

CO 432 Spring 2025:

Lecture Notes

1	Introduction	2
1.1	Entropy	2
1.2	Entropy as expected surprise	3
1.3	Entropy as optimal lossless data compression	4
2	Relative entropy	8
2.1	The boolean k -slice	11
2.2	Rejection sampling	14
3	Mutual information	21
3.1	Definition and chain rules	21
3.2	Markov chains, data processing, and sufficient statistics	24
3.3	Communication complexity	25
3.4	Parameter estimation	27
	Back Matter	28
	List of Named Results	28
	Index of Defined Terms	29

Lecture notes taken, unless otherwise specified, by myself during the Spring 2025 offering of CO 432, taught by Vijay Bhattiprolu.

Lectures			Lecture 6	May 22	14
			Lecture 7	May 27	16
Lecture 1	May 6	2	Lecture 8	May 29	18
Lecture 2	May 8	4	Lecture 9	June 3	21
Lecture 3	May 13	6	Lecture 11	June 10	24
Lecture 4	May 15	8	Lecture 10	June 5	25
Lecture 5	May 20	11			

Chapter 1

Introduction

Notation. I will be using my usual L^AT_EX typesetting conventions:

- $[n]$ means the set $\{1, 2, \dots, n\}$
- $\{0, 1\}^*$ means the set of bitstrings of arbitrary length (i.e., the Kleene star)
- \sum_i is implicitly $\sum_{i=1}^n$
- A, B, \dots, Z are random variables (in sans-serif)
- $X \sim (p_1, p_2, \dots, p_n)$ means X is a discrete random variable with n outcomes such that $\Pr[X = 1] = p_1, \Pr[X = 2] = p_2$, etc. (abbreviate further as $X \sim (p_i)$)

1.1 Entropy

↓ Lecture 1 adapted from Arthur ↓

Lecture 1
May 6

Definition 1.1.1 (entropy)

For a random variable $X \sim (p_i)$, the entropy $H(X)$ is

$$H(X) = - \sum_i p_i \log p_i = \sum_i p_i \log \frac{1}{p_i}.$$

Convention. By convention, we usually use \log_2 . Also, we define entropy such that $\log_2(0) = 0$ so that impossible values do not break the formula.

Example 1.1.2. If X takes on the values a, b, c, d with probabilities $1, 0, 0, 0$, respectively, then $H(X) = 1 \log 1 = 0$.

If X takes on those values instead with probabilities $\frac{1}{2}, \frac{1}{4}, \frac{1}{8}, \frac{1}{8}$, respectively, then $H(X) = \frac{7}{4}$.

Fact 1.1.3. $H(\mathbf{X}) = 0$ if and only if \mathbf{X} is a constant.

Proof. Suppose \mathbf{X} is constant. Then, $H(\mathbf{X}) = 1 \log 1 = 0$.

Suppose $H(\mathbf{X}) = 0$. Probabilities are in $[0, 1]$, so $p_i \log \frac{1}{p_i} \geq 0$. Since $H(\mathbf{X}) = \sum_i p_i \log \frac{1}{p_i} = 0$ and each term is non-negative, each term must be zero. Thus, each p_i is either 0 or 1. We cannot have $\sum p_i > 1$, so exactly one $p_i = 1$ and the rest are zero. That is, \mathbf{X} is constant. \square

Theorem 1.1.4 (Jensen's inequality)

Let $f : \mathbb{R} \rightarrow \mathbb{R}$ be concave. That is, for any a and b in the domain of f and $\lambda \in [0, 1]$, $f(\lambda a + (1 - \lambda)b) \geq \lambda f(a) + (1 - \lambda)f(b)$. For any discrete random variable \mathbf{X} ,

$$\mathbb{E}[f(\mathbf{X})] \leq f(\mathbb{E}[\mathbf{X}])$$

Proof. Consider a random variable \mathbf{X} with two values a and b , each with probabilities λ and $1 - \lambda$. Then, notice that

$$\mathbb{E}[f(\mathbf{X})] = \lambda f(a) + (1 - \lambda)f(b) \leq f(\lambda a + (1 - \lambda)b) = f(\mathbb{E}[\mathbf{X}])$$

by convexity of f .

TODO: This can be generalized by induction. \square

Fact 1.1.5. Assume \mathbf{X} is supported on $[n]$. Then, $0 \leq H(\mathbf{X}) \leq \log n$.

Proof. Start by claiming without proof that $\log n$ is concave, so we can apply [Jensen's inequality](#).

Let $\mathbf{X}' = \frac{1}{p_i}$ with probability p_i . Then,

$$\begin{aligned} H(\mathbf{X}) &= \sum_i p_i \log \frac{1}{p_i} \\ &= \mathbb{E}[\log(\mathbf{X}')] \\ &\leq \log(\mathbb{E}[\mathbf{X}']) \\ &= \log\left(\sum p_i \frac{1}{p_i}\right) \\ &= \log n \end{aligned}$$

\square

It is not a coincidence that $\log_2 n$ is the minimum number of bits to encode $[n]$.

1.2 Entropy as expected surprise

We want $S : [0, 1] \rightarrow [0, \infty)$ to capture how “surprised” we are $S(p)$ that an event with probability p happens. We want to show that under some natural assumptions, this is the only function we could have defined as entropy. In particular:

1. $S(1) = 0$, a certainty should not be surprising
2. $S(q) > S(p)$ if $p > q$, less probable should be more surprising
3. $S(p)$ is continuous in p
4. $S(pq) = S(p) + S(q)$, surprise should add for independent events. That is, if I see something twice, I should be twice as surprised.

↑ Lecture 1 adapted from Arthur ↑

Lecture 2
May 8

Proposition 1.2.1

If $S(p)$ satisfies these 4 axioms, then $S(p) = c \cdot \log_2(1/p)$ for some $c > 0$.

Proof. Suppose a function $S : [0, 1] \rightarrow [0, \infty)$ exists satisfying the axioms. Let $c := S(\frac{1}{2}) > 0$.

