output Y .

Let us assume that the two symbols of the input source $\{0, 1\}$ have probabilities $\{1/2, 1/2\}$ and so the source entropy is $H(X) = 1$ bit.

Note that the two output symbols also retain probabilities $\{1/2, 1/2\}$, independent of the error probability p , and so we also have $H(Y) = 1$ bit.

$$\begin{aligned}
 \text{The channel's conditional entropy } H(X|Y) &= - \sum_{x,y} p(x,y) \log p(x|y) \\
 &= -\frac{1}{2}(1-p) \log(1-p) - \frac{1}{2}p \log(p) - \frac{1}{2}p \log(p) - \frac{1}{2}(1-p) \log(1-p) \\
 &= -p \log(p) - (1-p) \log(1-p).
 \end{aligned}$$

Finally, the mutual information $I(X; Y)$ between the input and output is:

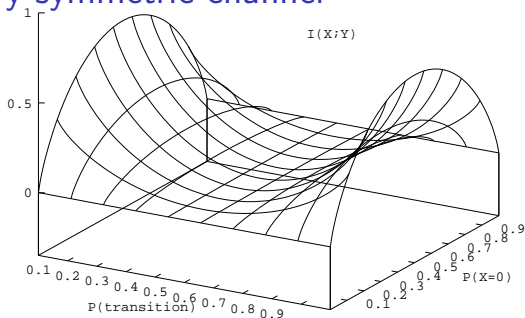
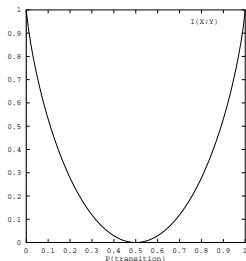
$$\begin{aligned}
 I(X; Y) &= H(X) - H(X|Y) \\
 &= 1 + p \log(p) + (1-p) \log(1-p).
 \end{aligned}$$

The **channel capacity**, denoted C , is defined as the maximum of its mutual information $I(X; Y)$ over all possible input source distributions.

$$C = \max_{\{p(x_j)\}} I(X; Y)$$

In this calculation we are assuming a “binary **symmetric** channel,” meaning that flips are equally likely for both input bits. Note also that having **equiprobable** $\{1/2, 1/2\}$ binary source symbols maximises $I(X; Y)$, which would be smaller had the input symbols not been equiprobable.

Capacity of the binary symmetric channel



The plot on the left shows $I(X; Y) = 1 + p \log(p) + (1 - p) \log(1 - p)$ as a function of the channel error (or transition) probability p , for the binary symmetric channel. Clearly $I(X; Y)$ is optimal in a channel with $p = 0$ or $p = 1$, meaning that a bit is never flipped or it is always flipped, leaving no uncertainty, and then the channel capacity is 1 bit per transmitted bit.

The surface on the right extends that same plot in a second dimension representing the variation of the two input probabilities. Over all the possible distributions $\{p(x_j)\}$ for the binary input, we see that channel capacity is always maximised for the equiprobable case $\{1/2, 1/2\}$, at a value that depends on the error probability p (again optimal if 0 or 1).

Schemes for acquiring immunity to channel noise

Adding various kinds of redundancy to messages can provide immunity to channel noise, including some remarkable cases of **error-correcting codes**. We begin with simple repetition codes, and with the observation that much redundancy already exists in natural language. The two sentences

1. “Bring reinforcements, we’re going to advance.”
2. “It’s easy to recognise speech.”

remain intelligible after the random corruption of 1 in 10 characters:

1. “Brizg reinforce ents, we’re goint to advance.”
2. “It’s easy mo recognis speech.”

But this intrinsic redundancy does not overcome the lack of orthogonality between messages. In audio terms, insufficient “distance” exists between the original two sentences above and the following spoken sentences:

1. “Bring three and fourpence, we’re going to a dance.”
2. “It’s easy to wreck a nice peach.”

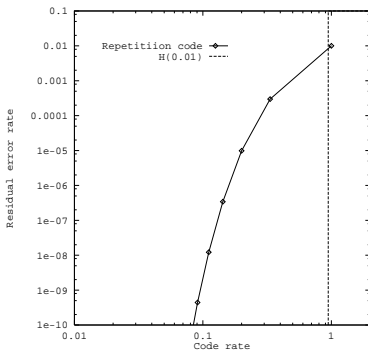
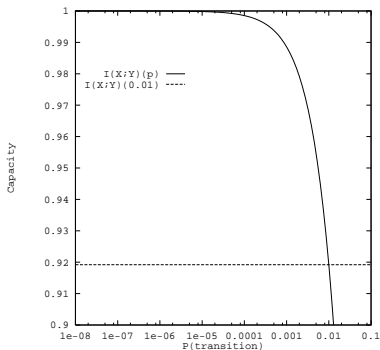
Repetition codes

A simple approach to overcoming channel noise might be just to repeat each message several times. Obviously the effective transmission rate is then diluted by a factor of N if there are N transmissions. We can analyze how much good this does in the case of the binary symmetric channel, with transition (error) probability p , which we saw for equiprobable inputs has a channel capacity $C = 1 + p \log(p) + (1 - p) \log(1 - p)$.

If we transmit every symbol an odd number $N = 2m + 1$ times and then perform majority voting, an error still persists if $m + 1$ or more bits are received in error. The probability P_e of this happening can be calculated with a binomial series. The summation is done over all modes of failed majority voting; the first factor is a combinatorial term (how many ways to choose i failed bits out of the $2m + 1$ bits), and the other factor is the probability that i of the bits were flipped while the remaining $2m + 1 - i$ bits were correctly transmitted:

$$P_e = \sum_{i=m+1}^{2m+1} \binom{2m+1}{i} p^i (1-p)^{2m+1-i}$$

Let us suppose the transition probability is $p = 0.01$, in which case the channel capacity is $C = 1 + p \log(p) + (1 - p) \log(1 - p) = 0.9192$ bit per bit. Capacity is plotted as a function of p in the left panel below.



The right panel plots the residual error probability P_e as a function of the code rate (channel capacity diluted by the number of transmissions).

It shows, for example, that if every transmission is repeated 7 times, giving a code rate of 0.13 , then for this channel with transition probability $p = 0.01$ there remains an error probability of about 1 in a million.

Channel Coding Theorem: error-correcting codes

We arrive at Shannon's second theorem, the **channel coding theorem**:

For a channel of capacity C and a symbol source of entropy H , provided that $H \leq C$, there exists a coding scheme such that the source is reliably transmitted through the channel with a residual error rate lower than any arbitrarily small ϵ .

Shannon's proof of this theorem is an existence proof rather than a means to construct such codes in the general case. In particular, the choice of a good code is dictated by the characteristics of the channel noise.

It is remarkable that it is possible to transmit data reliably through noisy channels, without using repetition. Methods of error-correcting encoding are pervasive not only in communications, but also within storage media. For example, a bubble or a scratch in CD or DVD media may obliterate many thousands of bits, but with no loss of data.

Today there is a vast literature around this subject. We will examine just one simple but efficient error-correcting code, with $H = C$.

A systematic (7/4) Hamming Code

Suppose symbols are encoded into a channel in blocks of seven bits. For example, it might be a block code for up to 128 ASCII characters. But the channel may randomly corrupt 1 bit in each block of 7 (or none).

Thus for each block of 7 bits $b_1 b_2 b_3 b_4 b_5 b_6 b_7$ encoding a particular input symbol x_j among the 128 in the alphabet, any one of 8 possible output symbols may emerge, all equiprobable (with probability $1/8$), having the following bit patterns:

$$\begin{aligned} &b_1 b_2 b_3 b_4 b_5 b_6 b_7 \\ &\bar{b}_1 b_2 b_3 b_4 b_5 b_6 b_7 \\ &b_1 \bar{b}_2 b_3 b_4 b_5 b_6 b_7 \\ &b_1 b_2 \bar{b}_3 b_4 b_5 b_6 b_7 \\ &b_1 b_2 b_3 \bar{b}_4 b_5 b_6 b_7 \\ &b_1 b_2 b_3 b_4 \bar{b}_5 b_6 b_7 \\ &b_1 b_2 b_3 b_4 b_5 \bar{b}_6 b_7 \\ &b_1 b_2 b_3 b_4 b_5 b_6 \bar{b}_7 \end{aligned}$$

where an overbar \bar{b}_i signifies that bit b_i has been flipped.

Let us calculate the information capacity of this channel per symbol, $C_S = \max_{\{p(x_j)\}} I(X; Y) = \max_{\{p(x_j)\}} (H(Y) - H(Y|X))$, allocated per bit C_b by dividing by 7. For convenience we observe that $H(Y) = 7$ because there are $N = 128 = 2^7$ equiprobable possible symbols, and also we observe that $p(y_k|x_j)$ is always $1/8$ because any one of 8 equiprobable possible output symbols y_k will emerge for any given input symbol x_j .

$$C_S = \max_{\{p(x_j)\}} (H(Y) - H(Y|X)) \text{ bits per symbol}$$

$$C_b = \frac{1}{7} \left(7 - \sum_j \sum_k p(y_k|x_j) \log \left(\frac{1}{p(y_k|x_j)} \right) p(x_j) \right) \text{ bits per bit}$$

$$= \frac{1}{7} \left(7 + \sum_j 8 \left(\frac{1}{8} \log \frac{1}{8} \right) \frac{1}{N} \right)$$

$$= \frac{1}{7} \left(7 + N \left(\frac{8}{8} \log \frac{1}{8} \right) \frac{1}{N} \right)$$

$$= \frac{4}{7} \text{ bits per bit}$$

Thus the information capacity of this channel is $\frac{4}{7}$ bit per bit encoded.

Syndromes

Can we develop an error-correcting code that reliably transmits data through this channel, if our source entropy $H \leq C = \frac{4}{7}$ bit per bit?

We construct new 7-bit codewords, each of which contains just 4 bits of symbol-encoding data, plus another 3 bits computed for error correction.

In our new codewords, bits b_3 , b_5 , b_6 , and b_7 are symbol-encoding data. But new bits b_4 , b_2 , and b_1 are computed from those 4 bits, as follows:

$$b_4 = b_5 \oplus b_6 \oplus b_7$$

$$b_2 = b_3 \oplus b_6 \oplus b_7$$

$$b_1 = b_3 \oplus b_5 \oplus b_7$$

Upon reception of these new 7-bit codewords (or corruptions of them), 3 further bits called **syndromes** s_4 , s_2 , and s_1 are then computed:

$$s_4 = b_4 \oplus b_5 \oplus b_6 \oplus b_7$$

$$s_2 = b_2 \oplus b_3 \oplus b_6 \oplus b_7$$

$$s_1 = b_1 \oplus b_3 \oplus b_5 \oplus b_7$$

If the syndromes computed upon reception are all 0, there was no error. Otherwise, the bit position $b_{s_4 s_2 s_1}$ is the bit in error.

Thus we can reliably transmit data through this noisy channel, whose capacity is $C = \frac{4}{7}$ bit per bit encoded, by embedding 3 error-correcting bits with every 4 useful data bits. This dilutes our source entropy to $H = \frac{4}{7}$ bit per bit of data, consistent with the requirement of Shannon's Channel Coding Theorem that $H \leq C$.

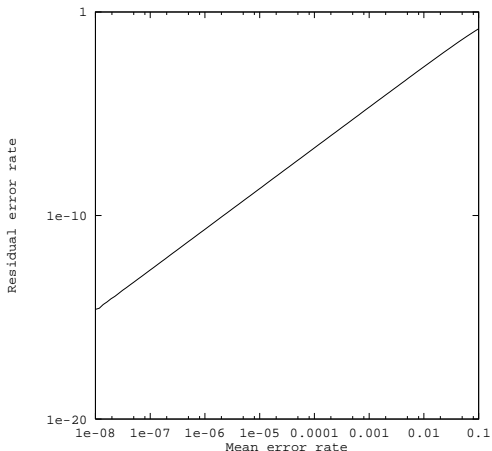
Hamming codes are called **perfect** because they use m bits to correct $2^m - 1$ error patterns (in this case 3 bits to correct 7 error patterns), and transmit $2^m - 1 - m$ (in this case 4) useful bits.

