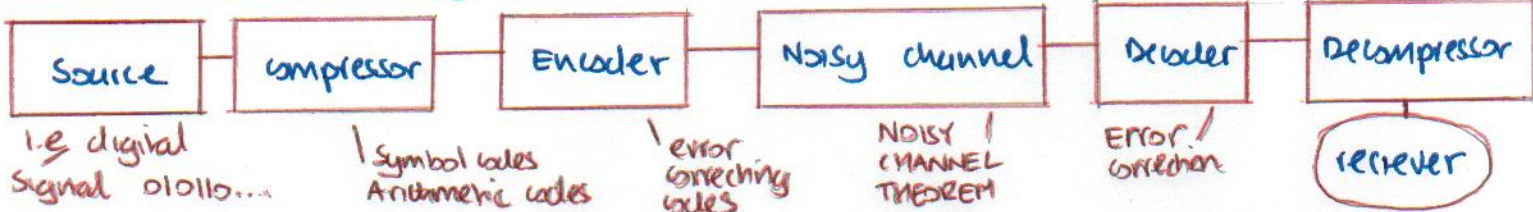


INFORMATION THEORY AND NEURAL NETWORKS

"The fundamental problem of communication is that of reproducing at one point either exactly or approximately a message selected at another point" (Claude Shannon, 1948).

1/ Picture synopsis of course: Source \rightarrow receiver block diagram. (well information theory part).

unique decodability



2/ Information theory is intimately linked with probability theory and Bayesian inference.

Notation and definitions

An ensemble X is a tuple (x, A_X, P_X) where the outcome x is the value of a random variable which takes on one of a set of possible values from ALPHABET A_X with probability $p_i \in P_X$. Now consider joint ensembles of random variables x, y . (Definitions below can easily be extended to more variables). Start from joint probability distribution $P(x, y)$. (which must have PROB property $\sum_{x, y} P(x, y) = 1$).

\Rightarrow BAYES THEOREM $P(y|x) = \frac{P(x, y) P(y)}{P(x)}$ POSTERIOR LIKELIHOOD / PRIOR EVIDENCE

$P(x) = \sum_y P(x, y)$ $P(y) = \sum_x P(x, y)$ $P(x, y) = P(x) P(y|x) = P(y) P(x|y)$

$H(x) = -\sum_x P(x) \log_2 \frac{1}{P(x)}$ $H(x, y) = -\sum_{x, y} P(x, y) \log_2 \frac{1}{P(x, y)}$

SHANNON INFORMATION CONTENT

ENTROPY

JOINT ENTROPY

$H(x|y) = -\sum_{x, y} P(x, y) \log_2 \frac{1}{P(x|y)}$

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

CONDITIONAL ENTROPY

MUTUAL INFORMATION

$D_{KL}(P||Q) = -\sum_x P(x) \log_2 \frac{P(x)}{Q(x)}$

RELATIVE ENTROPY OR KULLBACK-LEIBLER DIVERGENCE (between two probability distributions defined over the same alphabet A_X).

Theorems and inequalities

Jensen's inequality.

If $f(x)$ is a convex \cup function and x a random variable:

$E[f(x)] \geq f(E[x])$

with equality when $x = \text{constant}$.

If $f(x)$ is a concave \cap function and x a random variable:

$E[f(x)] \leq f(E[x])$

This can be used to prove Gibbs' inequality

$D_{KL}(P||Q) \geq 0$ with equality when $P=Q$.

Also: $H(x, y) = H(x) + H(y|x) = H(y) + H(x|y)$

$I(x; y) = I(y; x) \geq 0$

with equality when $X=Y$

Important probability distributions

"Probability of x occurrences of binary outcome with probability f "

BINOMIAL: $P(x|f, N) = \binom{N}{x} f^x (1-f)^{N-x}$

$\binom{N}{x} = \frac{N!}{(N-x)! x!}$ DISCRETE DISTRIBUTION $M=N$ $\sigma^2 = Nf(1-f)$

POISSON $P(x) = \frac{M^x}{x!} e^{-M}$ $\sigma^2 = M$
 GAUSSIAN (NORMAL) $P(x|M, \sigma) = \frac{1}{\sigma \sqrt{2\pi}} e^{-\frac{(x-M)^2}{2\sigma^2}}$
 Approximate Gaussian by Gaussian with $M=5^2$
 $\frac{e^{-M} M^x}{x!} \approx \frac{1}{\sqrt{2\pi M}} e^{-\frac{(x-M)^2}{2M}}$ let $x=M \Rightarrow$ STIRLING'S APPROXIMATION $M! \approx M^M e^{-M} \sqrt{2\pi M}$ INF ①

STIRLING'S APPROXIMATION CONT... $\Rightarrow \ln x \approx x \ln x - x + \frac{1}{2} \ln 2\pi x$ (for $x > 0$)
 $\therefore \ln \binom{N}{x} \approx (N-x) \ln \frac{N}{N-x} + x \ln \frac{N}{x}$ * Define BINARY ENTROPY FUNCTION $H_2(x)$:

$H_2(x) = x \log_2 \frac{1}{x} + (1-x) \log_2 \frac{1}{1-x}$. Note $\frac{d}{dx} H_2(x) = \log_2 \frac{1-x}{x}$. * value for any base. (since all terms are logarithms)
 $\therefore \log_2 \binom{N}{x} \approx NH_2(x/N) - \frac{1}{2} \log_2 [2\pi N \cdot \frac{N-x}{N} \cdot \frac{x}{N}]$.