By axiom 4 (addition), $S(\frac{1}{2^k}) = kS(\frac{1}{2})$. Likewise, $S(\frac{1}{2^{1/k}} \cdots \frac{1}{2^{1/k}}) = S(\frac{1}{2^{1/k}}) + \cdots + S(\frac{1}{2^{1/k}}) = kS(\frac{1}{2^{1/k}})$.

Then, $S(\frac{1}{2^{m/n}}) = \frac{m}{n}S(\frac{1}{2}) = \frac{m}{n} \cdot c$ for any rational m/n .

By axiom 3 (continuity), $S(\frac{1}{2^z}) = c \cdot z$ for all $z \in [0, \infty)$ because the rationals are dense in the reals. In particular, for any $p \in [0, 1]$, we can write $p = \frac{1}{2^z}$ for $z = \log_2(1/p)$ and we get

$$S(p) = S\left(\frac{1}{2^z}\right) = c \cdot z = c \cdot \log_2(1/p)$$

as desired. □

We can now view entropy as expected surprise. In particular,

$$\sum_i p_i \log_2 \frac{1}{p_i} = \mathbb{E}_{x \sim \mathbf{X}} [S(p_x)]$$

for a random variable $\mathbf{X} = i$ with probability p_i .

1.3 Entropy as optimal lossless data compression

Suppose we are trying to compress a string consisting of n symbols drawn from some distribution.

Problem 1.3.1

What is the expected number of bits you need to store the results of n independent samples of a random variable \mathbf{X} ?

We will show this is $nH(\mathbf{X})$.

Notice that we assume that the symbols we are drawn independently, which is violated by almost all data we actually care about.

Definition 1.3.2

Let $C : \Sigma \rightarrow (\Sigma')^*$ be a code. We say C is a uniquely decodable code (UDC) if there does not exist a collision $x, y \in \Sigma^*$, with identical encoding $C(x_1)C(x_2) \cdots C(x_k) = C(y_1)C(y_2) \cdots C(y_{k'})$.

Also, C is prefix-free (sometimes called instantaneous) if for any distinct $x, y \in \Sigma$, $C(x)$ is not a prefix of $C(y)$.

Proposition 1.3.3

Prefix-freeness is sufficient for unique decodability.

Example 1.3.4. Let $C : \{A, B, C, D\} \rightarrow \{0, 1\}^*$ where $C(A) = 11$, $C(B) = 101$, $C(C) = 100$, and $C(D) = 00$. Then, C is prefix-free and uniquely decodable.

We can easily parse 1011100001100 unambiguously as 101.11.00.00.11.00 (*BADDAD*).

Recall from CS 240 that a prefix-free code is equivalent to a trie, and we can decode it by traversing the trie in linear time.

Theorem 1.3.5 (Kraft's inequality)

A prefix-free binary code $C : \{1, \dots, n\} \rightarrow \{0, 1\}^*$ with codeword lengths $\ell_i = |C(i)|$ exists if and only if

$$\sum_{i=1}^n \frac{1}{2^{\ell_i}} \leq 1.$$

Proof. Suppose $C : \{1, \dots, n\} \rightarrow \{0, 1\}^*$ is prefix-free with codeword lengths ℓ_i . Let T be its associated binary tree and let W be a random walk on T where 0 and 1 have equal weight (stopping at either a leaf or undefined branch).

Define E_i as the event where W reaches i and E_\emptyset where W falls off. Then,

$$\begin{aligned} 1 &= \Pr(E_\emptyset) + \sum_i \Pr(E_i) \\ &= \Pr(E_\emptyset) + \sum_i \frac{1}{2^{\ell_i}} && \text{(by independence)} \\ &\geq \sum_i \frac{1}{2^{\ell_i}} && \text{(probabilities are non-negative)} \end{aligned}$$

Conversely, suppose the inequality holds for some ℓ_i . WLOG, suppose $\ell_1 < \ell_2 < \dots < \ell_n$.

Start with a complete binary tree T of depth ℓ_n . For each $i = 1, \dots, n$, find any unassigned node in T of depth ℓ_i , delete its children, and assign it a symbol.

Now, it remains to show that this process will not fail. That is, for any loop step i , there is still some unassigned node at depth ℓ_i .

Let $P \leftarrow 2^{\ell_n}$ be the number of leaves of the complete binary tree of depth ℓ_n . After the i^{th} step, we decrease P by $2^{\ell_n - \ell_i}$. That is, after n steps,

$$\begin{aligned} P &= 2^{\ell_n} - \sum_{i=1}^n \frac{2^{\ell_n}}{2^{\ell_i}} \\ &= 2^{\ell_n} - 2^{\ell_n} \sum_{i=1}^n \frac{1}{2^{\ell_i}} \\ &\geq 0 \end{aligned}$$

by the inequality. □

Recall the problem we are trying to solve:

Lecture 3
May 13

Problem 1.3.1

What is the expected number of bits you need to store the results of n independent samples of a random variable \mathbf{X} ?

Solution (Shannon & Fano). Consider the case where \mathbf{X} is symbol i with probability p_i . We want to encode independent samples $x_i \sim \mathbf{X}$ as $C(x_i)$ for some code $C : [n] \rightarrow \{0, 1\}^*$.