Finally we observe that in the more general case in which more than 1 bit might be corrupted in each block of 7 bits, then this scheme too will fail. The residual error probability P_e is again given by the remainder of a binomial series:

$$P_e = \sum_{i=2}^7 \binom{7}{i} p^i (1-p)^{7-i}$$

as $i = 2$ or more bits in each block of 7 are corrupted, and $7 - i$ are not.

The plot shows how P_e depends on unconstrained bit error probability p .



Many more robust block codes exist; for example the Golay code embeds 11 error-correcting bits into codeword blocks of 23 bits, enabling it to correct up to 3 bit errors in the block, and transmit 12 useful data bits.

5. Spectral properties of continuous-time channels

- ▶ *Continuous-time signals as bandlimited carriers of information.*
- ▶ *Signals represented as superpositions of complex exponentials.*
- ▶ *Eigenfunctions of channels modelled as linear time-invariant systems.*
- ▶ *Continuous-time channels as spectral filters with added noise.*

Both signals and the channels that transmit them are ultimately physical systems, with spectral and information properties determined by physics. Understanding these physical properties leads to important insights.

Information channels are typically an assigned **spectral band** of some medium, within which a **carrier signal** is modulated in certain of its parameters (sine frequency, amplitude, or phase) to encode information. The resulting complex $f(t)$ has a certain **bandwidth** Ω in Hertz, which is related to its information capacity C in bits/sec. By Fourier analysis we may regard $f(t)$ as a superposition of complex exponentials:

$$f(t) = \sum_n c_n e^{i\omega_n t}$$

We can understand the process of sending signals through channels in Fourier terms, because channels are (ideally) linear time-invariant systems whose **eigenfunctions** are in fact complex exponentials:

$$e^{(i\omega_n t)} \longrightarrow \boxed{h(t)} \longrightarrow \alpha e^{(i\omega_n t)}$$

Linear time-invariant systems (e.g. a coaxial cable, or the air through which spoken sounds pass, or the electromagnetic spectrum in space) obey the properties of superposition and proportionality, and can always be described by some **linear operator** $h(t)$. Examples of $h(t)$ may include: a derivative, or combination of them making a differential operator; or a convolution. The point is that a complex exponential is never changed in its form by being acted on by a linear operator; it is only multiplied by a complex constant, α . This is the **eigenfunction property**.

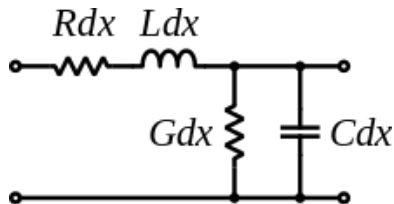
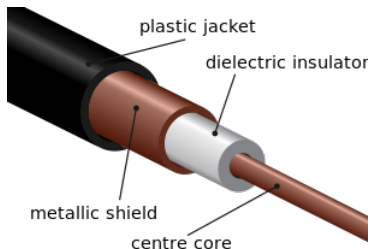
The consequence is that if we know the amplitude-and-phase changing values α_n for the relevant frequencies ω_n in a particular channel, then we simply incorporate them into the earlier Fourier series expansion of $f(t)$ in order to understand how $f(t)$ is spectrally affected by transmission through that channel:

$$f(t) = \sum_n \alpha_n c_n e^{(i\omega_n t)}$$

Bandwidth limitation of communication channels

A coaxial cable is one example of a communication channel. Its physical properties (distributed per unit length, hence the differential dx notation) include:

- ▶ a non-zero series resistance R ;
- ▶ a non-zero series inductance L ;
- ▶ a non-zero shunt conductance G (or non-infinite resistance) through the insulator separating signal from ground;
- ▶ and a non-zero shunt capacitance C between signal and ground.



This “equivalent circuit” for a physical communication channel makes it a **low-pass filter**: its ability to transmit signals is restricted to some finite bandwidth W , in Hertz. Frequency components higher than about $W \approx 1/(RC)$ are significantly attenuated by the signal pathway between conductor and shield (signal and ground). Although the attenuation is first-order and therefore gradual (6 dB per octave in terms of amplitude), still we must take into account this limited bandwidth W in Hertz.

How do you expect the band limitation W will influence the information carrying capacity of a channel?

Another physical limitation of communication channels is that they add some **wideband noise** to the signal during transmission. This is generally thermal noise, or **shot noise**, but it may come from various sources and generally has a broad spectral density across the channel's bandwidth W . If the noise spectrum is quite uniform, it is called **white noise**, in analogy with the spectral composition of white light. *Other colours are available.*

How do you expect the added noise, of a certain power, will influence the information carrying capacity of a channel?

6. Continuous information; Noisy Channel Coding Theorem

- ▶ *Extensions of discrete entropies and measures to continuous variables.*
- ▶ *Gaussian channels; signal-to-noise ratio; power spectral density.*
- ▶ *Relative significance of bandwidth and noise limits on channel capacity.*

We turn now to the encoding and transmission of information that is continuously variable in time or space, such as sound, or optical images, rather than discrete symbol sets. Using continuous probability densities, many of our metrics generalise from the discrete case in a natural way.

If the value X that a continuous signal may take (such as voltage, or sound pressure $x(t)$ as a function of time) has some probability density

$p(x)$ with $\int_{-\infty}^{+\infty} p(x)dx = 1$, then we define its **differential entropy** as:

$$h(X) = \int_{-\infty}^{+\infty} p(x) \log_2 \left(\frac{1}{p(x)} \right) dx .$$

(We use h for entropy of a continuous variable; H for discrete variables.)

Let $p(x, y)$ be the **joint probability distribution** of two continuous random variables X and Y . The marginal distribution of either one is obtained by integrating $p(x, y)$ over the other variable, just like the Sum Rule:

$$p(x) = \int_{-\infty}^{+\infty} p(x, y) dy$$

$$p(y) = \int_{-\infty}^{+\infty} p(x, y) dx .$$

The **joint entropy** $h(X, Y)$ of continuous random variables X and Y is

$$h(X, Y) = \int_{-\infty}^{+\infty} \int_{-\infty}^{+\infty} p(x, y) \log_2 \left(\frac{1}{p(x, y)} \right) dx dy$$

and their **conditional entropies** $h(X|Y)$ and $h(Y|X)$ are

$$h(X|Y) = \int_{-\infty}^{+\infty} \int_{-\infty}^{+\infty} p(x, y) \log_2 \left(\frac{p(y)}{p(x, y)} \right) dx dy$$

$$h(Y|X) = \int_{-\infty}^{+\infty} \int_{-\infty}^{+\infty} p(x, y) \log_2 \left(\frac{p(x)}{p(x, y)} \right) dx dy .$$

We also have the property that $h(X, Y) \leq h(X) + h(Y)$, with the upper bound reached in the case that X and Y are **independent**.

Finally, the **mutual information** $i(X; Y)$ between two continuous random variables X and Y is

$$i(X; Y) = \int_{-\infty}^{+\infty} \int_{-\infty}^{+\infty} p(x, y) \log_2 \left(\frac{p(x|y)}{p(x)} \right) dx dy = h(Y) - h(Y|X)$$

and as before, the **capacity** C of a continuous communication channel is determined by maximising this mutual information over all possible input distributions $p(x)$ for X , the signal entering the channel.

Thus we need to think about what class of continuous signals have “maximum entropy” for a given amplitude range of excursions, which corresponds essentially to the power level, or variance, of a signal.

We know that entropy is maximised with “equiprobable symbols,” which means in the continuous case that if signal value x is limited to some range v , then the probability density for x is uniformly $p(x) = 1/v$.

Gaussian noise maximises entropy for any given variance

For any continuous random variable X , the greater its variance (which corresponds to the power or volume level of a sound signal), the greater its differential entropy $h(X)$. But for any given power level or variance σ^2 of the signal's excursions around its mean value μ , it can be proven that the distribution $p(x)$ of those excursions which generates *maximum* $h(X)$ is the **Gaussian** distribution:

$$p(x) = \frac{1}{\sqrt{(2\pi)\sigma}} e^{-(x-\mu)^2/2\sigma^2}$$

in which case the differential entropy is maximised at

$$h(X) = \frac{1}{2} \log_2(2\pi e\sigma^2).$$

Such a signal sounds like “hiss” and it is called **white noise**, because its power spectrum is flat, just like the spectrum of white light.

The entropy-maximising property of Gaussian noise is important because we compute channel capacity as the mutual information $i(X; Y)$ of the channel when maximised over all possible input distributions $p(x)$.

Additive white Gaussian noise (AWGN) channel

We consider channels into which Gaussian noise (denoted N) is injected, independently of the signal, and added to it. The channel itself has a limited spectral bandwidth W in Hertz, and so both the signal and the added noise are strictly **lowpass**: they have no frequency components higher than W within the channel.

Over this bandwidth W , the white noise has **power spectral density** N_0 , which gives it a **total noise power** $N_0 W$ and variance $\sigma^2 = N_0 W$.

Because the noise N is independent of the input signal X , we have $h(Y|X) = h(N)$; and because the noise is Gaussian with $\sigma^2 = N_0 W$,

$$h(Y|X) = h(N) = \frac{1}{2} \log_2(2\pi e\sigma^2) = \frac{1}{2} \log_2(2\pi eN_0 W).$$

The input signal X itself has variance or power P over the spectrum W . As variances add, the channel output $Y = X + N$ has variance $P + N_0 W$.

Now we can ask the question:

What is the capacity C of such a channel, given P , N_0 , and W ?

Noisy Channel Coding Theorem

Maximising the mutual information of the channel over all possible input distributions requires this calculation to assume that the input X itself has a Gaussian $p(x)$. The mutual information $i(X; Y)$ between the input to the AWGN channel and the output signal Y transmitted is then:

$$\begin{aligned}i(X; Y) &= h(Y) - h(Y|X) \\ &= h(Y) - h(N) \\ &= \frac{1}{2} \log_2(2\pi e(P + N_0 W)) - \frac{1}{2} \log_2(2\pi e N_0 W) \\ &= \frac{1}{2} \log_2 \frac{2\pi e(P + N_0 W)}{2\pi e N_0 W} \\ &= \frac{1}{2} \log_2 \left(1 + \frac{P}{N_0 W} \right)\end{aligned}$$

which gives us the channel capacity C in bits per channel “symbol:”

$$C = \frac{1}{2} \log_2 \left(1 + \frac{P}{N_0 W} \right)$$

(Noisy Channel Coding Theorem, con't)

We will see from the **sampling theorem** that a strictly bandlimited signal whose lowpass bandwidth in Hertz is W (like the output of this channel) is completely specified by sampling it at a rate $2W$. Thus we can convert the channel capacity C from bits per “symbol” to bits per second, by multiplying the last expression for C by $2W$ symbols/sec:

$$C = W \log_2 \left(1 + \frac{P}{N_0 W} \right) \text{ bits/sec}$$

Note that the term inside the logarithm is $1 +$ the **signal-to-noise ratio**, often abbreviated SNR. Because of the logarithm, SNR is often reported in **decibels**, denoted dB: $10 \times \log_{10}(\text{SNR})$ if SNR is a ratio of power, or $20 \times \log_{10}(\text{SNR})$ if SNR is a ratio of amplitudes. Thus, for example, an amplitude signal-to-noise ratio of 100:1 corresponds to an SNR of 40 dB.

The critical insight about this theorem is that continuous-time channel capacity is dictated primarily by the signal-to-noise ratio!