Central limit theorem: All distributions $P(x|N, p) \rightarrow$ Gaussian as $|Ax| \rightarrow \infty$
 $|Ax| = \#$ values within ensemble.

Bayesian inference Two types of Bayesian inference problems
 1. Direct Hypothesis inference. A set of distinct hypothesis $\{H_i\}$ generate a set of data D . Relationship from H_i to D is given by the forward probability distribution $P(D|H_i)$. Apply Bayes' theorem to infer probability of H_i given observed data D .

$P(H_i|D) = \frac{P(D|H_i)P(H_i)}{P(D)}$ $P(D) = \sum_j P(D|H_j)P(H_j)$
 Examples:
 * CRIMINAL CASES
 * 3 DOORS PROBLEM.

2. PARAMETERISED Hypothesis inference. A particular hypothesis H is parameterised by a set of variables α . i.e. $H = H(\alpha)$. Forward probability distribution is $P(D|H, \alpha)$. Apply Bayes theorem to infer parameter α given observed data D , and hypothesis H . In this statement regard H as a general structure varied by parameter α .


$P(\alpha|D, H) = \frac{P(D|H, \alpha)P(\alpha|H)}{P(D|H)}$ $P(D|H) = \sum_{\alpha} P(D|H, \alpha)P(\alpha|H)$
 Example
 * Degenerate problem.

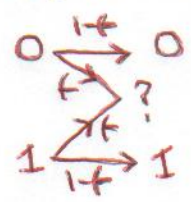
Note in many cases prior $P(\alpha|H)$ will be uniform and thus can be absorbed into the normalisation 'Evidence' $P(D|H)$.

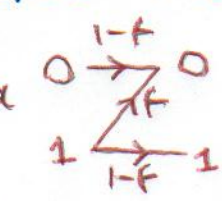
If $P(\alpha|H)$ constant $\Rightarrow P(\alpha|D, H) = \left[1 + \sum_{y \neq \alpha} \frac{P(D|H, y)}{P(D|H, \alpha)} \right]^{-1}$

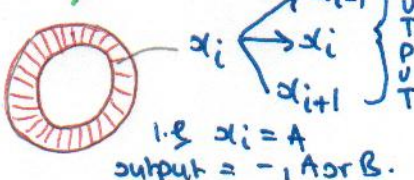
3/ communication channels. (Noisy). A discrete memoryless channel \mathcal{Q} is characterised by an input alphabet A_X , an output alphabet A_Y and a set of conditional probabilities $P(y|x)$ where $x \in A_X, y \in A_Y$.

Define $Q_{ji} \equiv \theta_{ji} = P(y=b_j|x=a_i)$. Write $\{a_i\}$ as a vector \underline{x} and $\{b_j\}$ as a vector \underline{y} .
 $\Rightarrow \underline{P}_y = \underline{Q} \underline{P}_x$ using $P(y) = \sum_x P(x, y) = \sum_x P(x)P(y|x)$. So $P_{y_j} = Q_{ji} P_{x_i}$

Model channels BINARY SYMMETRIC CHANNEL
 $A_X = \{0, 1\}$ $A_Y = \{0, 1\}$
 $P(y=0|x=0) = 1-f$ $P(y=0|x=1) = f$
 $P(y=1|x=0) = f$ $P(y=1|x=1) = 1-f$

 $\underline{Q} = \begin{pmatrix} 1-f & f \\ f & 1-f \end{pmatrix}$

BINARY ERASURE CHANNEL
 $A_X = \{0, 1\}$ $A_Y = \{0, ?, 1\}$
 $P(y=0|x=0) = 1-f$ $P(y=0|x=1) = 0$
 $P(y=?|x=0) = f$ $P(y=?|x=1) = f$
 $P(y=1|x=0) = 0$ $P(y=1|x=1) = 1-f$


Z CHANNEL
 $A_X = \{0, 1\}$ $A_Y = \{0, 1\}$
 $P(y=0|x=0) = 1$ $P(y=1|x=0) = 0$
 $P(y=0|x=1) = 0$ $P(y=1|x=1) = 1-f$

 $\underline{Q}_z = \begin{pmatrix} 1 & 0 \\ 0 & 1-f \end{pmatrix}$

NOISY TYPEWRITER
 $A_X = A_Y = 27$ letters $\{A, B, C, \dots, Z, -\}$
 letters arranged in a circle.

 $P(y=x_{i-1}|x_i) = \frac{1}{3}$
 $P(y=x_i|x_i) = \frac{1}{3}$
 $P(y=x_{i+1}|x_i) = \frac{1}{3}$
 i.e. $x_i = A$
 output = $-, A$ or B . INFO 2

Define CAPACITY $C = \max_{P_X} I(X;Y)$ for a channel described above $P_{Y|X} = \{Q_{ij}\}$

The capacity of a channel is a good measure of the amount of error free information that can be transmitted over the channel per unit time.