Suppose for simplification that $p_i = \frac{1}{2^{\ell_i}}$ for some integers ℓ_i . Since $\sum p_i = 1$, we must have $\sum \frac{1}{2^{\ell_i}} = 1$. Then, by [Kraft's inequality](#), there exists a prefix-free binary code $C : [n] \rightarrow \{0, 1\}^*$ with codeword lengths $|C(i)| = \ell_i$. Now,

$$\mathbb{E}_{x_i \sim \mathbf{X}} \left[\sum_i |C(x_i)| \right] = \sum_i p_i \ell_i = \sum_i p_i \log_2 \frac{1}{p_i} = H(\mathbf{X})$$

Proceed to the general case. Suppose $\log_2 \frac{1}{p_i}$ are non-integral. Instead, use $\ell'_i = \lceil \log_2 \frac{1}{p_i} \rceil$. We still satisfy Kraft since $\sum_i \frac{1}{2^{\ell'_i}} \leq \sum_i p_i = 1$. Then,

$$\mathbb{E}_{x_i \sim \mathbf{X}} \left[\sum_i |C(x_i)| \right] = \sum_i p_i \ell'_i = \sum_i p_i \left\lceil \log_2 \frac{1}{p_i} \right\rceil$$

which is bounded by

$$H(\mathbf{X}) = \sum_i p_i \log_2 \frac{1}{p_i} \leq \sum_i p_i \left\lceil \log_2 \frac{1}{p_i} \right\rceil < \sum_i p_i \left(1 + \log_2 \frac{1}{p_i} \right) = H(\mathbf{X}) + 1$$

We call the code C generated by this process the Shannon–Fano code. □

We can improve on this bound $[H(\mathbf{X}), H(\mathbf{X}) + 1]$ by amortizing over longer batches of the string.

Solution (batching). For \mathbf{Y} defined on $[n]$ equal to i with probability q_i , define the random variable $\mathbf{Y}^{(k)}$ on $[n]^k$ equal to the string $i_1 \dots i_k$ with probability $q_{i_1} \dots q_{i_k}$. That is, $\mathbf{Y}^{(k)}$ models k independent samples of \mathbf{Y} .

Apply the Shannon–Fano code to $\mathbf{Y}^{(k)}$ to get an encoding of $[n]^k$ as bitstrings of expected length ℓ

satisfying $H(\mathbf{Y}^{(k)}) \leq \ell \leq H(\mathbf{Y}^{(k)}) + 1$.

$$\begin{aligned}
 H(\mathbf{Y}^{(k)}) &= \mathbb{E}_{i_1 \dots i_k \sim \mathbf{Y}^{(k)}} \left[\log_2 \frac{1}{q_{i_1} \dots q_{i_k}} \right] && \text{(by def'n)} \\
 &= \mathbb{E}_{i_1 \dots i_k \sim \mathbf{Y}^{(k)}} \left[\log_2 \frac{1}{q_{i_1}} + \dots + \log_2 \frac{1}{q_{i_k}} \right] && \text{(log rules)} \\
 &= \sum_{j=1}^k \mathbb{E}_{i_1 \dots i_k \sim \mathbf{Y}^{(k)}} \left[\log_2 \frac{1}{q_{i_j}} \right] && \text{(linearity of expectation)} \\
 &= \sum_{j=1}^k \mathbb{E}_{i \sim \mathbf{Y}} \left[\log_2 \frac{1}{q_i} \right] && (q_{i_j} \text{ only depends on one character}) \\
 &= kH(\mathbf{Y}) && \text{(by def'n, no } j\text{-dependence in sum)}
 \end{aligned}$$

For every k symbols, we use ℓ bits, i.e., $\frac{\ell}{k}$ bits per symbol. From the Shannon–Fano bound, we have

$$\begin{aligned}
 \frac{H(\mathbf{Y}^{(k)})}{k} &\leq \frac{\ell}{k} < \frac{H(\mathbf{Y}^{(k)})}{k} + \frac{1}{k} \\
 H(\mathbf{Y}) &\leq \frac{\ell}{k} < H(\mathbf{Y}) + \frac{1}{k}
 \end{aligned}$$

Then, we have a code for \mathbf{Y} bounded by $[H(\mathbf{Y}), H(\mathbf{Y}) + \frac{1}{k}]$.

Taking a limit of some sort, we can say that we need $H(\mathbf{Y}) + o(1)$ bits. □

Chapter 2

Relative entropy

Definition 2.0.1 (relative entropy)

Given two discrete distributions $p = (p_i)$ and $q = (q_i)$, the relative entropy

$$D(p \parallel q) := \sum p_i \log_2 \frac{1}{q_i} - \sum_i p_i \log_2 \frac{1}{p_i} = \sum p_i \log_2 \frac{p_i}{q_i}$$

This is also known as the KL divergence.

The KL divergence works vaguely like a “distance” between distributions. (In particular, KL divergence is not a metric since it lacks symmetry and does not follow the triangle inequality, but it can act sorta like a generalized squared distance.)

*Lecture 4
May 15*

Fact 2.0.2. $D(p \parallel q) \geq 0$ with equality exactly when $p = q$.

Proof. Observe that

$$-D(p \parallel q) = \sum_i p_i (-\log_2 \frac{p_i}{q_i}) = \sum_i p_i \log_2 \frac{q_i}{p_i}$$

and then define $X' = \frac{q_i}{p_i}$ with probability p_i . By construction, we get

$$-D(p \parallel q) = \mathbb{E}[\log_2 X'] \leq \log_2(\mathbb{E}[X'])$$

by [Jensen's inequality](#) (as $f = \log_2$ is concave). Finally,

$$D(p \parallel q) \geq -\log_2(\mathbb{E}[X']) = \log_2 \left(\sum_i p_i \frac{q_i}{p_i} \right) = \log_2 1 = 0$$

□

Proposition 2.0.3

Any prefix-free code has an expected length at least $H(X)$.

Proof. Let $X \sim (p_i)$. Suppose C is a prefix-free code with codeword lengths ℓ_i .

Then, by [Kraft's inequality](#), $\sum_i 2^{-\ell_i} \leq 1$. We want to show that $\sum_i p_i \ell_i \geq H(\mathbf{X})$, and we will prove this by showing that $\sum_i p_i \ell_i - H(\mathbf{X}) = D(p \parallel q)$ for some distribution q (then apply fact [2.0.2](#)).

We will take q to be the random walk distribution corresponding to the binary tree associated to the candidate prefix-free code.

Let T be the binary tree associated to C . Consider the process of randomly going left/right at each node and stopping when either falling off the tree or hitting a leaf.