(Noisy Channel Coding Theorem, con't)

Thus we have arrived at Shannon's third theorem (also called the Shannon-Hartley Theorem), the **Noisy Channel Coding Theorem**:

The capacity of a continuous-time channel, bandlimited to W Hertz, perturbed by additive white Gaussian noise of power spectral density N_0 and bandwidth W , using average transmitted power P , is:

$$C = W \log_2 \left(1 + \frac{P}{N_0 W} \right) \text{ bits/sec}$$

It is noteworthy that increasing the bandwidth W in Hertz yields a monotonic, but asymptotically limited, improvement in capacity because inside the logarithm its effect on the total noise power $N_0 W$ is reciprocal. Using the series: $\log_e(1+x) = x - \frac{x^2}{2} + \frac{x^3}{3} - \frac{x^4}{4} + \dots$, so $\log_e(1+x) \approx x$ for $x \ll 1$, we see that the limit as $W \rightarrow \infty$ is: $C \rightarrow \frac{P}{N_0} \log_2 e$.

But improving the channel's signal-to-noise ratio SNR improves the channel capacity without limit.

Effect of “coloured” noise on channel capacity

The preceding analysis assumed a constant noise power spectral density N_0 over the bandwidth W . But often the noise spectrum is non-uniform. For example, **pink noise** (also called $1/f$ or **flicker noise**) has a power spectral density inversely proportional to the frequency. Other named noise “colours” include red, blue, violet, and grey noise.

As SNR thus varies across the frequency band W , one would expect the information capacity also to be non-uniformly distributed. Taking an infinitesimal approach, the channel capacity C in any small portion $\Delta\omega$ around frequency ω where the signal-to-noise ratio follows $\text{SNR}(\omega)$, is:

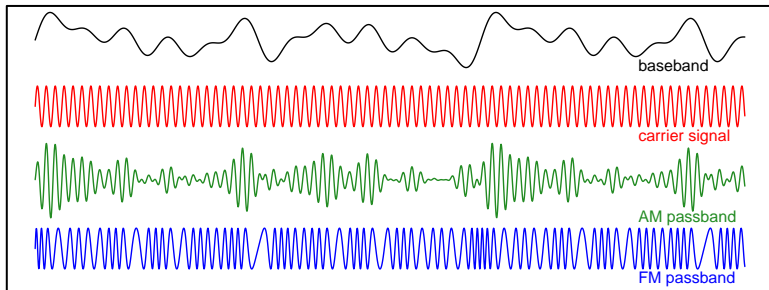
$$C_{(\Delta\omega)} = \Delta\omega \log_2(1 + \text{SNR}(\omega)) \text{ bits/sec.}$$

Integrating over all of these small $\Delta\omega$ bands in some available range from ω_1 to ω_2 , the information capacity of this variable-SNR channel is thus:

$$C = \int_{\omega_1}^{\omega_2} \log_2(1 + \text{SNR}(\omega)) d\omega \text{ bits/sec.}$$

7. Communications schemes exploiting Fourier theorems

Encoding and transmission of information in communication channels involves many schemes for the **modulation** of some parameters of a **carrier signal**, which is typically just a sinewave of some high frequency. Some familiar examples include amplitude modulation (**AM**), frequency modulation (**FM**), and phase modulation (**PM**, or **phase-shift keying**). Upon reception, the modulated carrier (*passband*) must be **demodulated** to recover the encoded information (the *baseband*). One reason for such schemes is that many different channels can share a common medium, such as a band of the electromagnetic spectrum. Selecting the carrier frequency is how one “tunes” into a given channel or mobile signal.



Modulation theorem

Conveying a message signal inside another signal that can be transmitted involves manipulations in the Fourier domain. Here we consider the basis of **amplitude modulation**, either SSB or DSB (single or double sideband).

Let $f(t)$ be the baseband message we wish to transmit; it might be the audio signal of somebody speaking. Let $F(\omega)$ be its Fourier transform:

$$F(\omega) = \mathcal{FT}\{f(t)\} = \frac{1}{2\pi} \int_{-\infty}^{\infty} f(t)e^{-i\omega t} dt$$

But what happens to $F(\omega)$ if we first multiply (modulate) $f(t)$ by a complex exponential **carrier signal** e^{ict} of frequency c ?

$$\begin{aligned}\mathcal{FT}\{e^{ict}f(t)\} &= \frac{1}{2\pi} \int_{-\infty}^{\infty} e^{ict}f(t)e^{-i\omega t} dt \\ &= \frac{1}{2\pi} \int_{-\infty}^{\infty} f(t)e^{-i(\omega-c)t} dt \\ &= F(\omega - c).\end{aligned}$$

Demodulation after single sideband modulation

We have seen that modulating a signal using a complex exponential of frequency c simply shifts its Fourier spectrum up by that frequency c , so we end up with $F(\omega - c)$ instead of $F(\omega)$. This enables the baseband signal $f(t)$ to be encoded into its own nominated slice of a shared broad communication spectrum, for transmission.

As c ranges into MHz or GHz, upon reception of this passband we must recover the original audio baseband signal $f(t)$ by shifting the spectrum back down again by the same amount, c . Clearly such **demodulation** can be achieved simply by multiplying the received signal by e^{-ict} , which is what a tuner does (*i.e.* the channel selector “dial” corresponds to c):

$$[e^{ict} f(t)] e^{-ict} = f(t).$$

But note both these operations of multiplying $f(t)$ by **complex-valued** functions e^{ict} and e^{-ict} are complicated: two modulating sinewaves are actually required, $\cos(ct)$ and $\sin(ct)$, in precise **quadrature phase**, and two circuits for multiplying them are needed. Both parts of the resulting complex-valued signal must be transmitted and detected.

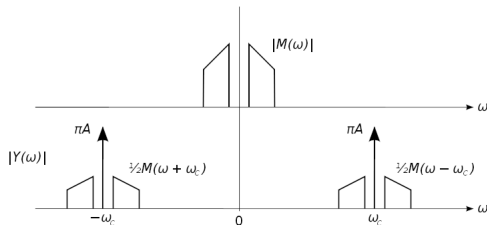
Double sideband amplitude modulation

To simplify the process and circuitry, but doubling the bandwidth needed, an alternative is just to multiply $f(t)$ by one (real-valued) cosine wave.

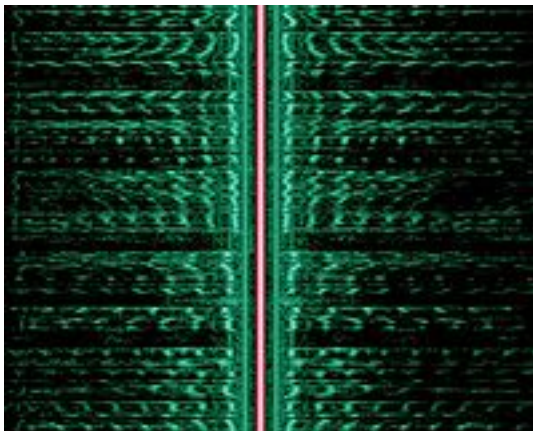
Since $\cos(ct) = \frac{e^{ict} + e^{-ict}}{2}$, by applying the Modulation theorem twice (adding the results) we see that if $F(\omega)$ is the Fourier transform of the original baseband $f(t)$, then after multiplying $f(t)$ by $\cos(ct)$, the new spectrum becomes:

$$\mathcal{FT}\{\cos(ct)f(t)\} = \frac{F(\omega - c) + F(\omega + c)}{2}.$$

Thus we end up with a passband whose spectrum consists of two copies of the baseband, shifted by both positive, and negative, frequencies c .



An actual passband broadcast using amplitude modulation



Spectrogram (frequency spectrum versus time) of an AM broadcast. Here, time is the vertical axis; frequency is the horizontal axis. The pink band of energy is the carrier itself (which can be suppressed). The two sidebands are displayed in green, on either side of the carrier frequency. Speech formants are visible.

Frequency Division Multiple Access (FDMA)

Radio waves propagate well through the atmosphere in a frequency range extending into GigaHertz, with specific bands allocated by government for commercial broadcasting, mobile phone operators, etc. The spectrum around 1 MegaHertz (0.3 to 3.0 MHz) is allocated for AM radio, and the UHF band around 1 GigaHertz (0.3 to 3.0 GHz) is for mobile phones, etc.

A human audio signal $f(t)$ occupies only about 1 KHz, so its spectrum $F(\omega)$ has a bandwidth that is tiny relative to the carrier frequency; thus a great many different mobile channels can be allocated by assignments of frequencies c . **Frequency Division Multiple Access (FDMA)** refers to one such regulated access protocol. (Other access protocols are available.)

Many alternative coding schemes exist, regulated by government agencies and by standards. **Phase modulation (PM)** illustrates another method for transmitting information. One natural application is in colour television, since perceived colour has a cyclical topology (the “colour wheel”), just like the phase ϕ of a carrier signal. PM resembles FM, because frequency ω is the time-derivative of phase: $\omega = \frac{d\phi}{dt}$. The detection of modulated functions as such can be decoded into symbols, as readily as colour.

8. The quantised degrees-of-freedom in a continuous signal

- ▶ *Nyquist sampling theorem; strict bandlimiting; aliasing and its prevention.*
- ▶ *Logan's theorem and the richness of zero-crossings in one-octave signals.*
- ▶ *The information diagram: how many independent quanta can it contain?*

Several independent results all illustrate the (perhaps surprising) fact that strictly **bandlimiting** a continuous function or signal causes it to have a finite, countable number of degrees-of-freedom. All the data contained in the continuous function can be regarded as **quantised**, into countable *quanta* of information, rather than having the density of real numbers.

- ▷ **Nyquist's sampling theorem** asserts that if a signal is **strictly bandlimited** to some highest frequency W , then simply sampling its values at a rate $2W$ specifies it completely everywhere, even *between* the sample points. Thus over some time interval T , it is fully determined by $2WT$ numbers.
- ▷ **Logan's theorem** asserts that if a function is bandlimited to one octave, then merely listing its **zero-crossings** fully determines it.
- ▷ **Gabor's information diagram** has a quantal structure, with a minimal area (bandwidth \times time) dictated by an **Uncertainty Principle** and occupied by Gaussian-attenuated complex exponentials: **Gabor wavelets**.

Nyquist's sampling theorem

Ideal sampling function

We define a **sampling function** $\text{comb}(t) = \delta_X(t)$ as an endless sequence of regularly spaced tines, separated by some **sampling interval** X :



Each “tine” is actually a Dirac δ -function, which can be regarded as the limit of a Gaussian whose width shrinks to 0. Multiplying a signal with one $\delta(t)$ samples its value at $t = 0$. The portrayed sequence of tines spaced by X is a sum of shifted δ -functions, making a sampling comb:

$$\delta_X(t) = \sum_n \delta(t - nX).$$

The sampling function $\delta_X(t)$ is **self-Fourier**: its Fourier transform $\Delta_X(\omega)$ is *also* a $\text{comb}(\omega)$ function, but with reciprocal interval $1/X$ in frequency:

$$\mathcal{FT}\{\delta_X(t)\} = \Delta_X(\omega) = \frac{1}{X} \sum_m \delta(\omega X - 2\pi m).$$



Properties of Dirac δ -functions

Conceptually,

$$\delta(t) = \begin{cases} \text{“}\infty\text{”} & t = 0 \\ 0 & t \neq 0 \end{cases}$$

and $\delta(t)$ contains unit area:

$$\int_{-\infty}^{\infty} \delta(t) dt = 1.$$

Multiplying any function $f(t)$ with a displaced δ -function $\delta(t - c)$ and integrating the product just picks out the value of $f(t)$ at $t = c$:

$$\int_{-\infty}^{\infty} f(t)\delta(t - c) dt = f(c)$$

which implies also that **convolving** any function with a δ -function simply reproduces the original function.

Having defined the “comb” sampling function as a sequence of displaced δ -functions with some sampling interval X , we can now understand the act of sampling a signal as just multiplying it with our $\text{comb}(t)$ function.