Recipe for working out the capacity of a channel. ① Write $P_X = \{P_1, P_2, P_3, \dots, P_{N-1}\}$ ($N = |P_X|$). ② Work out $I(X;Y)$. For mechanistic methods given $P_X = \{P_i\}$ and $P_Y = \{Q_j\}$: (write $\{Q_j\} = \{Q_{j1}, Q_{j2}, Q_{j3}, \dots, Q_{jM}\}$; $M = |P_Y|$) and $\|Q_{ji}\|$:

$I(X;Y) = H(X) - H(X|Y)$. $H(X) = \sum_x P(x) \log_2 \frac{1}{P(x)}$. $H(X|Y) = \sum_{x,y} P(x,y) \log_2 \frac{1}{P(x,y)}$
 $P(x,y) = P(x)P(y|x)$ $\therefore P(x,y) = P(y|x) \frac{P(x)}{P(y)}$ $\therefore H(X|Y) = \sum_{x,y} P(x)P(y|x) \log_2 \frac{P(y)}{P(y|x)P(x)}$

$\therefore I(X;Y) = \sum_x P(x) \log_2 \frac{1}{P(x)} - \sum_{x,y} P(x)P(y|x) \log_2 \frac{P(y)}{P(y|x)P(x)}$ (*)

OR (using P_i, Q_{ji}, Q_{ji}): $I(X;Y) = \sum_{i=1}^N P_i \log_2 \frac{1}{P_i} - \sum_{i=1}^N \sum_{j=1}^M P_i Q_{ji} \log_2 \frac{Q_{ji}}{Q_{ji} P_i}$

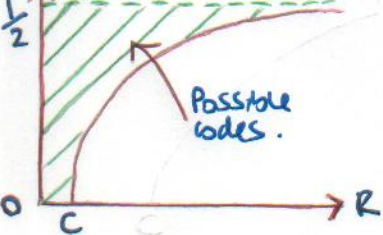
To find OPTIMAL INPUT DISTRIBUTION: $\{P_i^*\}$ Solve $\nabla I = 0$
 where $\nabla = (\frac{\partial}{\partial P_1}, \frac{\partial}{\partial P_2}, \dots, \frac{\partial}{\partial P_{N-1}})$ Note last term is not independent of the others in order to absorb constraint $\sum P_i = 1$ into our minimisation analysis. (Easier than Lagrange multiplier).

$\therefore I(X;Y)$ given $\{P_i^*\} =$ capacity C .

NOISY CHANNEL CODING THEOREM.

For every discrete memoryless channel the capacity $C = \max_{P_X} I(X;Y)$ has the following property: For any $\epsilon > 0$ and $R < C$ for large enough N , there exists a code of length N and RATE $\geq R$ and a decoding algorithm s.t. the maximal probability of bit error is $< \epsilon$.

If bit error probability P_b is acceptable $R(P_b) \leq \frac{C}{1 - H_2(P_b)}$



Definition of rate R : If # codewords $S = 2^k$, $k \in \mathbb{Z}$
 $R \equiv \frac{k}{N}$ where $N =$ length of the codewords.

① Suppose for non block codes i.e. of varying length use 'expected length' L ? \rightarrow see later for expected length)

- So can generate codes with zero bit error at non zero rate!
 Unfortunately the N.C.C.T. does not tell you how.... Derivation of N.C.C.T involves the notion of the typical set. If \underline{x} is typical of $P(x)$ $|\frac{1}{N} \log_2 \frac{1}{P(\underline{x})} - H(x)| < \beta$ where β is some (small) because and $N =$ codeword length. (similar definitions for \underline{x}, y jointly typical of $P(\underline{x}, y)$).

4/ Codes. Motivated by N.C.C.T we will look firstly at error correcting codes whose bit error probabilities and rates lie in shaded area of diagram above.

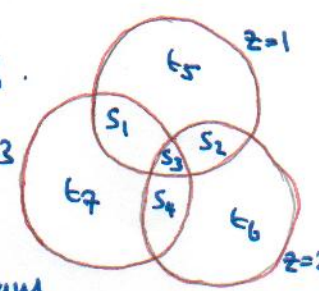
Two examples. (i) REPETITION CODES. (R_N) . These repeat each bit of a binary signal during transmission over a noisy channel (N times). i.e. for source $S = 01010$ under R_3 transmitted signal $T = 000111000111000$. Receiver reads in N bits at a time and takes a majority vote. i.e. $001 \Rightarrow 0$ for R_3 info ③

* Note (*) can be written in a more compact form
 $I(X;Y) = \sum_{x,y} P(x,y) \log_2 \frac{P(x,y)}{P(x)P(y)} = \sum_{x,y} P(x)P(y) \log_2 \frac{P(x,y)}{P(x)P(y)}$
 $\Rightarrow I(X;Y) = \sum_{i,j} P_i Q_{ji} \log_2 \frac{Q_{ji}}{Q_{ji} P_i}$
 In distribution limit for continuous distributions
 $I = \sum_{x,y} P(x)P(y) \log_2 \frac{P(y|x)}{P(y)}$ $\rightarrow \int dx dy P(x)P(y) \log_2 \frac{P(y|x)}{P(y)}$
 i.e. for Gaussian channel: $P(y|x) = \frac{1}{\sqrt{2\pi\sigma^2}} \exp[-(y-x)^2/2\sigma^2]$

Repetition codes RN have error probability $P_b = \sum_{n=(N+1)/2}^N \binom{N}{n} f^n (1-f)^{N-n}$ for a binary symmetric channel. If information stored on one typical hard disc (1 Gbyte) is coded with RN and no errors are expected during a ten year lifecycle $\Rightarrow P_b \sim 10^{-15} \Rightarrow N \sim 61$. So 61 hard drives are required to deliver adequate performance of 1! (use Shannons approximation on lengthy order $\binom{N}{(N+1)/2}$ term above).

A better code is the HAMMING (7,4) code.

