

LOGARITHMS

Defined implicitly by the equation

$$x = b^{\log_b(x)}$$

That is, the inverse of the exponential function. Many useful properties:

$$\begin{aligned} \log_b(xy) &= \log_b(x) + \log_b(y) \\ \log_b(x^r) &= r\log_b(x) \\ \log_b(x) &= \frac{\log_a(x)}{\log_a(b)} \\ \log_a(1) &= 0 \end{aligned}$$

POTENTIAL PITFALL: $\log_a(x+y)$ is in general not simplifiable.

CONCAVITY AND CONVEXITY

EXAMPLE: CONCAVITY OF THE LOGARITHM

The log and $-x \log x$ are concave-down (\cap) functions.

PROBABILITY DISTRIBUTIONS

THEOREM: JENSEN'S INEQUALITY

Given a convex-down (\cup) function $f(x)$ and a PDF $p(x)$.

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

Entropy

IMPORTANT DEFINITIONS

Consider a RV X taking values in an **alphabet** \mathcal{X} .

DEFINITION: INFORMATION CONTENT

The information content (a.k.a. surprise or surprisal) of a realisation x is

$$h(x) := \log \frac{1}{p(x)}$$

We will sometimes use the equivalent notation $h(p) := -\log p$.

DEFINITION: ENTROPY

The entropy of a RV X is given by

$$H(X) := \sum_{x \in \mathcal{X}} p(x) \log \frac{1}{p(x)}$$

Note that entropy is the expected value of surprisal: $H(X) = E[h(X)] = \sum_x p(x)h(x)$.

ENTROPY OF A BIASED COIN

► Mathematically, a biased coin flip X is a **Bernoulli distribution** with parameter p : $X \sim \mathcal{B}(p)$

► The units of $h(x)$ depend on the base of the log:

- \log_2 yields units of bits.
- \log_e yields units of nats.

THE TERM "BIT"

We use the term **bit** to denote two things:

- The units of information content, taken with \log_2 .
- A random variable with alphabet {0, 1}.

Details will be clear from context.

► Entropy of discrete distributions is non-negative. Entropy is additive for independent random variables.

DEFINITION: SHANNON AXIOMS

The survival of an event with probability p , $i(p)$ must satisfy:

1. Certain events are unsurprising: $i(1) = 0$.
2. Less probable events are more surprising: $i(p) \leq i(q)$.
3. Independent events yield the sum of their surprises: $i(p+q) = i(p) + i(q)$.

The only function satisfying these axioms is the negative logarithm:

$$i(p) = h(p) = -\log_2(p)$$

CONCLUSIONS

- For uniform distributions, $h(x)$ corresponds to the length of the binary representation of x (i.e. number of yes/no questions).
- $h(x)$ is a measure of the **intrinsic** information content of x , regardless of how that information is obtained.
- Information is given by probability mass exclusions (Hartley 1928).

The goal of experiments is to reduce uncertainty about something – The value of a parameter, the best model for some dataset, etc

DEFINITION: DIFFERENTIAL ENTROPY

The differential entropy of a continuous RV X with PDF $f(x)$ is given by

$$H(X) := -\int f(x) \log f(x) dx$$

EXAMPLE: DIFFERENTIAL ENTROPY IN UNIFORM DISTRIBUTIONS

Let X be uniform RV in the interval of length a .

$$X \sim U([0, a]) \quad \text{i.e. } p(x) = \frac{1}{a}$$

$$H(X) = -\int_0^a \frac{1}{a} \log \frac{1}{a} dx = -\log \frac{1}{a} \int_0^a dx = \log a$$

► Conclusions: Differential entropy...

- can be negative (e.g. for $a < 1$).
- grows with the volume of the distribution ($2^{H(X)} = a > 0$).

ENTROPY IN GAUSSIAN DISTRIBUTIONS

For the second term:

$$\begin{aligned} \mathbb{E}[(x - \mu)^T \Sigma^{-1} (x - \mu)] &= \text{tr}(\Sigma^{-1} \mathbb{E}[(x - \mu)(x - \mu)^T]) \\ &= \text{tr}(\Sigma^{-1} \Sigma) \\ &= D \end{aligned}$$

Back to the main expression:

$$H(X) = \frac{1}{2} \ln(2\pi e^2) + \frac{1}{2} D = \frac{1}{2} \ln((2\pi e)^2)$$

Overall:

$$H(X) = \frac{1}{2} \ln(2\pi e^2)$$

► Entropy has a closed-form expression for Gaussian distributions.