Strict bandlimiting, and spectral consequence of sampling

We now consider only signals $f(t)$ that have been **strictly bandlimited**, with some upper bound W on frequency. Thus their Fourier transforms $F(\omega)$ are 0 for all frequencies ω larger in absolute value than W :

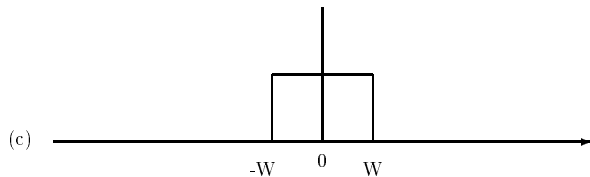
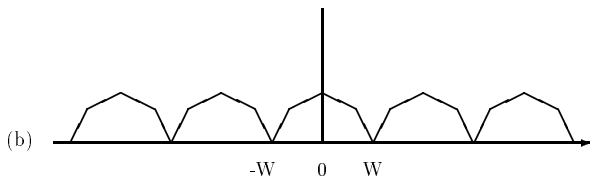
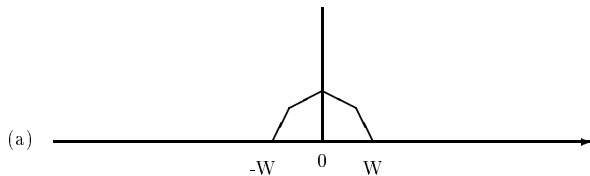
$$F(\omega) = 0 \quad \text{for } |\omega| > W.$$

If we have a signal $f(t)$ for which this is not already true, we make it so, by passing the signal through a strict **lowpass filter**, set for W . Its Fourier transform $F(\omega)$ becomes truncated, similar to trace (a) on the next slide.

Now if we sample $f(t)$ by multiplying it with our comb function $\delta_X(t)$, and use a **sampling rate** of at least $2W$ so that the sampling interval $X \leq 1/(2W)$, then we know from the **Convolution theorem** that the resulting sequence of samples will have a Fourier transform that is the convolution of $F(\omega)$ with $\Delta_X(\omega)$, the Fourier transform of $\delta_X(t)$.

Since convolution with a single δ -function reproduces function $F(\omega)$, and since $\Delta_X(\omega)$ is a sum of many shifted δ -functions $\sum_m \delta(\omega X - 2\pi m)$, our $F(\omega)$ becomes reproduced at every “tine” of $\Delta_X(\omega)$, as seen in trace (b).

(Nyquist's sampling theorem, continued)



Recovering the signal, even between its sampled points!

It is clear from trace (b) that the original spectrum $F(\omega)$ has completely survived the sampling process; but it just has been joined by multiple additional copies of itself. **As those did not overlap and superimpose**, we need only eliminate them (by another strict lowpass filtering action) in order to recover perfectly the Fourier transform $F(\omega)$ of our original (strictly bandlimited) signal, and thereby, that signal $f(t)$ itself. Visually, on the previous slide, multiplying trace(b) with trace(c) recovers trace(a).

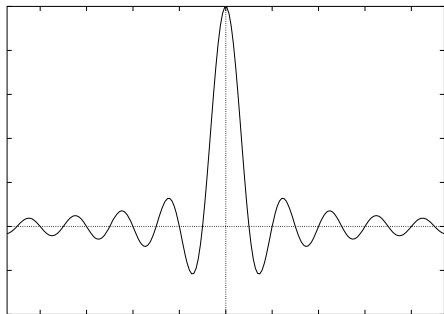
Why did those copied spectral lobes not overlap and superimpose?

Answer: we now see the critical importance of sampling the bandlimited signal $f(t)$ at a sampling rate, $2W$, at least twice as rapid as its highest frequency component W . Given the reciprocal spacing of the “tines” in $\delta_X(t)$ (namely $X \leq 1/(2W)$) and those in $\Delta_X(\omega)$ (namely $1/X \geq 2W$), it is clear that this **Nyquist sampling rate** $\geq 2W$ is the key constraint that ensures complete reconstructability of the signal $f(t)$ from just its discrete samples, even between the points where it was sampled.

But how can such complete reconstruction be achieved?

(Nyquist's sampling theorem, continued)

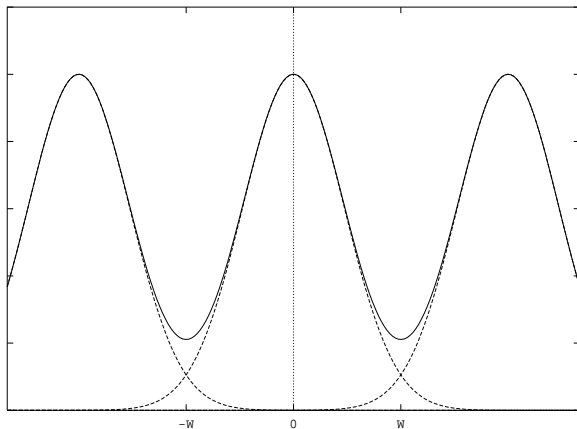
Again, we know from the Convolution theorem that strict lowpass filtering of the sequence $\delta_X(t)f(t)$ resulting from the sampling process equates to **convolving** this sample sequence by the inverse Fourier transform of that ideal lowpass filter. That is, of course, the **sinc** function: $\frac{\sin(2\pi Wt)}{2\pi Wt}$



Effectively, each sampled value get replaced by one of these, scaled in its (signed) amplitude by the sample value. These wiggly functions “fill in” all the space between the sample points, giving us back a **continuous** function $f(t)$; indeed exactly the one we started with before sampling it.

Aliasing

Failure to sample at a rate $\geq 2W$ causes **aliasing**: the added copies of the spectrum $F(\omega)$ overlap each other, because the tones in the frequency domain are now too close together. The extra lobes of $F(\omega)$ which they generate in the sampling process become superimposed, and now they can no longer be separated from each other by strict lowpass filtering.



(Nyquist's sampling theorem, con't): Effects of aliasing

This is the reason why “old western” movies with stagecoaches dashing away often seem to have their cartwheels rotating backwards. The film frame rate was only about 16 Hz, and in the interval between frames, each spoke in a cartwheel can be interpreted as moving backwards by a small amount, rather than forwards by the actual amount. Subtracting a small phase is equivalent to adding a larger phase to the cartwheel angle.

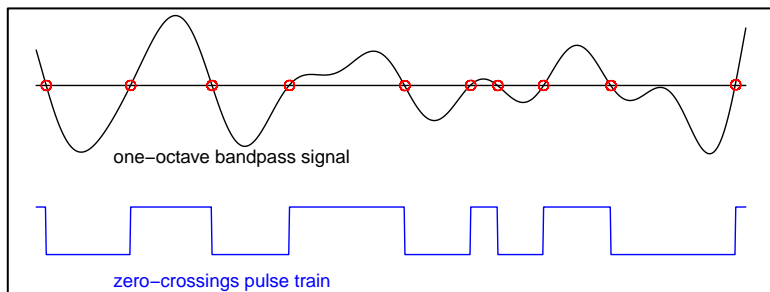


In terms of the previous slide showing superimposed spectral lobes after sampling at too low a frequency like $16 \text{ Hz} < 2W$, effectively the higher frequencies get “folded back” into the lower frequencies, producing an apparent, slow, reverse motion of the cartwheels.

Logan's theorem: reconstruction from zero-crossings alone

We saw in Nyquist's sampling theorem that strictly bandlimiting a signal (in that case to W) means that all the information it contains can be represented by a countable number of packets, or quanta of information (namely $2WT$ quanta over some time interval T).

Another result in the same spirit is **Logan's theorem**, which says that if a signal is **bandpass** limited to one-octave, so that the highest frequency component W_H in it is no higher than twice W_L its lowest frequency component, $W_H \leq 2W_L$, then merely listing the **zero-crossings** of this signal suffice for exact reconstruction of the signal (up to a scale factor in amplitude), even in the continuous spaces between those zero-crossings.



Does Logan's theorem explain how cartoons work?

It has been suggested that Logan's theorem might explain how cartoons work: cartoon sketches are highly impoverished images, with mere edges being drawn, and no real resemblance to the data contents of a picture. Perhaps human vision imparts so much richness to edges because they are the zero-crossings in bandpass filtered images, and retinal encoding is usually modeled in such terms (albeit with > 1.3 octave bandwidth).

The question can be posed in more abstract form, too. "Attneave's cat" is even more impoverished than Margaret Thatcher's caricature, because it consists only of some vertices joined by straight lines (no curvature). Yet like the former PM's sketch, it is instantly recognised as intended.



Limitations to Logan's theorem

After Logan proved his remarkable theorem at Bell Labs in 1977, there was enthusiasm to exploit it for extreme compression of speech signals. The zero-crossings are where the MSB (most significant bit) changes, as the waveform is bandpass; and indeed a pulse train of “1-bit speech” is remarkably intelligible. But the distortion is audible, and objectionable.

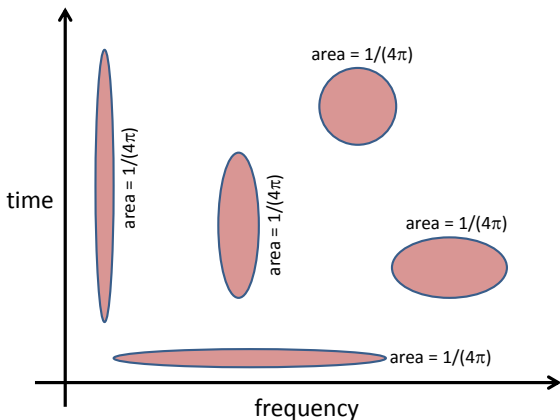
Clearly signals that are purely positive amplitude-modulations of a carrier $c(t)$, of the form $f(t) = [1 + a(t)]c(t)$ where $[1 + a(t)] > 0$, will encode nothing about $a(t)$ in the zero-crossings of $f(t)$: they will remain always just the zero-crossings of the carrier $c(t)$. Thus, “pure AM” signals are excluded from Logan's theorem by a condition that the signal must have **no complex zeroes in common with its Hilbert transform**.

Another limitation emerged in efforts to extend Logan's theorem to computer vision. The zero-crossings of bandpass filtered *images* are not discrete and countable, as for a 1D signal, but form continuous “snakes.”

Finally, numerical algorithms for actually reconstructing 1D signals from their zero-crossings are **unstable**: tiny changes in the position of an extracted zero-crossing have divergent effects on the reconstruction.

Quantal structure of Gabor's "information diagram"

We will see that a fundamental **uncertainty principle** limits simultaneous **localisability** of a signal in time and frequency. The shorter its duration, the broader its bandwidth becomes. The narrower its bandwidth, the longer it must persist in time. Thus the **information diagram** whose axes are time and frequency has a deeply **quantal structure**: no function can occupy an "area" smaller than an irreducible quantum, $1/(4\pi)$.



9. Gabor-Heisenberg-Weyl uncertainty principle; “logons”

We define the “effective width” (Δx) of a (complex-valued) function $f(x)$ in terms of its normalized variance, or normalized second-moment:

$$(\Delta x)^2 = \frac{\int_{-\infty}^{+\infty} f(x)f^*(x)(x - \mu)^2 dx}{\int_{-\infty}^{+\infty} f(x)f^*(x) dx}$$

where the purpose of the denominator is to normalise the amplitude or power in $f(x)$ so that we really just measure its width, and where μ is the mean value, or normalized first-moment, of the function $\|f(x)\|$:

$$\mu = \frac{\int_{-\infty}^{+\infty} xf(x)f^*(x) dx}{\int_{-\infty}^{+\infty} f(x)f^*(x) dx} .$$

Then, if we perform the same operations to measure the effective width ($\Delta\omega$) of the Fourier transform $F(\omega) = \mathcal{FT}\{f(x)\}$ of the function $f(x)$:

(Uncertainty principle, continued)

$$(\Delta\omega)^2 = \frac{\int_{-\infty}^{+\infty} F(\omega)F^*(\omega)(\omega - \nu)^2 d\omega}{\int_{-\infty}^{+\infty} F(\omega)F^*(\omega)d\omega}$$

where ν is the mean value, or normalized first-moment, of $\|F(\omega)\|$:

$$\nu = \frac{\int_{-\infty}^{+\infty} \omega F(\omega)F^*(\omega)d\omega}{\int_{-\infty}^{+\infty} F(\omega)F^*(\omega)d\omega}$$

then it can be proven (by Schwartz Inequality arguments) that there exists a **fundamental lower bound** on the product of these two “spreads,” *regardless* of the function $f(x)$:

$$\boxed{(\Delta x)(\Delta\omega) \geq \frac{1}{4\pi}}$$

This is the famous Gabor-Heisenberg-Weyl **Uncertainty Principle**.