If $\underline{s} = (s_1, s_2, s_3, s_4)$ transmit $\underline{t} = (s_1, s_2, s_3, s_4, t_5, t_6, t_7)$ $t_i \oplus s_i \in \{0, 1\}$. t_5, t_6, t_7 are parity bits. They are set so $\sum s_i \oplus t_i$ in each circle on the right is even.



Define SYNDROME $\underline{z} = \{z_1, z_2, z_3\}$ where z_i is the sum of all bits in each circle modulo 2. $+1 \Rightarrow \begin{pmatrix} 0 \\ 1 \end{pmatrix} + 1$. If only one bit of \underline{t} is flipped, computing syndrome and referring to circle diagram allows erroneous bit to be isolated. Hence Hamming (7,4) code is error correcting up to 1 bit error.

SYNDROME \underline{z}	000	001	010	011	100	101	110	111
UNFLIP THIS BIT	NONE	7	6	4	5	1	2	3

Symbol codes are designed to compress information before transmission over a noisy channel. In a symbol code each symbol has an associated binary string code word. i.e. $A_x = \{a, b, c\}$, $P_x = \{1/3, 1/3, 1/3\}$. $\Rightarrow C_x = \{00, 01, 1\}$ (by Huffman encoding scheme). Symbol code C for ensemble X with have EXPECTED LENGTH $L(C, P) = \sum P_i l_i$. $P_i \in \mathcal{P}_X$ and $l_i =$ length of code word C_i .

A symbol code is UNIQUELY DECODABLE if each code word has a distinct binary string and extended strings cannot be represented by more than one string of code symbols. i.e. if $a \Rightarrow 0, b \Rightarrow 00$ [00] could mean 'aa' or 'b'. All PREFIX CODES are uniquely decodable. A symbol code is called a prefix code if no code word is a prefix of any other code word. (Also known as 'instantaneous code' by virtue of 'decode as you receive').

A COMPLETE CODE is uniquely decodable and satisfies the KRAFT INEQUALITY with equality. All uniquely decodable symbol codes satisfy the Kraft inequality.

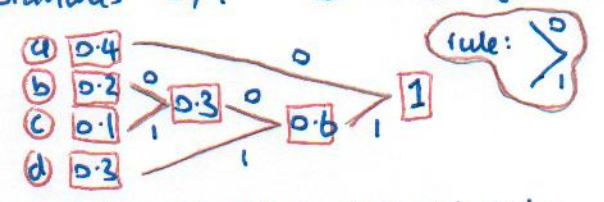
KRAFT INEQUALITY $\sum 2^{-l_i} \leq 1$. Define $z_i \equiv \frac{2^{-l_i}}{z}$ where $z = \sum 2^{-l_i}$

Now from Gibbs inequality $D_{KL}(P||Q) \geq 0 \Rightarrow \sum P_i \log_2 \frac{P_i}{z_i} \geq 0 \Rightarrow \sum P_i \log_2 \frac{1}{z_i} \geq \sum P_i \log_2 \frac{1}{P_i} = H(X)$
 Now $\log_2 \frac{1}{z_i} = \log_2 z + l_i$. $\therefore \log_2 z \sum P_i + \sum P_i l_i \geq H(X)$. Now Kraft inequality $\Rightarrow z \leq 1 \Rightarrow \log_2 z \leq 0$. $\therefore L(C, P) \geq H(X)$ i.e. can't compress further than the entropy of a source. Note also $L(C, P) - H(X) = D_{KL}(P||Q)$

where $Q = \{ \frac{2^{-l_i}}{\sum 2^{-l_i}} \}$. i.e. $D_{KL}(P||Q)$ is the difference between the actual code length and the optimum.

Huffman coding - algorithm for producing optimal complete codes. (1) order code words by probability. combine two lowest probabilities in a tree diagram. (2) repeat until tree is complete. (3) starting from top of tree label branches 0, 1 and thus generate the code.

Example: $A_x = \{a, b, c, d\}$ $P_x = \{0.4, 0.2, 0.1, 0.3\}$
 $\Rightarrow C_x = \{0, 100, 101, 11\}$

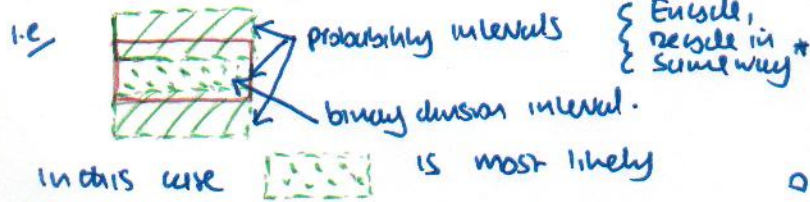


Arithmetic codes combine a probabilistic model with an encoding algorithm that identifies each string with a sub interval of [0,1] of size equal to the probability of that string under the model. The encoding is almost optimal in the sense that the compressed length of a string \approx closely matches the Shannon information content of x given the probabilistic model. Arithmetic codes fit with the philosophy that good compression requires data modelling, in the form of an adaptive Bayesian model.

Note arithmetic codes (and symbol codes) are NOT error correcting so are used for compression prior to transmission over a noisy channel. INFO (4)

Demonstration of Arithmetic coding

Let source string be $\underline{s} = [s_1 s_2 s_3 s_4 \dots]$
 * Superimpose probability model as intervals of $[0,1]$ on binary decision map.
 * Read off code as which is the largest interval within the binary decision map. i.e. 1...
 \Rightarrow b most likely, 10... \Rightarrow bb etc.