Let $q_i = 2^{-\ell_i}$ be the probability that this random walk reaches the leaf for the symbol i and let $q_{n+1} = 1 - \sum_i 2^{-\ell_i}$ be the probability that the random walk falls off the tree. Also, to keep ranges identical, let $\tilde{p}_i = p_i$ and $\tilde{p}_{n+1} = 0$. Now,

$$\begin{aligned} D(\tilde{p} \parallel q_C) &= \sum_{i=1}^{n+1} \tilde{p}_i \log_2 q_i^{-1} - \sum_{i=1}^{n+1} \tilde{p}_i \log_2 \frac{1}{p_i} \\ &= \sum_{i=1}^n p_i \log_2 2^{\ell_i} - \sum_{i=1}^n p_i \log_2 \frac{1}{p_i} \quad (\tilde{p}_{n+1} = 0) \\ &= \sum_{i=1}^n p_i \ell_i - H(\mathbf{X}) \end{aligned}$$

Therefore, by fact [2.0.2](#), $\sum_i p_i \ell_i \geq H(\mathbf{X})$. □

This proof technique generalizes. Recall the distinction between UDCs and prefix-free codes:

Definition 1.3.2

Let $C : \Sigma \rightarrow (\Sigma')^*$ be a code. We say C is a uniquely decodable code (UDC) if there does not exist a collision $x, y \in \Sigma^*$, with identical encoding $C(x_1)C(x_2) \cdots C(x_k) = C(y_1)C(y_2) \cdots C(y_{k'})$.

Also, C is prefix-free (sometimes called instantaneous) if for any distinct $x, y \in \Sigma$, $C(x)$ is not a prefix of $C(y)$.

Example 2.0.4. The code $C(1, 2, 3, 4) = (10, 00, 11, 110)$ is a uniquely decodable code.

The code $C'(1, 2, 3, 4) = (0, 10, 110, 111)$ is a prefix-free code.

Remark 2.0.5. A natural additional requirement for unique decodability is that for any $k \in \mathbb{N}$, $x \in [n]^k$, $y \in [n]^k$, $C(x) \neq C(y)$.

Theorem 2.0.6

For any uniquely decodable code $C : [n] \rightarrow \{0, 1\}^*$ of codeword lengths ℓ_i , there is also a prefix-free code $C' : [n] \rightarrow \{0, 1\}^*$ of lengths ℓ_i .

We will show that for any UDC C , the lengths $\sum_i 2^{-\ell_i} \leq 1$. Then, [Kraft's inequality](#) applies and we have a prefix-free code C' .

Partition the code's codomain $C([n]) = C_1 \cup C_2 \cup C_3 \cup \dots$ by the length of the codeword $C_j \subseteq \{0, 1\}^j$. We must instead show $\sum_j \frac{|C_j|}{2^j} \leq 1$.

Consider the easy case $C([n]) = C_2 \cup C_3$. If there are no collisions of length 5, we have

$$2 \cdot |C_2| \cdot |C_3| \leq 2^5$$

because every string in $\{xy : x \in C_2, y \in C_3\} \cup \{yx : x \in C_2, y \in C_3\}$ is unique in $\{0, 1\}^5$. That is, $|C_2| \cdot |C_3| \leq 2^4$.

Likewise, if there are no collisions of length $5k$, we get

$$\frac{(2k)!}{k! \cdot k!} \cdot |C_2|^k \cdot |C_3|^k \leq 2^{5k}$$

because the union $\bigcup_{\substack{\alpha \in \{2,3\}^{2k}, \\ \alpha_i=2 \text{ for} \\ k \text{ choices of } i}} C_{\alpha_i}$ consists of only unique strings.

In the limit, by [Sterling's approximation](#),

$$\begin{aligned} \frac{2^{2k}}{\sqrt{k}} \cdot |C_2|^k \cdot |C_3|^k &\leq 2^{5k} \\ |C_2| \cdot |C_3| &\leq \frac{2^5}{2^2} (\sqrt{k})^{1/k} \approx 1 + \mathcal{O}(\log k/k) \end{aligned}$$

I have no idea where this was going.

Proof. Fix a $k \geq 1$. Let $\ell_{max} = \max \ell_i$. Write $C^{(k)}$ to be the set of encoded k -length strings.

Consider the distribution: sample a length m uniformly from the set $[k \cdot \ell_{max}]$. Also, sample a uniformly random string $s \in \{0, 1\}^m$. For each $x \in C^{(k)}$, let E_x be the event where $s = x$.

Now, we can write

$$\sum_{x \in C^{(k)}} \Pr[E_x] \leq 1$$

because the events E_x are mutually exclusive. Then,

$$\begin{aligned} \sum_{x \in C^{(k)}} \frac{1}{k \cdot \ell_{max}} \cdot \frac{1}{2^{\ell(x)}} &\leq 1 \\ \sum_{x \in C^{(k)}} \frac{1}{2^{\ell(x)}} &\leq k \cdot \ell_{max} \end{aligned}$$

where $\ell(x)$ is the length of x . Since summing over each codeword $x \in C$ is the same as summing

over each codeword ℓ_i ,

$$\begin{aligned}
 \left(\sum_i \frac{1}{2^{\ell_i}} \right)^k &= \left(\sum_{x \in C} \frac{1}{2^{\ell(x)}} \right)^k \\
 &= \sum_{x_1, \dots, x_k \in C} \frac{1}{2^{\ell(x_1)}} \cdot \frac{1}{2^{\ell(x_2)}} \cdots \frac{1}{2^{\ell(x_k)}} \\
 &= \sum_{x_1, \dots, x_k \in C} \frac{1}{2^{\ell(x_1) + \ell(x_2) + \cdots + \ell(x_k)}} \\
 &= \sum_{x_1, \dots, x_k \in C} \frac{1}{2^{\ell(x_1 x_2 \cdots x_k)}} \\
 &= \sum_{x \in C^{(k)}} \frac{1}{2^{\ell(x)}}
 \end{aligned}$$

where we can take the last step by uniquely decoding $x_1 x_2 \cdots x_k$ into x . Combining,

$$\begin{aligned}
 \left(\sum_i \frac{1}{2^{\ell_i}} \right)^k &\leq k \cdot \ell_{\max} \\
 \sum_i \frac{1}{2^{\ell_i}} &\leq (k \cdot \ell_{\max})^{\frac{1}{k}} \\
 &\leq 1 + \mathcal{O}\left(\frac{\ell_{\max} \cdot \log_2 k}{k}\right)
 \end{aligned}$$

which tends to 1 as $k \rightarrow \infty$, as desired. □

Notation. Write $H(p)$ to denote $H(\mathbf{X})$ for $\mathbf{X} \sim \text{Bernoulli}(p)$.