(Uncertainty principle, continued)

Mathematically this is exactly equivalent to the uncertainty principle in quantum physics, where (x) would be interpreted as the position of an electron or other particle, and (ω) would be interpreted as its momentum or deBroglie wavelength⁻¹ (because of “wave-particle duality”).

$$(\Delta x)(\Delta \omega) \geq \frac{1}{4\pi}$$

We now see that this **irreducible joint uncertainty** is not just a law of nature, but it is a property of all functions and their Fourier transforms. It is thus another respect in which the information in continuous signals is quantised, since there is a limit to how sharply resolved they can be in the Information Diagram (e.g., their **joint resolution in time and frequency**).

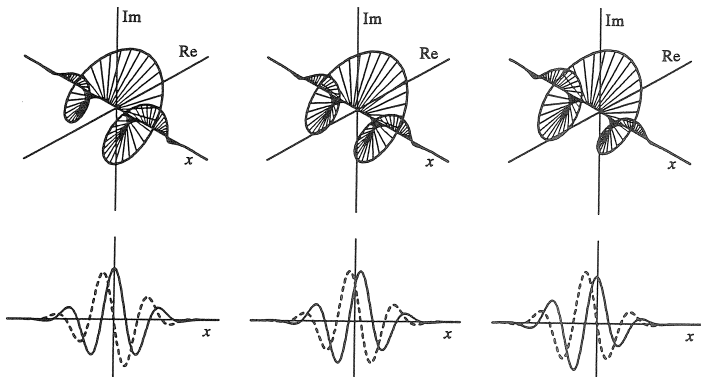
Dennis Gabor named such minimal areas “logons” from the Greek word for information, or order: *lōgos*. Their total number within a region of the Information Diagram specifies the maximum number of independent degrees-of-freedom available to a signal. Gabor went on to discover what family of functions actually achieve this minimal joint uncertainty.

Gabor wavelets, or “logons”

The unique family of functions that actually achieve the lower bound in joint uncertainty are the complex exponentials multiplied by Gaussians. Originally “logons”, today these are usually called **Gabor wavelets**:

$$f(x) = e^{-(x-x_0)^2/\alpha^2} e^{i\omega_0(x-x_0)}$$

localised at “epoch” x_0 , modulated with frequency ω_0 , and having a size or spread constant α . As **phasors** they are helical functions of x , and the lower trace plots the real and imaginary parts for different epochs x_0 .



(Gabor wavelets, continued)

Such wavelets are **self-Fourier**: their Fourier transforms $F(\omega)$ have the same functional form, but with the parameters just interchanged:

$$F(\omega) = e^{-(\omega-\omega_0)^2\alpha^2} e^{-ix_0(\omega-\omega_0)}$$

Note that for a wavelet whose epoch is $x_0 = 0$, its Fourier transform is just a Gaussian located at the modulation frequency ω_0 , and whose size is $1/\alpha$, the reciprocal of the centred wavelet's size or spread constant.

Because such wavelets have the greatest possible simultaneous resolution in time and frequency, Gabor proposed using them as an expansion basis to represent signals. Unfortunately, because these wavelets are mutually **non-orthogonal**, it is difficult to compute the expansion coefficients.

In the wavelets plotted (real parts only) in left column of the next slide, the Gaussian is always the same size α , while the frequency ω_0 increases. This reduces the wavelets' bandwidth in octave terms. It is more usual to use a **self-similar** scaling rule (illustrated in the right column), in which α and ω_0 are inversely proportional, so the wavelets are all dilated copies of a single **mother wavelet**.

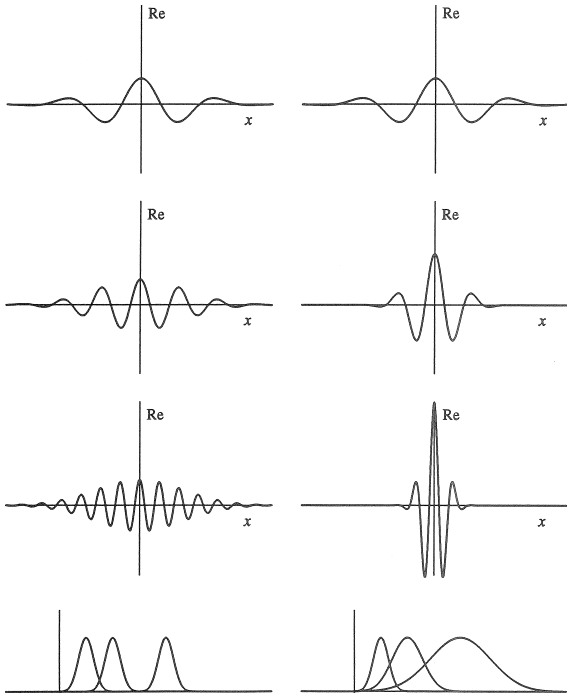
(Gabor wavelets, con't)



Left column: wavelets sharing a constant Gaussian size α but with increasing ω_0 .

Right column: self-similar wavelets with $\alpha^{-1} \sim \omega_0$.

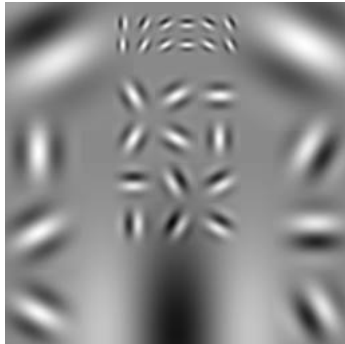
Bottom row: all of their Fourier transforms.



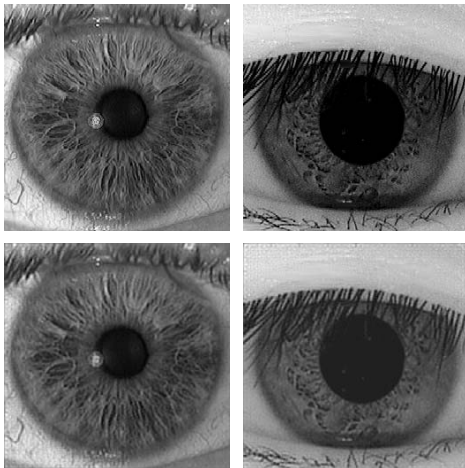
2D Gabor wavelets as encoders in computer vision

2D Gabor wavelets (defined as a complex exponential plane-wave times a 2D Gaussian windowing function) are extensively used in computer vision.

As multi-scale image encoders, and as pattern detectors, they form a complete basis that can extract image structure with a vocabulary of: location, scale, spatial frequency, orientation, and phase (or symmetry). This collage shows a 4-octave ensemble of such wavelets, differing in size (or spatial frequency) by factors of two, having five sizes, six orientations, and two quadrature phases (even/odd), over a lattice of spatial positions.



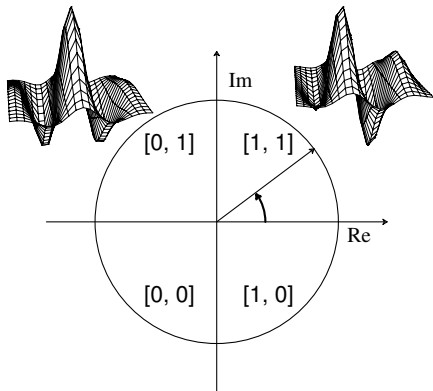
Complex natural patterns are very well represented in such terms. **Phase information** is especially useful to encode for pattern recognition e.g. iris. The upper panels show two iris images (acquired in near-infrared light); caucasian iris on the left, and oriental iris on the right.



The lower panels show the images reconstructed just from combinations of the 2D Gabor wavelets spanning 4 octaves seen in the previous slide.

Gabor wavelets are the basis for Iris Recognition systems

Phase-Quadrant Demodulation Code

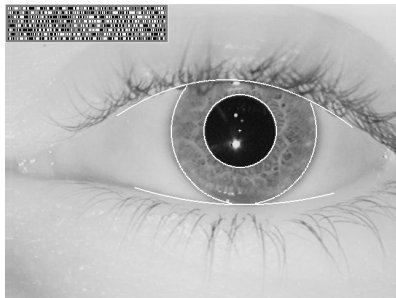


$$h_{Re} = 1 \text{ if } \text{Re} \int_{\rho} \int_{\phi} e^{-i\omega(\theta_0 - \phi)} e^{-(r_0 - \rho)^2 / \alpha^2} e^{-(\theta_0 - \phi)^2 / \beta^2} I(\rho, \phi) \rho d\rho d\phi \geq 0$$

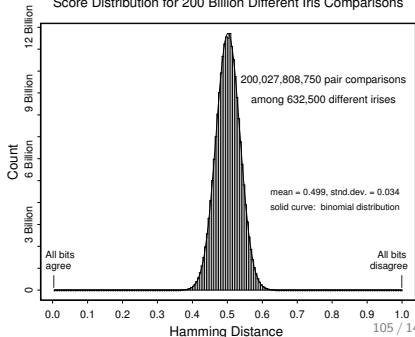
$$h_{Re} = 0 \text{ if } \text{Re} \int_{\rho} \int_{\phi} e^{-i\omega(\theta_0 - \phi)} e^{-(r_0 - \rho)^2 / \alpha^2} e^{-(\theta_0 - \phi)^2 / \beta^2} I(\rho, \phi) \rho d\rho d\phi < 0$$

$$h_{Im} = 1 \text{ if } \text{Im} \int_{\rho} \int_{\phi} e^{-i\omega(\theta_0 - \phi)} e^{-(r_0 - \rho)^2 / \alpha^2} e^{-(\theta_0 - \phi)^2 / \beta^2} I(\rho, \phi) \rho d\rho d\phi \geq 0$$

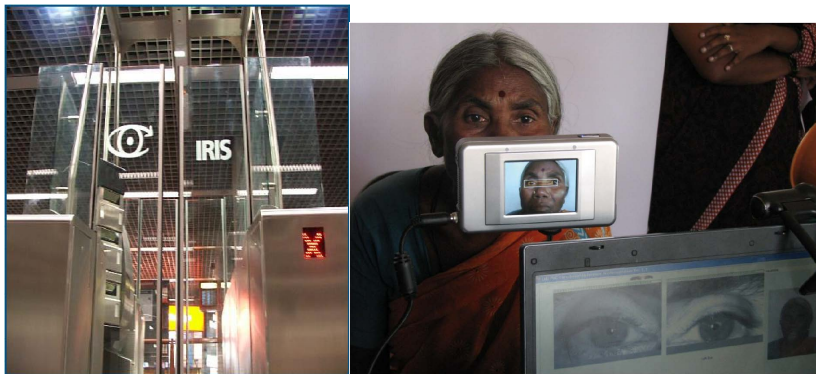
$$h_{Im} = 0 \text{ if } \text{Im} \int_{\rho} \int_{\phi} e^{-i\omega(\theta_0 - \phi)} e^{-(r_0 - \rho)^2 / \alpha^2} e^{-(\theta_0 - \phi)^2 / \beta^2} I(\rho, \phi) \rho d\rho d\phi < 0$$



Score Distribution for 200 Billion Different Iris Comparisons



Gabor wavelets are widely used for information encoding



At many airports worldwide, the **IRIS** system (Iris Recognition Immigration System) allows registered travellers to cross borders without having to present their passports, or make any other claim of identity. They just look at an iris camera, and (if they are already enrolled), the border barrier opens within seconds. Similar systems are in place for many other applications. The Government of India is currently enrolling the iris patterns of all its 1.2 Billion citizens as a means to access entitlements and benefits (the UIDAI slogan is “To give the poor an identity”), and to enhance social inclusion.