0	00	000	0000
	01	010	0100
1	10	100	1000
	11	110	1100
		111	1110
			1111

s_1 $s_1 s_2$ $s_1 s_2 s_3$ $s_1 s_2 s_3 s_4 \dots$

* i.e. decoder needs model.

Simple (Bayesian) Arithmetic coding probability model.

Consider alphabet $A_X = \{a, b, \square\}$ $\square \Rightarrow$ 'end of message'.
 Assume characters are selected RANDOMLY with probabilities $P_X = \{P_a, P_b, P_\square\}$.
 excluding strings with \square terminating them. Let:

BERNOULLI DISTRIBUTION $P(\underline{s} | P_a, F) = \frac{P_a^{F_a} (1-P_a)^{F_b}}{N}$ ← normalising constant where string so far (no \square yet) has F_a 'a's and F_b 'b's.

From Bayes theorem: $P(P_a | \underline{s}, F) = \frac{P(\underline{s} | P_a, F) P(P_a)}{P(\underline{s} | F)}$ = $\frac{P_a^{F_a} (1-P_a)^{F_b} P(P_a)}{P(\underline{s} | F)}$

lumping normalisation N into $P(\underline{s} | F)$. $\frac{P(\underline{s} | F)}{P(\underline{s} | F)} = \int_0^1 dP_a P_a^{F_a} (1-P_a)^{F_b} P(P_a) = \frac{F_a! F_b!}{(F_a + F_b + 1)!}$

If $P(P_a)$ is given by the distribution. (i.e. constant). ↑ Beta function

Now $P(a | \underline{s}, F) = \int_0^1 dP_a P(a | P_a) P(P_a | \underline{s}, F)$. $P(a | P_a) = P_a$

$\therefore P(a | \underline{s}, F) = \int_0^1 dP_a P_a^{F_a+1} (1-P_a)^{F_b} / P(\underline{s} | F)$ using above analysis $\Rightarrow P(a | \underline{s}, F) = \frac{(F_a+1)! F_b!}{(F_a+F_b+2)!}$

$\Rightarrow P(a | \underline{s}, F) = \frac{F_a+1}{F_a+F_b+2}$ Extending this result to larger alphabets $A_X = \{a_i, \square\}$

$\Rightarrow P(a_i | \underline{s}, F) = \frac{F_{a_i} + \alpha}{\sum_i (F_{a_i} + \alpha)}$ where $\alpha = 1$. LAPLACE'S RULE.
 If $\alpha \neq 1$ this corresponds to a DIRICHLET MODEL i.e. a more general case.

Note probability over the characters is likely to be non uniform and prior $P(P_{a_i})$ may not be uniform either. Note Arithmetic codes can generate random numbers if a point is selected randomly between 0 and 1 and decoded (well binary encoded...)

5) Sampling from probability distributions. Often probability distributions are complex s.t. $Z_p = \int P^*(x) d^N x$ (where $P^*(x)$ is unnormalised distribution - known) is hard to calculate. {Motivation for finding Z_p : Statistical physics - partition function etc...}

Laplace's method. * Assume $P^*(x)$ {1D example} has a peak at $x=x_0$. * Taylor expand $\ln P^*(x) |_{x=x_0}$. (Note first derivative = 0 since peak). $\ln P^*(x) \approx \ln P^*(x_0) - \frac{c}{2} (x-x_0)^2 + \dots$ where $c = -\frac{\partial^2}{\partial x^2} \ln P^*(x) |_{x=x_0}$. * Approximate $P^*(x)$ by unnormalised Gaussian $Q(x) = P^*(x) e^{-\frac{c}{2}(x-x_0)^2}$

$Z_Q = \int_{-\infty}^{\infty} Q(x) dx = P^*(x_0) \sqrt{\frac{2\pi}{c}}$. * Let $Z_p \approx Z_Q$. For $P^*(x)$ over k dimensional space:
 $Z_p \approx Z_Q = P^*(x_0) \frac{1}{\sqrt{\det \frac{1}{2\pi} A}}} = P^*(x_0) \frac{(2\pi)^{k/2}}{\sqrt{\det A}}$

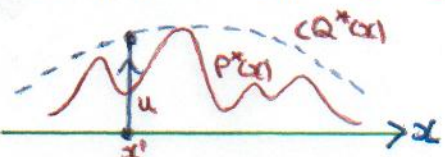
define matrix A s.t. $A_{ij} = -\frac{\partial^2}{\partial x_i \partial x_j} \ln P^*(x) |_{x=x_0}$; $Z_p \approx Z_Q = P^*(x_0) \frac{1}{\sqrt{\det \frac{1}{2\pi} A}}$

[Proof: $\int d^k x \exp[-\frac{1}{2} x^T A x] = \prod_{i=1}^k \int du_i \exp[-\frac{1}{2} \lambda_i u_i^2] = \prod_{i=1}^k \sqrt{\frac{2\pi}{\lambda_i}}$. $\prod_i \lambda_i = \det A$ ← known result.

- change basis x to one where A is diagonalised - elements are eigenvalues λ_i].