That is, $H(p) = p \log_2 \frac{1}{p} + (1-p) \log_2 \frac{1}{1-p}$.

Likewise, write $D(q \parallel p)$ to be $D(\mathbf{Y} \parallel \mathbf{X})$ where $\mathbf{Y} \sim \text{Bernoulli}(q)$.

Lecture 5
May 20

Recall Sterling's approximation (which we have used before):

Theorem 2.0.7 (Sterling's approximation)

$m!$ behaves like $\sqrt{2\pi m} \left(\frac{m}{e}\right)^m \left(1 + \mathcal{O}\left(\frac{1}{m}\right)\right)$

2.1 The boolean k -slice

Consider the boolean k -slice (also known as the Hamming k -slice) of the hypercube $\{0, 1\}^n$ defined by

$$B_k := \{x \in \{0, 1\}^n : x \text{ has exactly } k \text{ ones}\}$$

Remark 2.1.1.

$$|B_k| \approx 2^{H(\frac{k}{n}) \cdot n}$$

Proof. By [Sterling's approximation](#), knowing that $|B_k| = \binom{n}{k}$:

$$\begin{aligned} |B_k| &= \binom{n}{k} \\ &= \frac{n!}{k!(n-k)!} \\ &\approx \frac{\sqrt{2\pi n} \left(\frac{n}{e}\right)^n}{\sqrt{2\pi k} \left(\frac{k}{e}\right)^k \sqrt{2\pi(n-k)} \left(\frac{n-k}{e}\right)^{n-k}} \\ &= \sqrt{\frac{n}{2\pi k(n-k)}} \cdot \frac{n^k \left(\frac{n}{n-k}\right)^{n-k}}{k^k} \end{aligned}$$

Now, notice that $\left(\frac{n}{n-k}\right)^{n-k} = \left(1 + \frac{k}{n-k}\right)^{n-k} \approx e^k$ for $k \ll n-k$ because $1+x \approx e^x$ for small x . Then, $\left(1 + \frac{k}{n-k}\right)^{n-k} \approx (e^{k/(n-k)})^{n-k} = e^k$ and

$$\begin{aligned} |B_k| &\approx \left(\frac{ne}{k}\right)^k \\ &= 2^{k \log_2 \frac{ne}{k}} \\ &= 2^{k \log_2 \frac{n}{k} + k \log_2 e} \\ &= 2^{(\frac{k}{n} \log_2 \frac{n}{k})n + k \log_2 e} \\ &\approx 2^{(\frac{k}{n} \log_2 \frac{n}{k})n} \end{aligned} \tag{2.1}$$

for $1 \ll k \ll n$. Then, given that same assumption,

$$\begin{aligned} H\left(\frac{k}{n}\right) &= \frac{k}{n} \log_2 \frac{n}{k} + \left(1 - \frac{k}{n}\right) \log_2 \frac{1}{1 - \frac{k}{n}} \\ &\approx \frac{k}{n} \log_2 \frac{n}{k} \end{aligned}$$

because if $n \gg k$, $\frac{k}{n} \rightarrow 0$ and $1 \log_2 1 = 0$. Combining these approximations yields

$$|B_k| \approx 2^{H(\frac{k}{n})n} \quad \square$$

Let \mathbf{X} be a uniformly chosen point in B_k and $X_1, \dots, X_n \sim \text{Bernoulli}(\frac{k}{n})$.

This means that $H(\mathbf{X}) \approx H((X_1, \dots, X_n))$, which is remarkable because the latter could produce points in B_k or points with n ones or points with no ones.

This seems to imply that the majority of the mass of (X_1, \dots, X_n) lies within the boolean k -slice. Formally, we make the following claim about the concentration of measure:¹

¹cf. Dvoretzky–Milman theorem

Proposition 2.1.2

Fix any $\varepsilon > 0$. The probability

$$\Pr \left[(\mathbf{X}_1, \dots, \mathbf{X}_n) \notin \bigcup_{\ell=(1-\varepsilon)k}^{(1+\varepsilon)k} B_\ell \right] = \frac{1}{2^{n/\varepsilon^2}}$$

Informally, the probability of the randomly-drawn vector lying outside of the boolean k -slice is exponentially small.

We will prove a stronger claim:

Claim 2.1.3. Fix any $p \in (0, 1)$ and consider any $q > p$. Then,

$$\Pr[w((\mathbf{X}_i)) > q \cdot n] \leq 2^{-D(q\|p) \cdot n}$$

where $w((\mathbf{X}_i))$ is the number of ones. Likewise, consider any $q < p$. Then,

$$\Pr[w((\mathbf{X}_i)) < q \cdot n] \leq 2^{-D(q\|p) \cdot n}$$

Consider a toy example first. Let \mathbf{X} be the number of heads after n fair coin tosses.

Then, $\mathbb{E}[\mathbf{X}] = \frac{n}{2}$ and

$$\Pr[\mathbf{X} \geq 0.51n] = \frac{1}{2^n} \sum_{k \geq 0.51n} \binom{n}{k} \approx \frac{1}{2^n} \sum_{k \geq 0.51n} \left(\frac{ne}{k}\right)^k \rightarrow 0 \text{ very quickly}$$

by the same magic that we did in eq. (2.1) and because $\frac{1}{2^n}$ goes to 0 very quickly.

