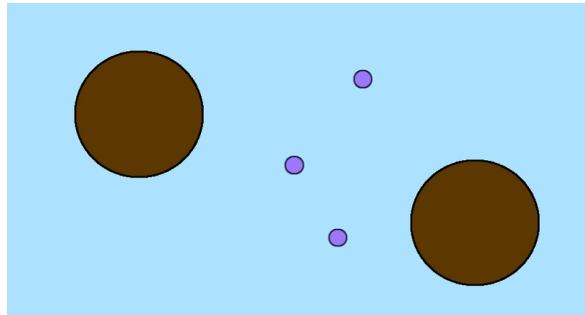


Entropy

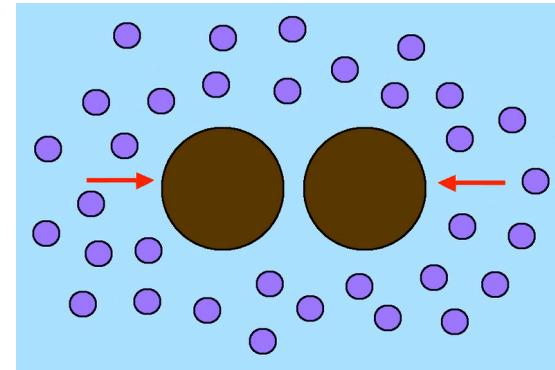
- 1 (Boltzmann) Physics (micro-macro)
- 2 (Shannon) Communication (compression)
- 3 (Bayes) ML - KL Divergence (prior to posterior)

Entropy

Communication

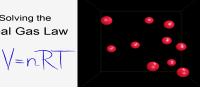


Why?



What is the Hope in Learning?

Solving the Ideal Gas Law
 $PV=nRT$





Avogadro constant $\equiv 6.02 \times 10^{23} \text{ particles mol}^{-1}$



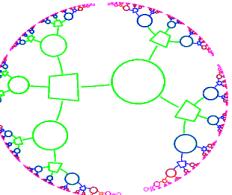
Kilo $= 10^3$ $= 2^{10}$
 Mega $= 10^6$ $= 2^{20}$
 Giga $= 10^9$ $= 2^{30}$
 Tera $= 10^{12}$ $= 2^{40}$
 Peta $= 10^{15}$ $= 2^{50}$
 Exa $= 10^{18}$ $= 2^{60}$
 Zetta $= 10^{21}$ $= 2^{70}$
 Yotta $= 10^{24}$ $= 2^{80}$

Large number to the rescue!