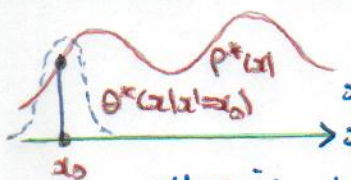
Monte Carlo Methods can show 'volume' of the typical set of an N dimensional distribution $\sim 2^{H(X)}$. So can become enormous. A similar approach to Laplace's method is IMPORTANCE SAMPLING where simpler distribution $Q(x) = \frac{P^*(x)}{Z_Q}$ is used to approximate $P(x)$ INFO (5)

IMPORTANCE SAMPLING CONT... Define estimator $\Phi[x] \approx \sum w_i \phi(x_i)$ where $w_i = \frac{P^*(x_i)}{Q^*(x_i)}$ and $\phi(x)$ is function whose expectation $\Phi[x]$ estimate. i.e. mean if $\phi(x) = x$, variance + mean² if $\phi(x) = x^2$. TWO DIFFICULTIES: (i) long time to sample from typical set unless $Q \approx P$ (ii) weights of samples within typical set will vary by large factors because probabilities of points differ by $\sim e^{\sqrt{N}}$ (Analysis involving N dimensional Gaussians).



REJECTION SAMPLING. Find c s.t. $cQ^*(x) > P^*(x)$ for all x . * Randomly (uniformly) sample from $\{x\}$. * Evaluate $cQ^*(x)$ at sample point x' . * $u \sim \text{unif}(0, 1)$ * Evaluate $P^*(x')$. * If $u > P^*(x')$ REJECT and start again. * Determine accept x' as valid sample. $c \rightarrow P \Rightarrow c$ large \Rightarrow rejection frequency high - so slow. Not good for high dimensions.

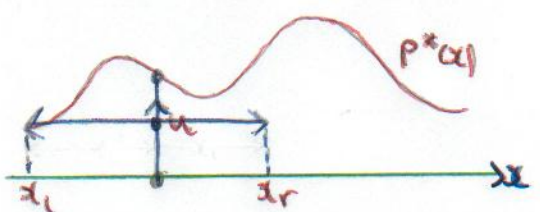
METROPOLIS - HASTINGS * Sample first point x_0 randomly from $\{x\}$ underneath $P^*(x)$. * Evaluate $P^*(x_0)$. $Q^*(x|x_0)$ is defined (i.e. a gaussian) and varies with x_0 . (i.e. mean of gaussian). * Sample next $x = x_1$ and compute: $\alpha = \frac{P^*(x_1) Q^*(x_0|x_1)}{P^*(x_0) Q^*(x_1|x_0)}$. If $\alpha \geq 1$ accept new point x_1 .



otherwise try again. Example of a MARKOV CHAIN where next step depends on the previous. * Random walk in sample space. i.e. if typical dimension of sample space is $L = NE$ where E is the average step size and N is an integer: # steps is (on average) N^2 . \therefore if τ is time interval between steps, total time to traverse L is $\sim N^2 \tau = \frac{L^2}{E^2} \tau$. Characteristic of all Monte Carlo methods used here (from Metropolis) and $\frac{L^2}{E^2}$ below.

GIBBS SAMPLING * Sample from (more tractable - assumed) CONDITIONAL DISTRIBUTION i.e. $x_i^{(t+1)} \sim P(x_i | \{x_j^{(t)}\}_{j \neq i})$. [where $x_i^{(t)}$ is i th component of state $x^{(t)}$]. Most (all?) Markov chain processes are ERGODIC i.e. \rightarrow to some invariant distribution as sample time $\rightarrow \infty$. If not they may tend to a limit cycle between periodic irreducible sets.

SLICE SAMPLING * Sample point x' from $\{x\}$ (randomly). * Sample $u \sim \text{unif}(0, P^*(x))$. * Draw horizontal interval $[x_l, x_r]$ and uniformly sample from it. $x'' \sim \text{unif}(x_l, x_r)$. * If $u > P^*(x'')$ modify x_l, x_r until $u \leq P^*(x'')$. Accept this x'' as a valid sample. PRO'S - No rejections. CON'S - random walk still, need to carefully adjust x_l, x_r s.t. DETAILED BALANCE is satisfied. i.e. uniform distribution of (x, u) under $P^*(x)$ is invariant. (in what?)



VARIATIONAL METHODS consider probability distributions of the form $P(x|\beta, \xi) = \frac{e^{-\beta E(x|\xi)}}{Z(\beta|\xi)}$ where $Z(\beta|\xi) = \sum e^{-\beta E(x|\xi)}$. As before Z is desirable but often intractable \approx to calculate. As before approximate P by simpler distribution Q .

Define $\tilde{F}(\theta) = \sum_x Q(x|\theta) \ln \frac{Q(x|\theta)}{e^{-\beta E(x|\xi)}} = D_{KL}(Q||P) + \beta F$ where $\beta F \equiv -\ln Z$. VARIATIONAL FREE ENERGY By Gibbs inequality $D_{KL}(Q||P) \geq 0 \Rightarrow \tilde{F} - F \geq 0$ if $\beta > 0$. $\Rightarrow \tilde{F} \geq F$. So vary parameters $\{\theta\}$ (possibly more than one) to minimize \tilde{F} . Minimize \tilde{F} and you will get a better estimate of F . Note $Z = e^{-\beta F}$ so can calculate Z this way.