10. Data compression codes and protocols

- ▶ *Run-length encoding; predictive coding; dictionaries; vector quantisation.*
- ▶ *Image compression; Fast Fourier Transform algorithm; JPEG; wavelets.*

In general, a basic perspective from information theory is that where there is **redundancy**, there is opportunity for compression. The limit to compressibility is a representation in which no redundancy remains.

Run-length encoding (RLE) is an obvious example: if part of a data string contains (say) 1023 zeroes, one alternative would be to list them, in 1023 character bytes. An RLE code would instead summarise this substring more compactly as: “1023 × 0”.

Predictive coding refers to a broad family of methods in which just the **deviations** from a prediction are encoded, rather than the actual sample values themselves. Today’s temperature (or equity index value, etc) is usually a good prediction about tomorrow’s, and so it can be more efficient to encode the data as a string of (often small) Δ variations.

Dictionary-based compression such as various lossless schemes for text, images, and audio originated by **Lempel** and **Ziv**, exploit the fact that **strings of symbols** have probabilities that vary much more than the probabilities of the individual symbols. **Sparseness** is exploitable.

Dictionary methods

Practical dictionary methods (LZW, gif, etc) first construct a dictionary by scanning through the data file and adding an entry for every new word (or byte) encountered. Because of repeated appearances, the dictionary is much smaller than the data file. Its words are sorted by frequency. For example, the word “the” is 7% of all written and spoken English.

An index is constructed, and then the data file is scanned a second time, replacing each word with a pointer to its position in the index.

Note that the most commonly occurring words will have the shortest indices, because the dictionary was sorted in order of descending word frequency. Decoding reads the words from the dictionary according to the index pointers. (Note that the dictionary itself is part of the payload.)

Adaptive variants of such dictionary methods monitor the compression level achieved and make an assessment about how effective the dictionary is being. A dictionary may be abandoned when it ceases to be a compact compendium for the source, and construction of a new one begun.

Vector quantisation

A lossy method that also uses dictionary look-up, for compressing strings (vectors) by finding good approximations to them, is **vector quantisation**. The key insight is that the space of possible combinations of samples, or symbols, is only **very sparsely populated** with those that actually occur.

For example, suppose we have a “coding budget” of only 15 bits that we wish to use to encode English words. If we encode them in terms of their component letters, then naïvely we might spend 5 bits on each letter, and find that we only get as far as all possible 3-letter strings, most of which are nonsense.

On the other hand, if we use the 15 bits to build **pointers** to entries in a **codebook** or dictionary of actual English words, then our address space spans $2^{15} = 32,768$ actual words. This is much more than most people use or even know. (The average US university graduate has a vocabulary of about 10,000 words.)

(Vector quantisation, continued)

This idea generalises to encoding (for example) image structure using a **codebook of pixel combinations** that can be the basis of making at least good approximations to the kinds of local image structure (edges, etc) that actually occur. For example, instead of encoding (4×4) pixel tiles as 16 pixels, consuming 16 bytes = 128 bits, a codebook with an address space of $2^{128} \approx 10^{39}$ possible **image tiles** could be “constructed.” That is more than a *trillion-trillion-trillion* possible image tiles; ...surely enough...

Taking this codebook concept to an absurd length: across all of human history, about 10^{10} humans have lived, each for about 10^9 seconds, with visual experiences resolvable in time (when their eyes were open) at a rate of about 10 “frames”/second. Thus, our collective human history of distinguishable visual experiences could be encoded into codebook having about 10^{20} different image frames as entries. Such an address space is spanned by pointers having a length of just 66 bits.

Compare 66 bits with the length of a single SMS text message...

Obviously, **VQ** methods pose certain problems in terms of the codebook memory space requirements.

JPEG compression; Fast Fourier Transform algorithm

The Fast Fourier Transform (of which there are several variants) exploits some clever efficiencies in computing a discrete Fourier transform (DFT).

Since the explicit definition of each Fourier coefficient in the DFT is

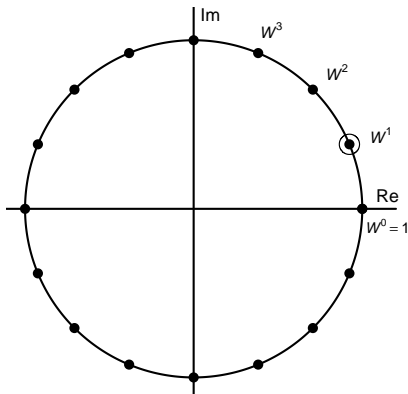
$$\begin{aligned} F[k] &= \sum_{n=0}^{N-1} f[n]e^{-2\pi ink/N} \\ &= f[0] + f[1]e^{-2\pi ik/N} + \dots + f[N-1]e^{-2\pi ik(N-1)/N} \end{aligned}$$

we can see that in order to compute one Fourier coefficient $F[k]$, using the complex exponential having frequency k , we need to do N (complex) multiplications and N (complex) additions. To compute all the N such Fourier coefficients $F[k]$ in this way for $k = 0, 1, 2, \dots, N-1$ would thus require $2N^2$ such operations. Since the number N of samples in a typical audio signal (or pixels in an image) whose DFT we may wish to compute may be $\mathcal{O}(10^6)$, clearly it would be very cumbersome to have to perform $\mathcal{O}(N^2) = \mathcal{O}(10^{12})$ multiplications. Fortunately, very efficient **Fast Fourier Transform (FFT)** algorithms exist that instead require only $\mathcal{O}(N \log_2 N)$ such operations, vastly fewer than $\mathcal{O}(N^2)$ if N is large.

(Fast Fourier Transform algorithm, review of IB, con't)

Recall that all the multiplications required in the DFT involve the N^{th} roots of unity, and that these in turn can all be expressed as powers of the primitive N^{th} root of unity: $e^{2\pi i/N}$.

Let us make that explicit now by defining this constant as $W = e^{2\pi i/N}$ (which is just a complex number that depends only on the data length N which is presumed to be a power of 2), and let us use W to express all the other complex exponential values needed, as the $(nk)^{\text{th}}$ powers of W : $e^{2\pi ink/N} = W^{nk}$.



(Fast Fourier Transform algorithm, review of IB, con't)

Or going around the unit circle in the opposite direction, we may write:

$$e^{-2\pi ink/N} = W^{-nk}$$

The same N points on the unit circle in the complex plane are used again and again, regardless of which Fourier coefficient $F[k]$ we are computing using frequency k , since the different frequencies are implemented by skipping points as we hop around the unit circle.

Thus the lowest frequency $k = 1$ uses all N roots of unity and goes around the circle just once, multiplying them with the successive data points in the sequence $f[n]$. The second frequency $k = 2$ uses every second point and goes around the circle twice for the N data points; the third frequency $k = 3$ hops to every third point and goes around the circle three times; etc.

Because the hops keep landing on points around the unit circle from the same set of N complex numbers, and the set of data points from the sequence $f[n]$ are being multiplied repeatedly by these same numbers for computing the various Fourier coefficients $F[k]$, it is possible to exploit some clever arithmetic tricks and an efficient recursion.

(Fast Fourier Transform algorithm, review of IB, con't)

Let us re-write the expression for Fourier coefficients $F[k]$ now in terms of powers of W , and divide the series into its first half plus second half. (“Decimation in frequency;” there is a “decimation in time” variant.)

$$\begin{aligned} F[k] &= \sum_{n=0}^{N-1} f[n] e^{-2\pi i n k / N} = \sum_{n=0}^{N-1} f[n] W^{-nk} \\ &= \sum_{n=0}^{N/2-1} f[n] W^{-nk} + \sum_{n=N/2}^{N-1} f[n] W^{-nk} \\ &= \sum_{n=0}^{N/2-1} (f[n] + W^{-kN/2} f[n + N/2]) W^{-kn} \\ &= \sum_{n=0}^{N/2-1} (f[n] + (-1)^k f[n + N/2]) W^{-kn} \end{aligned}$$

where the last two steps exploit the fact that advancing halfway through the cycle(s) of a complex exponential just multiplies value by $+1$ or -1 , depending on the parity of the frequency k , since $W^{-N/2} = -1$.

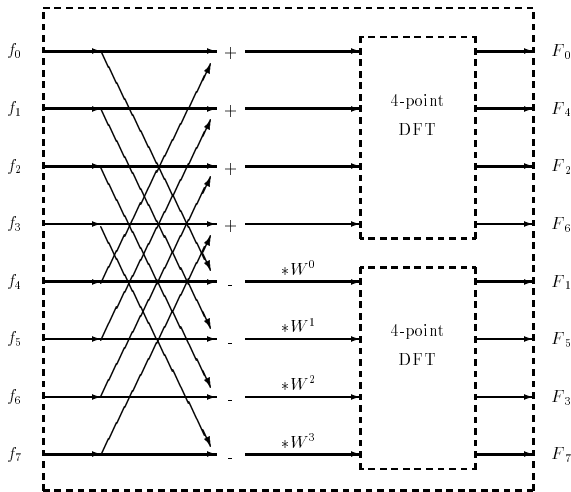
(Fast Fourier Transform algorithm, review of IB, con't)

Now, separating out even and odd terms of $F[k]$ we get $F_e[k]$ and $F_o[k]$:

$$F_e[k] = \sum_{n=0}^{N/2-1} (f[n] + f[n + N/2])W^{-2kn}, k = 0, 1, \dots, N/2 - 1$$

$$F_o[k] = \sum_{n=0}^{N/2-1} (f[n] - f[n + N/2])W^{-n}W^{-2kn}, k = 0, 1, \dots, N/2 - 1$$

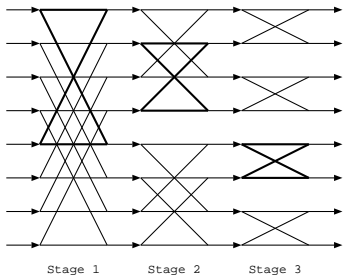
The beauty of this “divide and conquer” strategy is that we replace a Fourier transform of length N with two of length $N/2$, but each of these requires only one-quarter as many multiplications. The wonderful thing about the **Danielson-Lanczos Lemma** is that this can be done recursively: each of the half-length Fourier transforms $F_e[k]$ and $F_o[k]$ that we end up with can further be replaced by two quarter-length Fourier transforms, and so on down by factors of 2. At each stage, we combine input data halfway apart in the sequence (adding or subtracting), *before* performing any complex multiplications.



To compute the N Fourier coefficients $F[k]$ using this recursion we are performing N complex multiplications every time we divide length by 2, and given that the data length N is some power of 2, we can do this $\log_2 N$ times until we end up with just a trivial 1-point transform. Thus, the complexity of this algorithm is $\mathcal{O}(N \log_2 N)$ for data of length N .

The repetitive pattern formed by adding or subtracting pairs of points halfway apart in each decimated sequence has led to this algorithm (popularized by Cooley and Tukey in 1965) being called **the Butterfly**.

This pattern produces the output Fourier coefficients in bit-reversed positions: to locate $F[k]$ in the FFT output array, take k as a binary number of $\log_2 N$ bits, reverse them and treat as the index into the array. Storage requirements of this algorithm are only $\mathcal{O}(N)$ in space terms.



Discrete Cosine Transform; JPEG compression

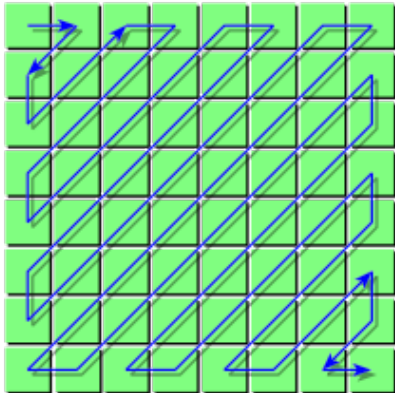
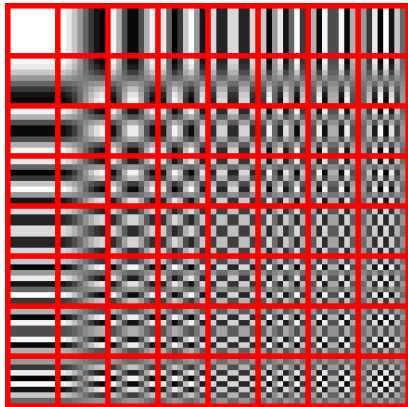
Fourier transforms have vast applications in information processing, but one area central to information theory is data and **image compression**. Images are compressible because neighbouring pixels tend to be highly **correlated**. Projecting such data onto a Fourier (or Fourier-like) basis produces highly **decorrelated** coefficients, and a great many that are 0 or are so small that they need not be encoded.