Now we can prove the claim.

Proof. Let $\theta_p(x)$ denote the probability of sampling a vector $x \in \{0, 1\}^n$ where each bit is IID Bernoulli(p). Then,

$$\begin{aligned} \frac{\theta_p(x)}{\theta_q(x)} &= \frac{p^k(1-p)^{n-k}}{q^k(1-q)^{n-k}} \\ &= \frac{(1-p)^n}{(1-q)^n} \left(\frac{p}{\frac{q}{1-q}} \right)^k \\ &\leq \frac{(1-p)^n}{(1-q)^n} \left(\frac{p}{\frac{q}{1-q}} \right)^{qn} \end{aligned}$$

for any $k \geq qn$ because (1) if $q \geq p$, then $\frac{q}{1-q} \geq \frac{p}{1-p}$ and the ugly fraction is greater than 1 and (2) increasing the exponent increases the quantity if the base is greater than 1.

Let $B_{\geq k} := \bigcup_{\ell \geq k} B_\ell$. Then, for all $x \in B_{\geq qn}$, we must show that

$$\theta_p(x) \leq \frac{(1-p)^n}{(1-q)^n} \left(\frac{p}{\frac{q}{1-q}} \right)^{qn} \cdot \theta_q(x) = 2^{-nD(q\|p) \cdot \theta_q(x)}$$

Consider the right-hand expression:

$$\begin{aligned} 2^{n \cdot D(q\|p)} &= 2^{n \cdot (q \log_2 \frac{1}{p} + (1-q) \log_2 \frac{1}{1-p} - q \log_2 \frac{1}{q} - (1-q) \log_2 \frac{1}{1-q})} \\ &= \left(\frac{1}{p^q} \cdot \frac{1}{(1-p)^{1-q}} \cdot q^q \cdot (1-q)^{1-q} \right)^n \end{aligned}$$

and the left-hand expression:

$$\begin{aligned} \frac{(1-p)^n}{(1-q)^n} \left(\frac{\frac{p}{1-p}}{\frac{q}{1-q}} \right)^{qn} &= \left(\frac{(1-p)^{1-q} p^q}{(1-q)^{1-q} q^q} \right)^n \\ &= \left(p^q \cdot (1-p)^{1-q} \cdot \frac{1}{q^q} \cdot \frac{1}{(1-q)^{1-q}} \right)^n \end{aligned}$$

which is clearly the reciprocal of the right-hand expression.

Now, we know that $\theta_p(x) = 2^{-nD(q\|p)} \theta_q(x)$, so

$$\begin{aligned} &\Pr_{\mathbf{X}_1, \dots, \mathbf{X}_n \sim \text{Bernoulli}(p)} [(\mathbf{X}_1, \dots, \mathbf{X}_n) \in B_{\geq qn}] \\ &= \sum_{x \in B_{\geq qn}} \theta_p(x) \\ &\leq 2^{-nD(q\|p)} \sum_{x \in B_{\geq qn}} \theta_q(x) \\ &\leq 2^{-nD(q\|p)} \end{aligned}$$

since the sum of the probabilities of x being any given entry in $B_{\geq qn}$ must be at most 1. \square

2.2 Rejection sampling

The KL divergence can give us a metric of how accurately we can sample one distribution using another distribution.

Example 2.2.1. Suppose $\mathbf{X} = \begin{cases} 0 & p = \frac{1}{2} \\ 1 & p = \frac{1}{2} \end{cases}$ and $\mathbf{Y} = \begin{cases} 0 & p = \frac{1}{4} \\ 1 & p = \frac{3}{4} \end{cases}$.

How can we sample \mathbf{Y} using \mathbf{X} ?

Solution (naive). Take IID \mathbf{X}_1 and \mathbf{X}_2 . Return 0 if $x_1 = x_2 = 0$ and 1 otherwise. \square

Solution (fancy). Take an infinite IID queue $\mathbf{X}_1, \mathbf{X}_2, \dots$

Starting at $i = 1$, if $\mathbf{X}_i = 0$, then output 0 with probability $\frac{1}{2}$, otherwise increment i until $\mathbf{X}_i = 1$. \square

↓ Lecture 6 adapted from Arthur ↓

Problem 2.2.2 (rejection sampling)

Given access to a distribution $Q = (Q(x))_{x \in \mathcal{X}}$, how efficiently can you simulate $P = (P(x))_{x \in \mathcal{X}}$?

Lecture 6
May 22

Example 2.2.3. Suppose $Q = (\frac{1}{3}, \frac{1}{3}, \frac{1}{3})$ and $P = (\frac{1}{2}, \frac{1}{2})$. We want to obtain the P distribution from Q .

Solution. Since Q and P are both uniform, we can just keep sampling from Q until we get something in P . That is, for $i = 1, \dots, \infty$:

1. Sample $X_i \sim Q$.
2. If $X_i \in \{1, 2\}$, accept and output $Y \leftarrow X_i$.
3. Otherwise, $i \leftarrow i + 1$.

This works because

$$\Pr[Y = 1] = \Pr[X_i = 1 \mid X_i = 1 \vee X_i = 2] = \frac{1/3}{2/3} = \frac{1}{2}$$

for the final round i , and similarly for $Y = 2$. □

Example 2.2.4. Consider a slightly more complex distribution $P = (\frac{1}{3}, \frac{2}{3})$ and $Q = (\frac{1}{2}, \frac{1}{2})$.

Solution. We will create a more complex rejection sampling protocol with some cheating.

Again, iterate and draw independent X_i :

- If $X_1 = 1$, accept with probability $\frac{2}{3}$. Otherwise, reject and continue to X_2 with probability $\frac{1}{3}$.
- If $X_1 = 2$, accept.
- For $i \geq 2$, accept if $X_i = 1$ and reject if $X_i = 2$.