ISING MODELS - Array of spins with couplings and applied field H that modifies $E(x|\xi)$. let x be state vector of a system with N spins where $x_n \in \{-1, +1\}$.

is:
$$E(\alpha; \xi) = - \left[\sum_{m, n \in N} J_{mn} \alpha_m \alpha_n + \sum_n H_n \alpha_n \right]$$

$$C = \frac{\partial^2 \bar{E}}{\partial T^2} \text{ where } \bar{E} = \frac{1}{Z} \sum_{\alpha} e^{-\beta E(\alpha)} E(\alpha)$$

$$Z = \sum_{\alpha} e^{-\beta E(\alpha)} \therefore \frac{\partial \ln Z}{\partial \beta} = \frac{1}{Z} \sum_{\alpha} (-E) e^{-\beta E} = -\bar{E}$$

$$\frac{\partial^2 \ln Z}{\partial \beta^2} = \text{Var}(E)$$

can generalise: $J \rightarrow J_{mn}, H \rightarrow h_n$ {spin glasses, Hopfield networks...}

ISING MODEL produces phase transitions as J, h, β are varied. { C_p curve temp}. Note energy fluctuations have a peak. As an aside note relationship between heat capacity C and $\text{Var}(E)$.

Now $\frac{\partial \bar{E}}{\partial T} = - \frac{\partial}{\partial T} \left(\frac{\partial \ln Z}{\partial \beta} \right) = - \frac{\partial^2 \ln Z}{\partial \beta^2} \frac{\partial \beta}{\partial T} = - \text{Var}(E) \left(-\frac{1}{k_B T^2} \right) \Rightarrow C = \frac{\text{Var}(E)}{k_B T^2}$

Now variational free energy (defined above) $\tilde{F} = \frac{1}{\beta} \sum_{\alpha} Q(\alpha) \ln Q(\alpha) - \sum_{\alpha} Q(\alpha) E(\alpha)$
 $\Rightarrow \beta \tilde{F} = \sum_{\alpha} (Q \ln Q - Q \ln e^{-\beta E}) = - \sum_{\alpha} Q \ln \frac{1}{Q} + \beta \sum_{\alpha} Q E = \beta \langle E \rangle_Q - H_Q$

Consider separable Q : $Q = \frac{1}{Z_Q} \exp \left[\sum_n a_n \alpha_n \right] \Rightarrow H_Q = \sum_n H_2(a_n)$

where $H_2(a_n) = \ln \sum_{\alpha_n} e^{a_n \alpha_n} = \ln (e^{a_n} + e^{-a_n})$ and $z_n = \frac{e^{a_n}}{e^{a_n} + e^{-a_n}} = \frac{1}{1 + e^{-2a_n}}$

using ξ convention $\langle E \rangle_Q = (J_{mn} \bar{\alpha}_m \bar{\alpha}_n + h_n \bar{\alpha}_n) (-1)$
 where $\bar{\alpha}_n = \frac{e^{a_n} - e^{-a_n}}{e^{a_n} + e^{-a_n}} = \tanh(a_n) = 2z_n - 1$. So $\beta \tilde{F} = -\beta (J_{mn} \bar{\alpha}_m \bar{\alpha}_n + h_n \bar{\alpha}_n) - \sum_n H_2(a_n)$

$\Rightarrow \frac{\partial \beta \tilde{F}}{\partial a_m} = -\beta \left\{ J_{mn} \left[\bar{\alpha}_m \frac{\partial \bar{\alpha}_n}{\partial a_m} + \bar{\alpha}_n \frac{\partial \bar{\alpha}_m}{\partial a_m} \right] + h_n \frac{\partial \bar{\alpha}_n}{\partial a_m} \right\} - \sum_n \frac{\partial H_2(a_n)}{\partial a_m}$

Now $\frac{\partial \bar{\alpha}_n}{\partial a_m} = \frac{\partial \bar{\alpha}_n}{\partial z_m} \frac{\partial z_m}{\partial a_m} = 2 \delta_{nm} \left(\frac{\partial z_m}{\partial a_m} \right)$. $\frac{\partial H_2(a_n)}{\partial a_m} = \ln \left(\frac{1-z_n}{z_n} \right) \delta_{nm} \left(\frac{\partial z_m}{\partial a_m} \right)$

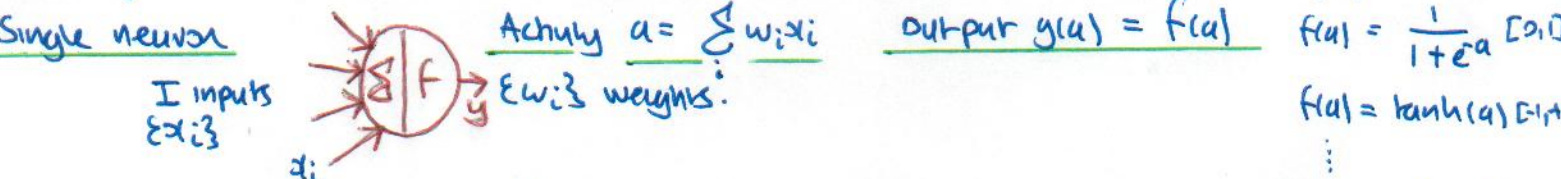
$\therefore \frac{\partial \beta \tilde{F}}{\partial a_m} = -\beta \left\{ J_{mn} \left[2 \bar{\alpha}_m \delta_{nm} \left(\frac{\partial z_m}{\partial a_m} \right) + 2 \bar{\alpha}_n \delta_{nm} \left(\frac{\partial z_m}{\partial a_m} \right) \right] + 2 h_n \delta_{nm} \left(\frac{\partial z_m}{\partial a_m} \right) \right\} - \sum_n \ln \left(\frac{1-z_n}{z_n} \right) \delta_{nm} \left(\frac{\partial z_m}{\partial a_m} \right)$