Entropy

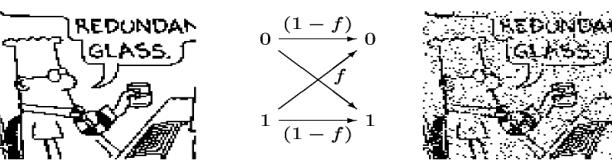
Communication

Information theory,
pattern recognition,
and neural networks



Lecture 1:
Introduction to
Information Theory

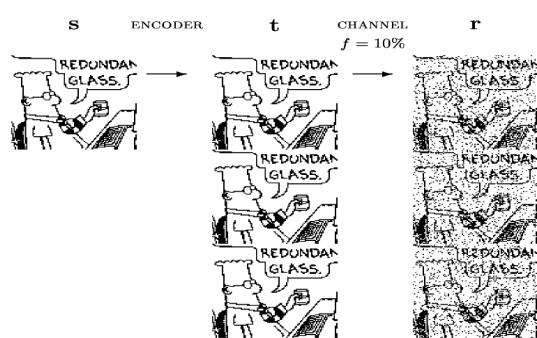
David Mackay



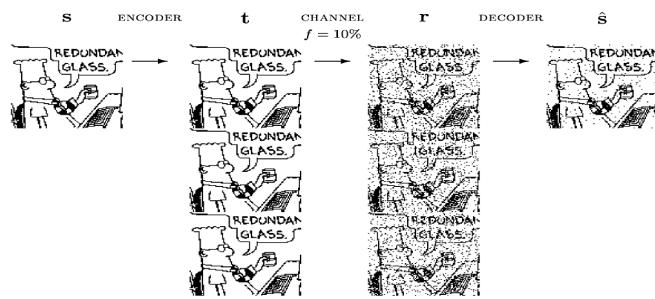
$f = 0.1$

Error rate 10% mistakes

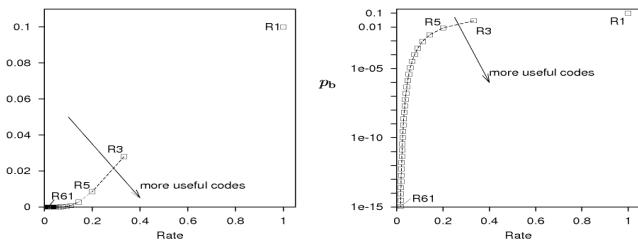
Repetition code 'R3'



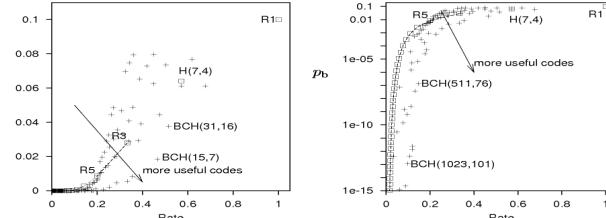
Repetition code 'R3'



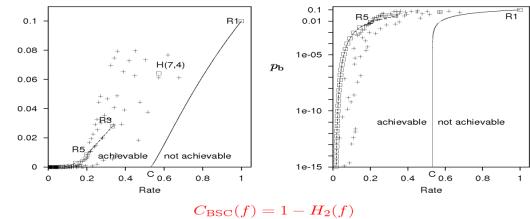
Performance of repetition codes



What's achievable?



Shannon's noisy-channel coding theorem



If $f=1/2$ (totally random, no real information), then $C=0$ (zero rate of transmission)!

Entropy

$$H[x] = - \sum_x p(x) \log_2 p(x)$$

- Important quantity in
- coding theory
 - statistical physics
 - machine learning

Entropy

- Coding theory: x discrete with 8 possible states; how many bits to transmit the state of x ?
- All states equally likely

$$H[x] = -8 \times \frac{1}{8} \log_2 \frac{1}{8} = 3 \text{ bits.}$$

Entropy

- In how many ways can N identical objects be allocated M bins?

$$W = \frac{N!}{\prod_i n_i!}$$

$$H = \frac{1}{N} \ln W \simeq - \lim_{N \rightarrow \infty} \sum_i \left(\frac{n_i}{N} \right) \ln \left(\frac{n_i}{N} \right) = - \sum_i p_i \ln p_i$$

- Entropy maximized when

$$\forall i : p_i = \frac{1}{M}$$

Entropy

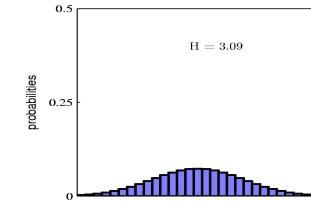
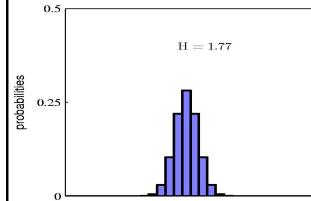
x	a	b	c	d	e	f	g	h
$p(x)$	$\frac{1}{2}$	$\frac{1}{4}$	$\frac{1}{8}$	$\frac{1}{16}$	$\frac{1}{64}$	$\frac{1}{64}$	$\frac{1}{64}$	$\frac{1}{64}$
code	0	10	110	1110	111100	111101	111110	111111

$$\begin{aligned} H[x] &= -\frac{1}{2} \log_2 \frac{1}{2} - \frac{1}{4} \log_2 \frac{1}{4} - \frac{1}{8} \log_2 \frac{1}{8} - \frac{1}{16} \log_2 \frac{1}{16} - \frac{4}{64} \log_2 \frac{1}{64} \\ &= 2 \text{ bits} \end{aligned}$$

$$\begin{aligned} \text{average code length} &= \frac{1}{2} \times 1 + \frac{1}{4} \times 2 + \frac{1}{8} \times 3 + \frac{1}{16} \times 4 + 4 \times \frac{1}{64} \times 6 \\ &= 2 \text{ bits} \end{aligned}$$

If not evenly distributed, the entropy is getting smaller.