KEY TAKEAWAY: ENTROPY

- Entropy quantifies the average information content obtained from an observation x .
- Entropy is a **generalised variance** (e.g. how hard is it to predict X).
- In discrete distributions, entropy...
- is bounded $0 \leq H(X) \leq \log |\mathcal{X}|$.
- is related to the number of yes/no questions needed to determine x .
- In continuous distributions, differential entropy is defined, but does not satisfy some properties of discrete entropy (e.g. can be negative).

DEFINITION: CODE AND CODE LENGTH

Given a RV X with alphabet \mathcal{X} and an alphabet \mathcal{D} , a **code** is a mapping $C: \mathcal{X} \rightarrow \mathcal{D}^*$, where \mathcal{D}^* is the set of all finite-length strings of symbols in \mathcal{D} .

The quantity $l(x)$ is the **code length** of $C(x)$, and $L = E[l(x)]$ the average code length.

DEFINITION: EXTENDED CODE

We've seen that (in uniform distributions), entropy corresponds to the number of **yes/no questions** needed to guess x .

$35 \rightarrow 100011 \quad 6 \rightarrow 000110 \quad 17 \rightarrow 01001$

► This forms a **binary code** for X with length

$$L = H(X) = \log |\mathcal{X}|$$

► More generally, $\log |\mathcal{X}|$ is referred to as the **raw bit content** of X .

THEOREM: AEP

Let X_1, \dots, X_n be i.i.d. RVs with $X_i \sim p(x)$. Then

$$\lim_{n \rightarrow \infty} -\frac{1}{n} \log \prod_{i=1}^n p(X_i, \dots, X_n) = H(X)$$

All samples are equally probable!

► Proof:

$$\begin{aligned} -\frac{1}{n} \log p(X_1, \dots, X_n) &= -\frac{1}{n} \log \prod_{i=1}^n p(x_i) = -\frac{1}{n} \sum_{i=1}^n \log p(x_i) \\ &\stackrel{\text{AEP}}{\rightarrow} -E[\log p(x)] = H(X) \end{aligned}$$

DEFINITION: TYPICAL SET

The typical set $A_{\varepsilon}^{(n)}$ for some $\varepsilon > 0$ is the set of strings x_1, \dots, x_n with the property

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

► Informally: events that almost always occur when sampling X_1, \dots, X_n . In other words, events outside $A_{\varepsilon}^{(n)}$ are very unlikely.

THEOREM: KRAFT INEQUALITY

The codeword lengths l_i of any prefix-free code must satisfy

$$\sum_i l_i \leq 1$$

Conversely, given a set of lengths that satisfy this inequality, there exists a prefix code with those lengths.

DEFINITION: ASYMPTOTIC EQUIPARTITION PRINCIPLE

The AEP says that if X is a discrete RV with $X_i \sim p(x)$, then

$$\lim_{n \rightarrow \infty} -\frac{1}{n} \log \prod_{i=1}^n p(X_i, \dots, X_n) = H(X)$$

All samples are equally probable!

DEFINITION: UNIQUELY DECODEABLE CODE

A code C is uniquely decodable if its extension is nonsingular.

DEFINITION: PREFIX CODE

A prefix (or prefix-free, or instantaneous, or self-puncturing) code is one where no codeword is the start of another codeword.

EXAMPLE

	$C(x)$
x	0
a	10
b	110

DEFINITION: NON SINGULAR

An **invertible** code is a code where no codeword is the start of another codeword.

DEFINITION: MAXIMUM LIKELIHOOD ESTIMATION

Given a dataset X and a family of PDFs $p(x, \theta)$, the **maximum likelihood estimate** of θ is

$$\hat{\theta}_{MLE} = \arg\max_{\theta} p(x, \theta)$$

DEFINITION: KL DIVERGENCE

The **Kullback-Leibler (KL) divergence** is a non-negative measure of the difference between two PDFs.

- Compression interpretation: the "extra cost" of using a wrong code.
- Geometric interpretation: grows as the two distributions get far.