Then, the probability of accepting $X_1 = 1$ is $\frac{1}{3}$, $X_1 = 2$ is $\frac{1}{2}$, and rejecting X_1 is $\frac{1}{6}$.

Since later rounds only output 1, we output 1 with probability $\frac{1}{3} + \frac{1}{6} = \frac{1}{2}$ and 2 with probability $\frac{1}{2}$. □

Definition 2.2.5 (rejection sampler)

A rejection sampler is a procedure that reads sequentially independent random samples $X_i \sim Q$ and in each round i either

- accepts the value of X_i and terminates with an index i^* , or
- rejects and continues.

The iteration we terminated on i^* is a random variable since it is a function of other random variables. It satisfies $X_{i^*} \sim P$, which is weird since for all fixed i , $X_i \sim Q$.

An interesting application is communication complexity. Suppose Alice has some hidden distribution P . Alice and Bob have access to a shared random IID sequence $X_i \sim Q$.

Alice can send an encoding of i^* to Bob who outputs $X_{i^*} \sim P$. This encoding i^* can be encoded using $\log i^*$ bits.

We will show that $\mathbb{E}[\log i^*] \leq D(P \parallel Q) + \mathcal{O}(1)$. You can also show that $D(P \parallel Q) \leq \mathbb{E}[\log i^*]$.

For each round i and symbol x , we need to know whether x was sampled before round i , i.e., the probability assigned to x in previous rounds.

For round $i \geq 1$, define:

- $\alpha_i(x)$ to denote the probability that the procedure accepts X_i and that $X_i = x$
- $p_i(x)$ to denote the probability that the procedure halts at round $i^* \leq i$ and $X_{i^*} = x$

We want to construct our procedure such that

- for all x , $P(x) = \sum_{i=1}^n \alpha_i(x)$
- for all x and i , $p_i(x) = \sum_{k=1}^i \alpha_k(x)$
- the probability that we halt on or before round i is $p_i^* := \sum_{x \in \mathcal{X}} p_i(x)$

↑ Lecture 6 adapted from Arthur ↑

Algorithm 1 REJECTIONSAMPLING(P, Q)

Require: $\forall x \in \mathcal{X}, Q(x) > 0 \iff D(P \parallel Q) < \infty$

```

1: for  $x \in \mathcal{X}$  do  $p_0(x) \leftarrow 0$ 
2:  $p_0^* \leftarrow 0$ 
3: for  $i = 1, \dots, \infty$  do
4:   sample  $X_i \sim Q$ 
5:   if  $P(X_i) - p_{i-1}(X_i) \leq (1 - p_{i-1}^*) \cdot Q(X_i)$  then
6:     with probability  $\beta_i(X_i) = \frac{P(X_i) - p_{i-1}(X_i)}{(1 - p_{i-1}^*)Q(X_i)}$  do
7:       ▷ so that the net probability of sampling  $X_i$  will be  $\alpha_i(X_i) = P(X_i) - p_{i-1}(X_i)$     ◁
8:       return  $X_i$ 
9:   else
10:    with probability  $\beta_i(X_i) = 1$  do
11:      ▷ so that the net probability of sampling  $X_i$  is  $\alpha_i(1 - p_{i-1}^*) \cdot Q(X_i)$     ◁
12:      return  $X_i$ 

```

Lecture 7
May 27

In this case, for all x and for all i :

- the probability of accepting x in round i is $\alpha_i(x) = \min\{P(x) - p_{i-1}(x), (1 - p_{i-1}^*)Q(x)\}$
- the probability of accepting x on or before round i is $p_i(x) = p_{i-1}(x) + \alpha_i(x)$
- the probability of terminating on or before round i is $p_i^* = p_{i-1}^* + \sum_{x \in \mathcal{X}} \alpha_i(x) = \sum_{x \in \mathcal{X}} p_i(x)$

Example 2.2.6. Let $P = (\frac{1}{2}, \frac{3}{8}, \frac{1}{8})$ and $Q = (\frac{1}{3}, \frac{1}{3}, \frac{1}{3})$. Do the procedure.

Solution. In round 1, sample $X_1 \sim Q$.

- If $X_1 = 1$, accept with probability 1.
- If $X_1 = 2$, accept with probability 1.
- If $X_1 = 3$, accept with probability $\frac{3}{8}$.

Then, $p_1(1) = \frac{1}{3}$, $p_1(2) = \frac{1}{3}$, $p_1(3) = \frac{1}{8}$, and $p_1^* = \frac{19}{24}$.

In round 2, sample $X_2 \sim Q$.

- If $X_2 = 1$, accept with probability 1. There is a $\frac{5}{72}$ chance of getting here, but deficit probability is $\frac{1}{6}$, so no need to reduce.
- If $X_2 = 2$, accept with probability $\frac{3}{5}$. There is a $\frac{5}{72}$ chance of getting here and deficit probability is $\frac{3}{8} - \frac{1}{3} = \frac{1}{24}$. For equality, use probability $\frac{3}{5} \cdot \frac{5}{72} = \frac{1}{24}$.
- If $X_3 = 3$, accept with probability 0. We already fulfilled $P(3) = p_1(3)$.

Then, $p_2(1) = \frac{29}{72}$, $p_2(2) = \frac{3}{8}$, $p_3(2) = \frac{1}{8}$, and $p_2^* = \frac{19}{24} + \frac{5}{24} \cdot (\frac{1}{3} + \frac{3/5}{5}) = \frac{65}{72}$.

In round 3, sample $X_3 \sim Q$.

- If $X_3 = 1$, accept with probability 1.
- If $X_3 = 2$ or 3, accept with probability 0.