Now $\frac{1-z_n}{z_n} = \frac{1 + e^{-2a_n} - 1}{1 + e^{-2a_n}} = e^{-2a_n}$. $\therefore \ln \left(\frac{1-z_n}{z_n} \right) = -2a_n$

$\therefore \frac{\partial \beta \tilde{F}}{\partial a_m} = -2\beta \left(\frac{\partial z_m}{\partial a_m} \right) \left\{ 2 J_{mn} \bar{\alpha}_n + h_m \right\} + 2 \frac{\partial z_m}{\partial a_m} a_m$

$\Rightarrow \frac{\partial \beta \tilde{F}}{\partial a_m} = 0$ when $a_m = \beta \sum_n (J_{mn} \bar{\alpha}_n + h_m)$ where $\bar{\alpha}_n = \tanh(a_n)$.
 MEAN FIELD EQUATIONS FOR SPIN SYSTEM.

6) NEURAL NETWORKS Unlike conventional CPUs and memory, organic brains are robust to error, associative (content addressed), distributed, massively // ...
 Idea of neural networks is to emulate this structure.



ARCHITECTURE: $\{w_i\}; f(a)$ ACTIVITY RULE: Dynamics of a with time.
 LEARNING RULE: Dynamics of $\{w_i\}$.

Idea is to train \neq neuron to give desired output t given input. (Neuron \rightarrow classifier). Simple learning rule: $\Delta w_i = \eta x_i \Delta y$ where $t = y + \Delta y$. Capacity of single neuron with I inputs is $\sim 2I$ (2 bits / connection). i.e. single neuron will almost certainly fail to memorize $> 2I$ states, regardless of sophistication of learning rule.

Learning as inference: Good learning rule is to minimize $M(\underline{w})$ by changing $\underline{w} = \|\underline{w}\|$. $M(\underline{w}) = G(\underline{w}) + \alpha E(\underline{w})$. $E(\underline{w}) = \frac{1}{2} \sum w_i^2$. $G(\underline{w}) = - \sum_n [t^{(n)} \log y^{(n)} + (1-t^{(n)}) \log (1-y^{(n)})]$ where $t^{(n)}$ is nth desired output.

* Interpret output $y(x; \underline{w})$ of neuron as (when \underline{w} are specified) defining the probability that an input x belongs to a class $t=1$ rather than alternative $t=0$. Thus $y(x; \underline{w}) \equiv P(t=1 | x, \underline{w})$. Then each value of \underline{w} defines a different hypothesis about the probability of class 1 relative to class 0 as a function of x .

So: $P(t=1 | \underline{w}, x) = y$ $P(t=0 | \underline{w}, x) = 1-y \Rightarrow P(t | \underline{w}, x) = y^t (1-y)^{1-t} = e^{t \log y + (1-t) \log (1-y)}$
 Hence $P(t | \underline{w}, x) = e^{-G(\underline{w})}$ in general. $P(\underline{w} | x) = \frac{1}{Z(\underline{w}, x)} e^{-\alpha E(\underline{w})}$
 \Rightarrow (Bayes theorem) $P(\underline{w} | D, x) = \frac{P(D | \underline{w}) P(\underline{w} | x)}{P(D | x)} = \frac{1}{Z} e^{-M(\underline{w})}$
 So minimizing $M(\underline{w}) \Rightarrow$ maximizing $P(\underline{w} | D, x)$.

Hopfield Network

- Example of a feedback network (as opposed to a feedforward one). Learning rule is HEBB rule



$\frac{dw_{ij}}{dt} \sim \text{correlation}(x_i, x_j)$ (associative memory).

Capacity: $w_{ij} = x_i^{(n)} x_j^{(n)} + \sum_{m \neq n} x_i^{(m)} x_j^{(m)}$
 $\Rightarrow a_i = \sum_j x_j^{(n)} x_j^{(n)} x_j^{(m)} + \sum_{j: m \neq n} x_j^{(m)} x_j^{(m)} x_j^{(n)}$
 $\Rightarrow E(a_i) = (I-1)x_i^{(n)}$
 $\text{var}(a_i) = (I-1)(N-1)$



Convention: w_{ij} denotes connection from neuron j to neuron i
 Architecture: I neurons, all connected together through symmetric, bidirectional connections i.e. $w_{ij} = w_{ji}$.

Activity rule: $\sigma(a) = \Theta(a) \equiv \begin{cases} 1 & a \geq 0 \\ -1 & a < 0 \end{cases}$

Synchronous updates: All neurons compute their activation $a_i = \sum w_{ij} x_j$ then update their states simultaneously to $x_i = \sigma(a_i)$.

OR asynchronous updates: one neuron at a time computes its activation and updates its state. (Fixed or random sequence).

Learning rule: Intentional MEMORIES ARE STABLE STATES of network's activity rule. Each memory is a binary pattern $\underline{x} \in \{-1, 1\}$.

$w_{ij} = \eta \sum_n x_i^{(n)} x_j^{(n)}$ $\eta = \frac{1}{N}$ is often convenient (prevents largest possible weight growing with N).

CONTINUOUS Hopfield network has:

$a_i = \sum_j w_{ij} x_j$ where $x_i = \tanh(\beta a_i)$ or fix η and write $x_i = \tanh(\beta a_i)$
 (beta = gain). NOTE SIMILARITY TO ISING MODEL!! INFO (9)