JPEG image compression is based on running a **discrete cosine transform** (DCT) on small “tiles” of an image, usually (8×8) pixels. For simplicity, discrete cosine functions whose 2D frequencies also form an (8×8) array are used for the projections. The resulting DCT coefficients are quantised, more coarsely for the higher frequencies so fewer bits are used to encode them, and there are long runs of zeroes captured by **run-length coding**.

Compression factors $> 10:1$ are achieved with minimal perceived loss. JPEG compression depends on aspects of human perception, which isn't bothered about the amplitudes of the higher frequencies being accurate. The lower frequencies do matter, but they are far fewer in number. A **quantisation table** specifies how severely quantised the coefficients are for different frequencies; e.g. 2 bits for high frequencies, 7 bits for low.

(Discrete Cosine Transform; JPEG compression, con't)

The 2D cosine patterns used in the DCT (as images are real-valued only) are shown on the left. Their 2D frequency vectors (ω_x, ω_y) form an (8×8) array. Pixel tiles that are also (8×8) arrays are projected onto these, and linear combinations of them can represent any such image. The read-out sequence on the right ensures that there are long runs of 'zeroes' coefficients towards the end, suitable for efficient RLE coding.



Example of JPEG compression by 20:1



Left: An uncompressed monochrome image, encoded at 8 bits per pixel (**bpp**).
Right: JPEG compression by 20:1 (Qual = 10), **0.4 bpp**. The foreground water shows some blocking artifacts at 20:1, and some patches of the water texture are obviously represented by a single vertical cosine in an (8×8) pixel block.

JPEG: compression factor (CF) vs quality factor (QF)

Illustration of progressive variation (left-to-right) in CF \downarrow and QF \uparrow



11. Kolmogorov complexity

- ▶ *Minimal description length of data.*
- ▶ *Algorithmic complexity; computability. Fractals.*

Any set of data can be created by a program, even if (in the worst case) that program simply consists of data statements. The length of such a program defines its **algorithmic complexity**.

A fundamental idea is the measure known as **Kolmogorov complexity**: the complexity of a string of data can be defined as the length of the shortest executable program for computing the string. Thus the complexity is the data string's "**minimal description length**."

It is an amazing fact that the Kolmogorov complexity K of a string is approximately equal to the entropy H of the distribution from which the data string is a randomly drawn sequence. Thus Kolmogorov descriptive complexity is intimately connected with information theory, and indeed K is a way to define ultimate data compression.

Reducing the data to a program that generates it exactly is obviously a way of compressing it. Running the program is a way of decompressing it.

(Kolmogorov complexity, continued)

Fractals are examples of entities that look very complex but in fact are generated by very simple programs (*i.e.*, iterations of a mapping). Therefore, the Kolmogorov complexity of fractals is nearly zero.



In general, Kolmogorov complexity is **not computable**: you never know for sure that you have found the shortest possible description of a pattern.

A sequence $x_1, x_2, x_3, \dots, x_n$ of length n is **algorithmically random** if its Kolmogorov complexity is at least n (*i.e.*, the shortest program that can generate the sequence is a listing of the sequence itself):

$$K(x_1x_2x_3\dots x_n|n) \geq n$$

(Kolmogorov complexity, continued)

An infinite string is defined to be *K-incompressible* if its Kolmogorov complexity, in the limit as the string gets arbitrarily long, approaches the length n of the string itself:

$$\lim_{n \rightarrow \infty} \frac{K(x_1 x_2 x_3 \dots x_n | n)}{n} = 1$$

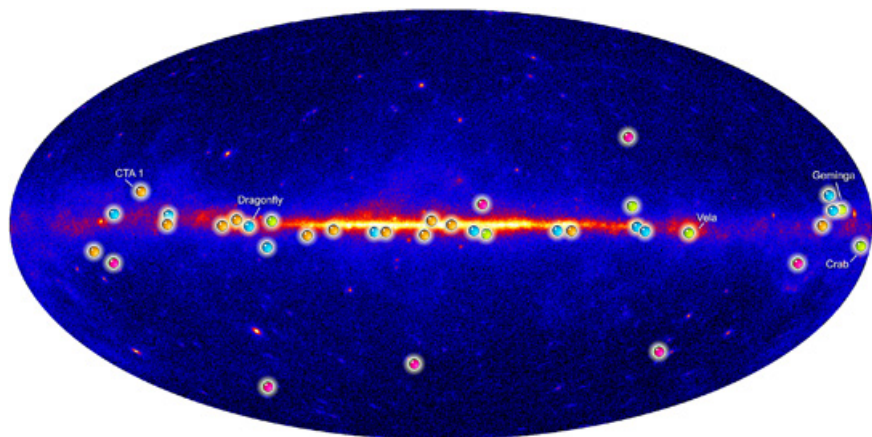
An interesting theorem, called the **Strong Law of Large Numbers for Incompressible Sequences**, asserts that the proportions of 0's and 1's in any incompressible string must be nearly equal.

Moreover, any incompressible sequence must satisfy all computable **statistical tests for randomness**. (Otherwise, identifying the statistical test for randomness that the string failed would reduce the descriptive complexity of the string, which contradicts its incompressibility.)

In this sense the algorithmic test for randomness is the ultimate test, since it includes within it all other computable tests for randomness.

12. Some scientific applications of information theory

- ▶ *Astrophysics; pulsar detection. Extracting signals buried in noise.*
- ▶ *Information theory perspectives in genomics and in neuroscience.*
- ▶ *Biometric pattern recognition.*

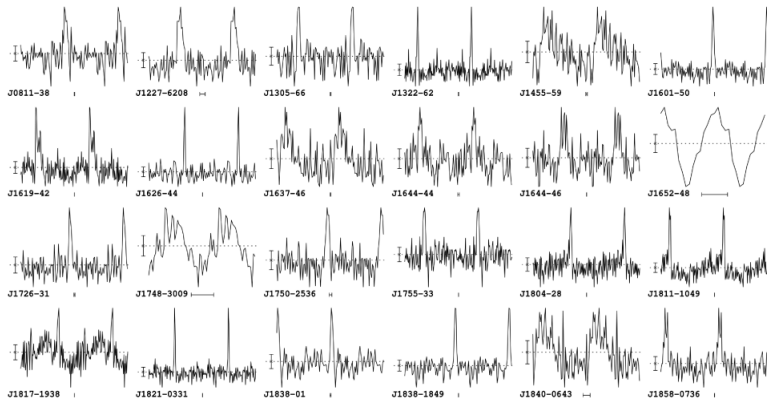


Fermi Pulsar Detections

- New pulsars discovered in a blind search
- Millisecond radio pulsars
- Young radio pulsars
- Pulsars seen by Compton Observatory EGRET instrument

Interpreting astrophysical signals buried in noise

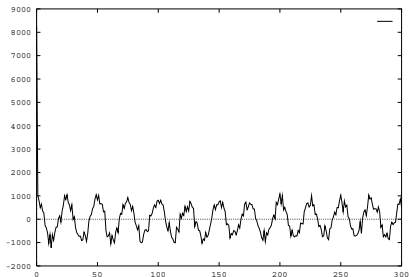
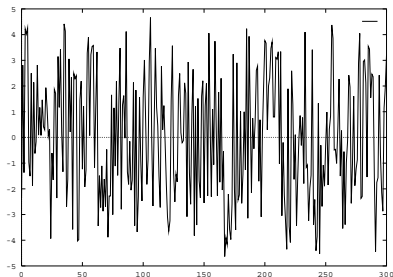
Pulsars are collapsed neutron stars, the dark remnants of a supernova. Because of their tiny size but high mass density, they spin very rapidly (up to 1,000 rotations/sec). Having a strong magnetic dipole, a pulsar's spinning causes emission of an electromagnetic radio beam. If the earth happens to be in the cone of the beam, the signal arrives pulsed at a very precise frequency, the rotation frequency, rather like a lighthouse beacon. But these signals are very faint and buried in the "galactic noise."



Extracting signals buried in noise: auto-correlation function

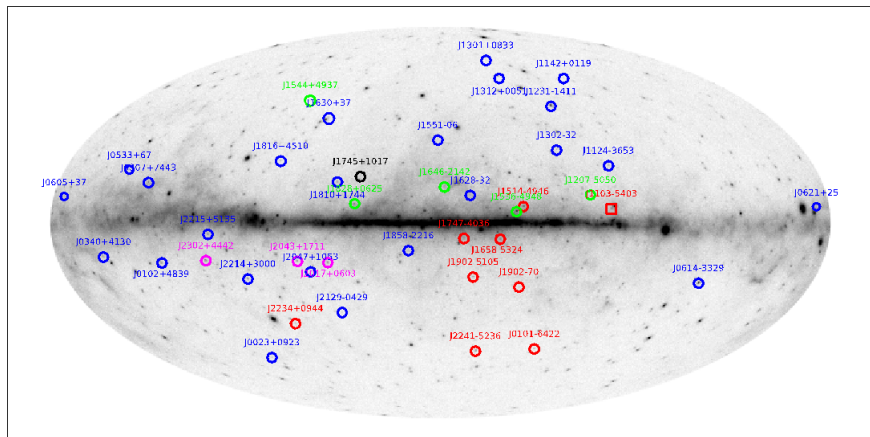
Although the noise may be much larger, the signal has the advantage of being **coherent** (periodic with a precise frequency). The **auto-correlation** integral $\rho_f(t)$ extracts this periodic component from the combined $f(t)$ (left panel), as the incoherent noise tends to average out against itself:

$$\rho_f(t) = \int_{-\infty}^{\infty} f(\tau)f(t + \tau)d\tau \quad (\text{right panel}).$$



Note the resemblance to the convolution of a function with itself, except without the flip: we have $f(t + \tau)$ instead of $f(t - \tau)$ inside the integral.

Pulsar detections within our galaxy (viewed “on-edge”)



It is convenient to compute auto-correlation using Fourier transforms, because if the combined signal $f(t)$ has $\mathcal{FT}\{f(t)\} = F(\omega)$, then the Fourier transform of the auto-correlation function $\rho_f(t)$ that we seek is simply the **power spectral density** of $f(t)$: $\mathcal{FT}\{\rho_f(t)\} = F(\omega)F^*(\omega)$.

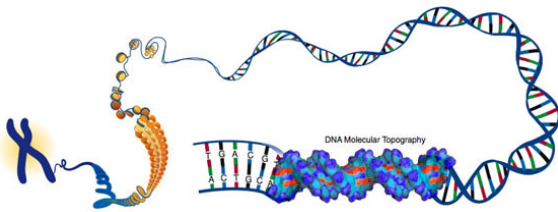
Information theory perspectives in genomics

- The information in DNA is coded in terms of ordered triples of four *nucleotide bases* (Adenine, Guanine, Cytosine, Thymine; or A,G,C,T).
- The ordered triples are called *codons*. There are 64 possible permutations (4 bases taken 3 at a time: $4^3 = 64$ permutations).
- Codons specify amino acids, the building blocks of proteins. There are only 20 primary amino acids, so any given one may correspond to more than one codon. (Example: AGC, AGU, UCG, UCU all specify *serine*.)
- The human genome consists of about 5 billion nucleotide bases, or about 1.7 billion codons, each able to specify ~ 6 bits of information.
- Only about 3% of the human genome specifies genes, of which there are some 24,000. Average gene length ~ 8,000 nucleotide base pairs.