DEFINITION: ENTROPY RATE

A stochastic process is an index set of random variables $\{X_1, \dots, X_n\}$. Here we will consider semi-infinite processes $\{X_1, \dots, X_n, \dots\}$. A stochastic process is stationary if for all t and s ,

$$p(X_t, X_{t+1}, \dots, X_s) = p(X_1, X_{2+1}, \dots, X_{s+1})$$

► i.e. a stochastic process is stationary if it's **invariant to time shift**. Let's introduce two interesting quantities:

- **Entropy rate**: how quickly $H(X_1, \dots, X_n)$ grows as n increases.
- $H(X) = \lim_{n \rightarrow \infty} \frac{1}{n} H(X_1, \dots, X_n)$
- **Innovation**: how much new information is introduced at time n .
- $h(X) = \lim_{n \rightarrow \infty} H(X|X_1, \dots, X_{n-1})$

THEOREM: ENTROPY RATE

For stationary stochastic processes, these two quantities are equal:

$$H(X) = \lim_{n \rightarrow \infty} \frac{1}{n} H(X_1, \dots, X_n) = \lim_{n \rightarrow \infty} H(X_1, X_{2+1}, \dots, X_{n+1})$$

STOCHASTIC PROCESSES

The more things you condition on, the less entropy you have:

$$H(X_1, X_2, \dots, X_n) \leq H(X_1, X_2, \dots, X_{n-1})$$

Intuitively: X is a time series, and $H(X_1, \dots, X_{n-1})$ is the performance of a predictor with $n-1$ timesteps of "memory".

Seen in the system's entropy convergence curve:

- Without memory ($t = 0$), we get the usual $H(X_1)$.
- As we include more memory, conditional entropy decreases...
- ...Until it reaches a minimum at the **entropy rate** ($H(X)$).
- If it happens at some finite t , this is the **Markov order** p .

THEOREM: ENTROPY

For stationary stochastic processes, these two quantities are equal:

$$h(X) = \lim_{n \rightarrow \infty} \frac{1}{n} H(X_1, \dots, X_n) = \lim_{n \rightarrow \infty} H(X_1 | X_2, \dots, X_{n-1})$$

MUTUAL INFORMATION

MULTIVARIATE PROBABILITIES

The marginal probability $p(X = x)$ is the probability of x happening, regardless of y :

$$p(X = x) = \sum_{y \in \mathcal{Y}} p(X = x, Y = y)$$

The conditional probability $p(X = x | Y = y)$ is the probability of x happening, given we know y happens:

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

CONDITIONING ON IMPOSSIBLE EVENTS

If $p(Y = y) = 0$, then $p(X | Y = y)$ is undefined.

JOINT PROBABILITY RECAP

Joint probabilities obey some fundamental properties:

- **Product rule**:
$$p(x, y) = p(x)p(y)$$
- **Sum rule**:
$$p(x) = \sum_{y \in \mathcal{Y}} p(x, y)$$
- **Bayes' theorem**:
$$p(x | y) = \frac{p(x, y)}{p(y)}$$

WARNING: NOTATION OVERLOAD

Technically, the symbol p in the expressions $p(x)$ and $p(y)$ represents two different functions. The precise way to write would be $p_X(x) = p(x)$. When it is unambiguous from the context, we will use simply $p(x)$.

GRAPHICAL MODELS

► Let's look at a particularly interesting Bayesian network:

Assign codewords by traversing the tree from right to left:

	a → 00	b → 10	...	e → 011
--	--------	--------	-----	---------

KEY TAKEAWAY: SYMBOL CODES

► **Symbol codes** map each individual symbol onto a codeword.

► **Prefix codes** are those where no codeword is the start of another codeword. They can be represented with trees.

► **Kraft inequality**: prefix code $\rightarrow \sum_i l_i \leq 1$.

► **Source coding theorem for symbol codes**. For any $p(x)$:

- There exists a code that achieves length close to $H(X)$. (Achievability)
- There are no codes that achieve length $< H(X)$. (Converse)

DEFINITION: CROSS ENTROPY

The cross-entropy of a PDF q relative to a PDF p is defined as:

$$H(p, q) := -\sum_{x \in \mathcal{X}} p(x) \log q(x)$$

It represents the expected length of a q -optimal code used on samples from p .