Entropy



Learning, in a way, is trying to narrow down the choices of "parameters" (the unknown) using the information obtained in the data set observed, or decreasing the entropy that is consistent with the observation. (also recall that the choice of logic functions for 3-variable problem: from 256 to 8 such functions).

C.G. Shannon: The Bell System Technical Journal, Vol 27, pp 399-415, July-Oct. 1948.

$H(p_1, p_2, \dots, p_n)$

Condition ① H should be continuous in p_i .

Condition ② If all the p_i are equal, $p_i = \frac{1}{n}$, then H should be a monotonic increasing function of n , with equally likely events there is more choices, or uncertainty, when there are more possible events.

Condition ③ If a choice be broken down into two successive choices, the original H should be the weighted sum of the individual

values of H .

**Another example illustrates the same

use the same sequence of events though

if differently seen variable 1

but this can be asked

if the variable zero or not?

equally likely when there are many choices, the individual iterates for m steps to a choice from S^m ident of different.

$A(S) = A(n)$, we can decompose a choice from S^m into a series of m choices from S equally likely. $A(S^m) = mA(S)$. Similarly, $A(t^m) = nA(t)$.

Given, we can choose n arbitrarily large and still get a result of $\left[\frac{ln t^n}{ln n}\right] = m$, or equivalently $\Rightarrow \frac{m}{n} \leq \frac{ln t^n}{ln n} \leq \frac{m}{n+1} \Rightarrow \left|\frac{m}{n} - \frac{ln t^n}{ln n}\right| < \varepsilon$ infinitesimal small by increasing n .

monic function, so $\Rightarrow A(t^m) \leq A(S^m)$

$A(S^m) \leq (mt)^t A(t) \Rightarrow \frac{m}{n} \leq \frac{A(t)}{A(t)} \leq \frac{m}{n+1}$.

If suppose we have n and t possible events whose probabilities are P_1, P_2, \dots, P_n . Their probabilities are known but that's concerning which event will occur. Can we find a measure $A(t)$ involved in this thus $\frac{A(t)}{A(t)} = \frac{P_1}{P_1}$ or $A(t) = \frac{A(t)}{A(t)} P_1 = K P_1$ with K being a positive constant.

(b) Now suppose we have a choice from n possible (comensurable) refined probabilities $P_1 = \frac{n_1}{\sum n_i}$ the same choices, we can break down a choice possibilities into a choice from n possibilities with them, if t the one who chooses, a choice from n with probabilities.

(example) $n=3$ $\begin{cases} P_1 = \frac{n_1}{n} = \frac{3}{6} \\ P_2 = \frac{n_2}{n} = \frac{3}{6} \\ P_3 = \frac{n_3}{n} = \frac{3}{6} \end{cases} \quad \begin{matrix} n_1 \\ n_2 \\ n_3 \end{matrix}$

General Case

$$\begin{cases} P_1 = \frac{n_1}{\sum n_i} \\ P_2 = \frac{n_2}{\sum n_i} \\ \vdots \\ P_k = \frac{n_k}{\sum n_i} \end{cases} \Rightarrow \begin{matrix} n_1 \\ n_2 \\ \vdots \\ n_k \\ \text{such branches} \end{matrix} \rightarrow \begin{matrix} n_1 \\ n_2 \\ \vdots \\ n_k \\ \text{total branches} \end{matrix}$$

$\Rightarrow \frac{n_1}{\sum n_i} + \frac{n_2}{\sum n_i} + \dots + \frac{n_k}{\sum n_i} = \frac{\sum n_i}{\sum n_i} = 1$

$H(Y_{n_1}, Y_{n_2}, \dots, Y_{n_k}) = H(P_1, P_2, \dots, P_k) + \sum_{i=1}^k P_i H(Y_{n_i}, Y_{n_i}, \dots)$

$A(\sum n_i) = K P_1 (\sum n_i)$

$\therefore H(P_1, P_2, \dots, P_k) = K \left[P_1 (\sum n_i) - \frac{1}{P_1} P_1 H(Y_{n_1}, Y_{n_1}, \dots) \right]$

$= K \left[\sum n_i P_1 \log \left(\frac{n_1}{P_1} \right) - \frac{1}{P_1} P_1 H(Y_{n_1}, Y_{n_1}, \dots) \right]$

$\Rightarrow \text{If } P_i \text{ are irrational, we could approximate them to rational, and the same expression must hold true by our continuity assumption (Goal)}$

Entropy in ML?