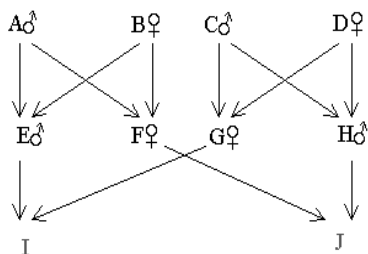
(Information theory perspectives in genomics, con't)

- Therefore on average, a typical gene could specify $\sim 16,000$ bits. But this space is very sparsely populated with meaningful alternatives.
- Many (but not all) of your genes could be passed on by either of your parents (\square, \circ) with equal probability, so genes inheritance entropy is: $H(\text{genes inheritance} \mid \square, \circ) \approx 10,000$ bits per generation. But many genes are indistinguishable in effects, whether maternal or paternal.
- At each generation, the genetic influence of any ancestor is diluted by 50% (as about half of genes come from either parent, each generation).
- Thus (for example): you share $\frac{1}{4}$ of your genes with any grandparent; $\frac{1}{2}$ with either parent, and with full siblings; $\frac{1}{2}$ with any double cousins; 100% with a monozygotic twin; and $\frac{1}{2^N}$ with N^{th} generation ancestors.



(Information theory perspectives in genomics, con't)

Diverse forms of genetic relatedness



Catastrophe: WHERE are all your missing ancestors??

- Going back N generations in your family tree, you must have 2^N ancestors.
- Consider just 950 years back (about 30 generations ago), around the time of William the Conqueror (1066). Your family tree must be populated by $2^{30} \approx 1$ billion ancestors from that time. (Many more if faster reproduction.)
- But the total population of Europe around that time was only a few million.
- WHERE will you find a billion ancestors, among just some million persons?
- The inevitable answer: EVERYBODY then living (who had descendants) was your ancestor; -- and on average, each about 1,000 times over.
- Conclusion: almost everyone today will become the co-ancestors of all of your descendants, within just a few generations. (Why, then, do you think it matters so much whether and with whom you decide to have children?)

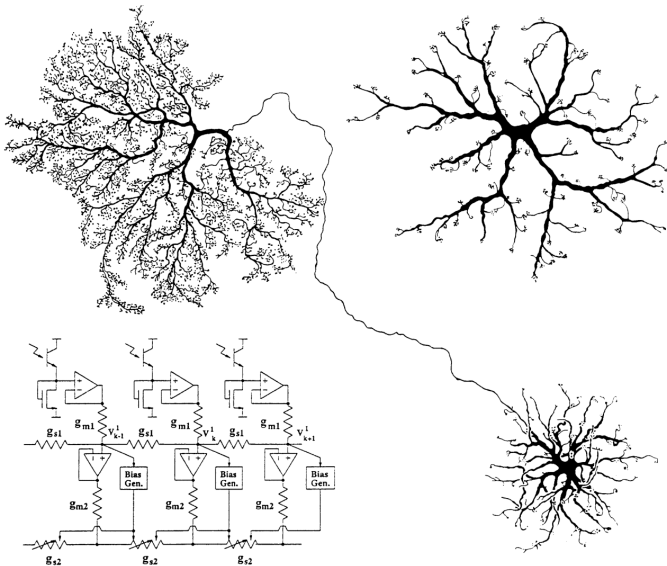
Communal co-ancestry of descendants after N generations



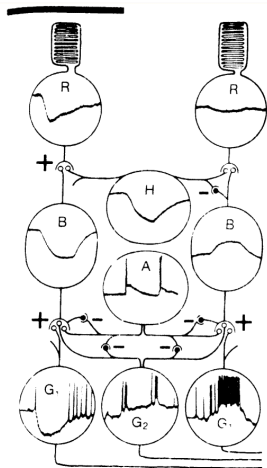
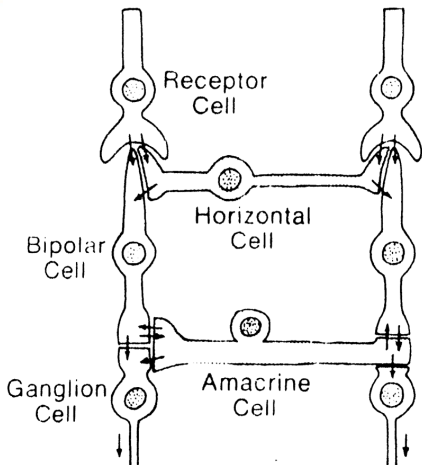
Ancestral state-space entropy $H = -\sum_1^{2^N} 2^{-N} \log_2 2^{-N} = N$ bits. The essence of sexual reproduction is perpetual mixing, and re-mixing, of the gene pool.

Information theory perspectives in neuroscience

- At a synapse, what is the typical rate of information transfer?



(Information theory perspectives in neuroscience, con't)

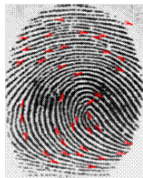


Estimated transmission rate at a graded synapse: 1,650 bits per second.

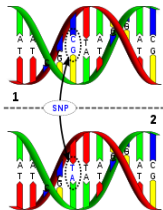
Based on SNR analysis as in Exercise 11. (de Ruyter van Steveninck and Laughlin, *Nature*, 1996.)

Biometric pattern recognition

Some examples of biometric methods and applications



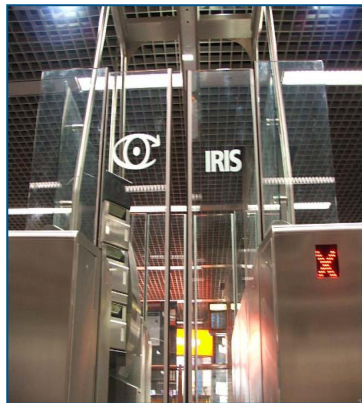
Forensics



"IrisKids" (US) missing children registration and identification



Face recognition ??



Entropy: the key to biometric collision avoidance

- The discriminating power of a biometric depends on its entropy
- Entropy measures the amount of random variation in a population:
 - the number of different states or patterns that are possible;
 - the probability distribution across those possible states
- Entropy H (in bits) corresponds to 2^H discriminable states or patterns
- Surviving large database searches requires large biometric entropy
- Epigenetic features (not genetically determined) make best biometrics

About 1 percent of persons have a monozygotic (“identical”) twin



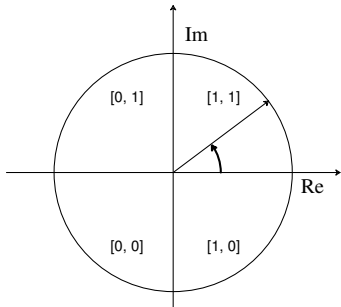
Setting Bits in an IrisCode by Wavelet Demodulation

$$h_{Re} = 1 \text{ if } \operatorname{Re} \int_{\rho} \int_{\phi} e^{-i\omega(\theta_0-\phi)} e^{-(r_0-\rho)^2/\alpha^2} e^{-(\theta_0-\phi)^2/\beta^2} I(\rho, \phi) \rho d\rho d\phi \geq 0$$

$$h_{Re} = 0 \text{ if } \operatorname{Re} \int_{\rho} \int_{\phi} e^{-i\omega(\theta_0-\phi)} e^{-(r_0-\rho)^2/\alpha^2} e^{-(\theta_0-\phi)^2/\beta^2} I(\rho, \phi) \rho d\rho d\phi < 0$$

$$h_{Im} = 1 \text{ if } \operatorname{Im} \int_{\rho} \int_{\phi} e^{-i\omega(\theta_0-\phi)} e^{-(r_0-\rho)^2/\alpha^2} e^{-(\theta_0-\phi)^2/\beta^2} I(\rho, \phi) \rho d\rho d\phi \geq 0$$

$$h_{Im} = 0 \text{ if } \operatorname{Im} \int_{\rho} \int_{\phi} e^{-i\omega(\theta_0-\phi)} e^{-(r_0-\rho)^2/\alpha^2} e^{-(\theta_0-\phi)^2/\beta^2} I(\rho, \phi) \rho d\rho d\phi < 0$$



Why phase is a good variable for biometric encoding

- Phase encodes structural information, independent of contrast
- Phase encoding thereby achieves some valuable invariances
- Phase information has much higher entropy than amplitude
- In harmonic (Fourier) terms, phase “does all the work”
- Phase can be very coarsely quantised into a binary string
- Phase classification is equivalent to a clustering algorithm
- Question: what is the best quantisation of phase (2, 4, 8... sectors)?
- Phase can be encoded in a scale-specific, or a scale-invariant, way

Gabor wavelets encode phase naturally, but in a scale- (or frequency)-specific way

Alternatives exist that encode phase in a total way (independent of scale/frequency), such as the Analytic function (the signal minus its Hilbert Transform $i f_{Hi}(x)$ cousin):

$f(x) - i f_{Hi}(x)$, which is a complex function whose 2 parts are “in quadrature”

Why IrisCode matching is so fast, parallelisable, scalable

Bit streams A and B are data words of two IrisCodes.

Bit streams C and D are their respective mask words.

(data) A	1	0	0	1	0	1	1	0	0	0	1	0	1	1	1	0	...
(data) B	0	1	0	1	1	1	0	0	1	0	0	1	0	1	1	0	...
$A \oplus B$	1	1	0	0	1	0	1	0	1	0	1	1	1	0	0	0	...
(mask) C	1	1	1	0	1	0	1	1	0	0	1	1	1	0	1	1	...
(mask) D	0	1	1	1	1	1	0	1	0	1	1	1	0	1	1	1	...
$C \cap D$	0	1	1	0	1	0	0	1	0	0	1	1	0	0	1	1	...
$(A \oplus B) \cap C \cap D$	0	1	0	0	1	0	0	0	0	0	1	1	0	0	0	0	...

Note that for these 16 bit chunks, only 8 data bits were mutually unmasked by $C \cap D$.

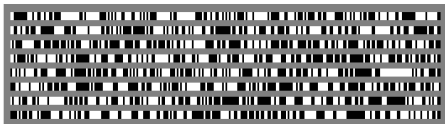
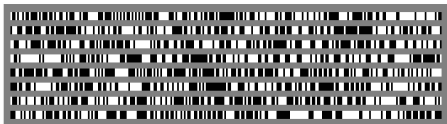
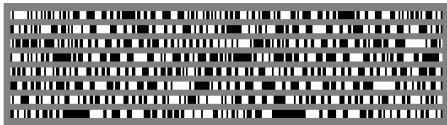
Of those 8, they agreed in 4 and disagreed in 4, so raw Hamming distance is $4/8 = 0.5$ which is typical for comparisons between “Impostors” (unrelated IrisCodes).

Bit-parallel logic programming allows all of this to be done in a single line of C-code, operating on word lengths up to the word-length of the CPU (*e.g.* 64 bits at once):

```
result = (A ^ B) & C & D;
```

Each of the 3 logical parallel operators executes in a single “clock tick” (*e.g.* at 3 GHz).

High entropy gives resistance against False Matches



The probability of two different people colliding by chance in so many bits (e.g. disagreeing in only one-third of their IrisCode bits) is infinitesimal. Thus the False Match Rate is easily made minuscule.

The Doctrine of Suspicious Coincidences



When the recurrence of patterns just by chance is a highly improbable explanation, it is unlikely to be a coincidence.



**UNIVERSITY OF
CAMBRIDGE**

Schedule of Examples Classes, in lieu of supervisions

Students should prepare these Exercises *before* each Examples Class:

1. Monday 19 October, 14:30–15:30, Room FW26: Exercises 1 – 4
2. Monday 26 October, 14:30–15:30, Room FW26: Exercises 5 – 8
3. Tuesday 3 November, 14:30–15:30, Room FW26: Exercises 9 – 12

The complete set of Exercises may be found at this URL:

<http://www.cl.cam.ac.uk/teaching/1516/InfoTheory/materials.html>

where Model Answers will also be published after each Examples Class.