DEFINITION: JOINT ENTROPY

The joint entropy of two RVs X and Y with alphabets \mathcal{X} and \mathcal{Y} , is given by

$$H(X, Y) = -\sum_{x \in \mathcal{X}, y \in \mathcal{Y}} p(x, y) \log p(x, y)$$

All the properties of entropy apply to joint entropy:

- Expected value of joint survival, $\mathbb{E}[p(X, Y)] = \int p(x, y) dx dy$.
- Uncertainty when predicting X and Y jointly.
- Minimum code length when encoding X and Y simultaneously.
- Note that joint entropy is **symmetric**: $H(X, Y) = H(Y, X)$.

DEFINITION: CONDITIONAL ENTROPY

The conditional entropy of X given Y is given by

$$H(X|Y) = -\sum_{y \in \mathcal{Y}} p(y) H(X|Y=y),$$

where

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

WARNING: EMPTY CONDITIONAL

Written this way, $H(X|Y)$ is not defined when $p(Y) = 0$ for any y . To fix this, we can either take the sum over y with $p(y) > 0$, or define $\delta_0^0 = 0$.

CONDITIONAL ENTROPY

Some properties of conditional entropy in discrete RVs:

- $H(X|Y) \geq 0$.
- $H(X|Y) = 0 \iff p(X|Y) = 1$. Since 0 $\leq p(X|Y) \leq 1$, $-\log p(X|Y) \geq 0$.
- $H(X|Y) \leq \log |\mathcal{X}|$.
- Proof: Similar to the proof of $H(X) \leq \log |\mathcal{X}|$.
- $H(X|Y) = 0 \iff X$ is a deterministic function of Y .
- Proof: Similar to the proof of $H(X) = 0$ for the Kronecker delta.
- (Applies to discrete and continuous). If $X \perp Y$, then $H(X|Y) = H(X)$.
- Proof: If $X \perp Y$, then $H(X|Y) = -\sum_{y \in \mathcal{Y}} p(y) \log p(y) = -\sum_{y \in \mathcal{Y}} p(y) \log p(X|Y=y) = H(X)$.

DEFINITION: CONDITIONAL MUTUAL INFORMATION

The conditional mutual information between X and Y given Z is defined as

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

CMI can be rewritten in terms of conditional entropy:

$$I(X; Y | Z) = H(X|Z) - H(X|Y|Z) = H(X|Z) - H(X|Y)$$

► In other words: MI measures the average reduction of uncertainty in X after knowing Z (or vice versa).

► You can't get out more than you put in: $I(X; Y) \leq \min(H(X), H(Y))$.

- Follows from $H(X|Y) \geq 0$ (discrete variables only).
- Follows from $H(X|Y) \geq 0$ (continuous variables only).

DEFINITION: INFORMATION DOESN'T HURT

Conditioning reduces entropy:

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

PROOF:

► Note: this holds for continuous variables, although $H(X)$ and $H(X|Y)$ can both be negative.

► In general, $H(X|Y) \geq H(X)$.

DEFINITION: CONDITIONAL ENTROPY

The conditional entropy of X given Y is given by

$$H(X|Y) = -\sum_{y \in \mathcal{Y}} p(y) H(X|Y=y),$$

where

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

DEFINITION: ENTROPY CHAIN RULE

The entropy chain rule states that

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

DEFINITION: DATA PROCESSING INEQUALITY

If $X - Y - Z$ form a Markov chain,

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

DEFINITION: INVARIANCE OF MUTUAL INFORMATION

Mutual information $I(X; Y)$ is invariant under isomorphisms in X and Y .

DEFINITION: CONNECTIONS WITH SET THEORY

The relationship between H and I can be understood with a **Venn diagram**:

This motivates a **set-theoretic interpretation of entropy**:

- $H(X) \leftrightarrow \mu(\mathcal{A})$
- $H(X, Y) \leftrightarrow \mu(\mathcal{A} \cup \mathcal{B})$
- $H(X|Y) \leftrightarrow \mu(\mathcal{A} \setminus \mathcal{B})$
- $I(X; Y) \leftrightarrow \mu(\mathcal{A} \cap \mathcal{B})$

Tremendously useful property.