KL Divergence

Foundation for Approximation

KL Divergence

Foundation for Approximation

The **relative entropy or Kullback–Leibler divergence** between two probability distributions $P(x)$ and $Q(x)$ that are defined over the same alphabet \mathcal{A}_X is

$$D_{\text{KL}}(P||Q) = \sum_x P(x) \log \frac{P(x)}{Q(x)}. \quad (2.45)$$

The relative entropy satisfies *Gibbs' inequality*

$$D_{\text{KL}}(P||Q) \geq 0 \quad (2.46)$$

with equality only if $P = Q$. Note that in general the relative entropy is not symmetric under interchange of the distributions P and Q : in general $D_{\text{KL}}(P||Q) \neq D_{\text{KL}}(Q||P)$, so D_{KL} , although it is sometimes called the ‘KL distance’, is not strictly a distance. The relative entropy is important in pattern recognition and neural networks, as well as in information theory.

2.7 Jensen's inequality for convex functions

The words ‘convex \smile ’ and ‘concave \frown ’ may be pronounced ‘convex-smile’ and ‘concave-frown’. This terminology has useful redundancy: while one may forget which way up ‘convex’ and ‘concave’ are, it is harder to confuse a smile with a frown.

Convex \smile functions. A function $f(x)$ is convex \smile over (a, b) if every chord of the function lies above the function, as shown in figure 2.10; that is, for all $x_1, x_2 \in (a, b)$ and $0 \leq \lambda \leq 1$,

$$f(\lambda x_1 + (1 - \lambda)x_2) \leq \lambda f(x_1) + (1 - \lambda)f(x_2). \quad (2.47)$$

A function f is strictly convex \smile if, for all $x_1, x_2 \in (a, b)$, the equality holds only for $\lambda = 0$ and $\lambda = 1$.

Similar definitions apply to concave \frown and strictly concave \frown functions.

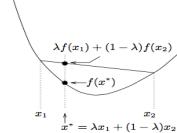


Figure 2.10. Definition of convexity.

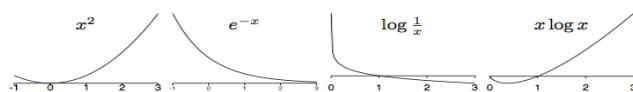
Jensen's inequality. If f is a convex \smile function and x is a random variable then:

$$\mathbb{E}[f(x)] \geq f(\mathbb{E}[x]), \quad (2.48)$$

where \mathbb{E} denotes expectation. If f is strictly convex \smile and $\mathbb{E}[f(x)] = f(\mathbb{E}[x])$, then the random variable x is a constant.

Some strictly convex \smile functions are

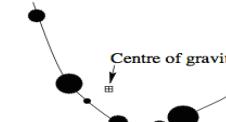
- x^2 , e^x and e^{-x} for all x ;
- $\log(1/x)$ and $x \log x$ for $x > 0$.



Proof (more of a argument for it)

A physical version of Jensen's inequality runs as follows.

If a collection of masses p_i are placed on a convex \smile curve $f(x)$ at locations $(x_i, f(x_i))$, then the centre of gravity of those masses, which is at $(\mathbb{E}[x], \mathbb{E}[f(x)])$, lies above the curve.



The relative entropy or Kullback–Leibler divergence between two probability distributions $P(x)$ and $Q(x)$ that are defined over the same alphabet \mathcal{A}_X is

$$D_{\text{KL}}(P||Q) = \sum_x P(x) \log \frac{P(x)}{Q(x)}. \quad (2.45)$$

The relative entropy satisfies Gibbs' inequality

$$D_{\text{KL}}(P||Q) \geq 0 \quad (2.46)$$

1st

$$D_{\text{KL}}(P||Q) = \sum_x P(x) \log \frac{P(x)}{Q(x)}. \quad (2.94)$$

We prove Gibbs' inequality using Jensen's inequality. Let $f(u) = \log 1/u$ and $u = \frac{Q(x)}{P(x)}$. Then

$$D_{\text{KL}}(P||Q) = \mathbb{E}[f(Q(x)/P(x))] \quad (2.95)$$

$$\geq f\left(\sum_x P(x) \frac{Q(x)}{P(x)}\right) = \log\left(\frac{1}{\sum_x Q(x)}\right) = 0, \quad (2.96)$$

with equality only if $u = \frac{Q(x)}{P(x)}$ is a constant, that is, if $Q(x) = P(x)$. \square

The relative entropy or Kullback–Leibler divergence between two probability distributions $P(x)$ and $Q(x)$ that are defined over the same alphabet \mathcal{A}_X is

$$D_{\text{KL}}(P||Q) = \sum_x P(x) \log \frac{P(x)}{Q(x)}. \quad (2.45)$$

The relative entropy satisfies Gibbs' inequality

$$D_{\text{KL}}(P||Q) \geq 0 \quad (2.46)$$

2nd

Second solution. In the above proof the expectations were with respect to the probability distribution $P(x)$. A second solution method uses Jensen's inequality with $Q(x)$ instead. We define $f(u) = u \log u$ and let $u = \frac{P(x)}{Q(x)}$. Then

$$D_{\text{KL}}(P||Q) = \sum_x Q(x) \frac{P(x)}{Q(x)} \log \frac{P(x)}{Q(x)} = \sum_x Q(x) f\left(\frac{P(x)}{Q(x)}\right) \quad (2.97)$$

$$\geq f\left(\sum_x Q(x) \frac{P(x)}{Q(x)}\right) = f(1) = 0, \quad (2.98)$$

with equality only if $u = \frac{P(x)}{Q(x)}$ is a constant, that is, if $Q(x) = P(x)$. \square

8.2.1 The Kullback-Leibler Divergence $\text{KL}(q|p)$

The Kullback-Leibler divergence $\text{KL}(q|p)$ measures the ‘difference’ between distributions q and p [68].

Definition 8.11 (KL divergence). For two distributions $q(x)$ and $p(x)$

$$\text{KL}(q|p) \equiv \langle \log q(x) - \log p(x) \rangle_{q(x)} \geq 0 \quad (8.2.28)$$

where $\langle f(x) \rangle_{r(x)}$ denotes average of the function $f(x)$ with respect to the distribution $r(x)$.

3rd Proof (Chap8 Barber)

The KL divergence is ≥ 0

The KL divergence is widely used and it is therefore important to understand why the divergence is positive.

To see this, consider the following linear bound on the function $\log(x)$

$$\log(x) \leq x - 1 \quad (8.2.29)$$

as plotted in the figure on the right. Replacing x by $p(x)/q(x)$ in the above bound

$$\frac{p(x)}{q(x)} - 1 \geq \log \frac{p(x)}{q(x)} \quad (8.2.30)$$

Since probabilities are non-negative, we can multiply both sides by $q(x)$ to obtain

$$p(x) - q(x) \geq q(x) \log p(x) - q(x) \log q(x) \quad (8.2.31)$$

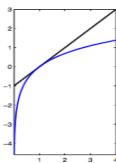
We now integrate (or sum in the case of discrete variables) both sides. Using $\int p(x)dx = 1$, $\int q(x)dx = 1$,

$$1 - 1 \geq \langle \log p(x) - \log q(x) \rangle_{q(x)} \quad (8.2.32)$$

Rearranging gives

$$\langle \log q(x) - \log p(x) \rangle_{q(x)} \equiv \text{KL}(q|p) \geq 0 \quad (8.2.33)$$

The KL divergence is zero if and only if the two distributions are exactly the same.



How to use KL Divergence

$$\begin{aligned} \text{KL}(p\|q) &= - \int p(\mathbf{x}) \ln q(\mathbf{x}) d\mathbf{x} - \left(- \int p(\mathbf{x}) \ln p(\mathbf{x}) d\mathbf{x} \right) \\ &= - \int p(\mathbf{x}) \ln \left\{ \frac{q(\mathbf{x})}{p(\mathbf{x})} \right\} d\mathbf{x} \end{aligned}$$

↓

This step is an approximation that gets better and better as N gets larger and larger.

$$\text{KL}(p\|q) \simeq \frac{1}{N} \sum_{n=1}^N \{-\ln q(\mathbf{x}_n|\boldsymbol{\theta}) + \ln p(\mathbf{x}_n)\}$$

$$\text{KL}(p\|q) \geq 0$$

Howto “Try my best!”

Variational Method (Mean Field Approximation) etc ...