Keep repeating until we accept a 1. □

Proposition 2.2.7

$(p_i(x))_{x \in \mathcal{X}}$ converges to $P(x)$ as $i \rightarrow \infty$. In fact, the residual decays exponentially fast

$$P(x) - p_i(x) \leq P(x) \cdot (1 - Q(x))^i.$$

Proof. Begin with the claim that the probability of reaching round i is at least the residual at i for any x :

$$1 - p_{i-1}^* \geq P(x) - p_{i-1}(x) \quad \forall x$$

Intuitively, either you returned prior to round i (i.e., p_{i-1}^*) or you did not (i.e., the residual).

$$\begin{aligned} 1 - p_{i-1}^* &= \sum_{x \in \mathcal{X}} P(x) - \sum_{x \in \mathcal{X}} p_{i-1}(x) \\ &= \sum_{x \in \mathcal{X}} (P(x) - p_{i-1}(x)) \end{aligned} \tag{2.2}$$

Also, claim that

$$\alpha_i \geq (P(x) - p_{i-1}(x)) \cdot Q(x) \tag{2.3}$$

If $\alpha_i = P(x) - p_{i-1}(x)$, then clearly $\alpha_i \geq \alpha_i Q(x)$. Otherwise, if $\alpha_i = (1 - p_{i-1}^*)Q(x)$, then eq. (2.2) applies.

Proceed by induction.

Base case: exercise.

Inductive step: suppose that $P(x) - p_i(x) \leq P(x) \cdot (1 - Q(x))^i$. Then,

$$\begin{aligned} P(x) - p_{i+1}(x) &= P(x) - p_i(x) - \alpha_{i+1}(x) \\ &\leq (P(x) - p_{i-1}(x))(1 - Q(x)) && \text{(by eq. (2.3))} \\ &\leq (P(x) \cdot (1 - Q(x))^i)(1 - Q(x)) && \text{(by supposition)} \\ &\leq P(x) \cdot (1 - Q(x))^{i+1} \end{aligned} \quad \square$$

Now, we will prove that this is related to relative entropy.

Proposition 2.2.8

Let i^* be the iteration at which the procedure returns. Then, $\mathbb{E}[\log_2 i^*] \leq D(P \parallel Q) + 2 \log_2 e$.

Proof. First, claim that for all $x \in \mathcal{X}$ and any $i \geq 2$ such that $\alpha_i(x) > 0$,

$$i \leq \frac{P(x)}{(1 - p_{i-1}^*) \cdot Q(x)} + 1 \quad (2.4)$$

That is, if we reach a particular round i , the probability mass left must be sufficiently large.

We know that $P(x) \geq p_{i-1}(x)$ since we increase to $P(x)$. Then,

$$\begin{aligned} P(x) &\geq p_{i-1}(x) \\ &= \alpha_1(x) + \dots + \alpha_{i-1}(x) \\ &\geq (1 - p_1^*) \cdot Q(x) + \dots + (1 - p_{i-1}^*) \cdot Q(x) \\ &\geq (1 - p_{i-1}^*) \cdot Q(x) + \dots + (1 - p_{i-1}^*) \cdot Q(x) \\ &= (i-1)(1 - p_{i-1}^*) \cdot Q(x) \\ i &\leq \frac{P(x)}{(1 - p_{i-1}^*) \cdot Q(x)} + 1 \end{aligned}$$

as long as $\alpha_{j-1} < \alpha_j$ for all j .

Do a gigantic algebra bash:

Lecture 8
May 29

$$\begin{aligned} \mathbb{E}[\log_2 i^*] &= \sum_{i=1}^{\infty} (p_i^* - p_{i-1}^*) \cdot \log_2 i \\ &= \sum_{i=1}^{\infty} \sum_{x \in \mathcal{X}} \alpha_i(x) \cdot \log_2 i \\ &\leq \sum_{i=1}^{\infty} \sum_{x \in \mathcal{X}} \alpha_i(x) \cdot \log_2 \left[\frac{P(x)}{(1 - p_{i-1}^*) Q(x)} + 1 \right] \quad (\text{by eq. (2.4)}) \\ &\leq \sum_{i=1}^{\infty} \sum_{x \in \mathcal{X}} \alpha_i(x) \cdot \log_2 \left[\frac{1}{(1 - p_{i-1}^*)} \left(\frac{P(x)}{Q(x)} + 1 \right) \right] \\ &= \underbrace{\sum_{i=1}^{\infty} \sum_{x \in \mathcal{X}} \alpha_i(x) \log_2 \frac{1}{(1 - p_{i-1}^*)}}_A + \underbrace{\sum_{i=1}^{\infty} \sum_{x \in \mathcal{X}} \alpha_i(x) \log_2 \left(\frac{P(x)}{Q(x)} + 1 \right)}_B \end{aligned}$$

Consider the first term A :

$$\begin{aligned} A &= \sum_{i=1}^{\infty} \sum_{x \in \mathcal{X}} \alpha_i(x) \log_2 \frac{1}{(1 - p_{i-1}^*)} \\ &= \sum_{i=1}^{\infty} (p_i^* - p_{i-1}^*) \log_2 \frac{1}{(1 - p_{i-1}^*)} \end{aligned}$$

Notice that this is a left-handed Riemann sum of $\log_2 \frac{1}{1-x}$:

$$\begin{aligned} A &\leq \int_0^1 \log_2 \frac{1}{1-x} dx \\ &= \log_2 e \end{aligned}$$

Now, consider the second term B :

$$\begin{aligned}
B &= \sum_{i=1}^{\infty} \sum_{x \in \mathcal{X}} \alpha_i(x) \log_2 \left(\frac{P(x)}{Q(x)} + 1 \right) \\
&= \sum_{x \in \mathcal{X}} \sum_{i=1}^{\infty} \alpha_i(x) \log_2 \left(\frac{P(x)}{Q(x)} + 1 \right) && \text{(Fubini?)} \\
&= \sum_{x \in \mathcal{X}} P(x) \log_2 \left(\frac{P(x)}{Q(x)} + 1 \right) && (P(x) = \sum_i \alpha_i(x)) \\
&= \sum_{x \in \mathcal{X}} P(x) \log_2 \left(\frac{P(x)}{Q(x)} \cdot \left(1 + \frac{Q(x)}{P(x)} \right) \right) \\
&= \sum_{x \in \mathcal{X}} P(x) \log_2 \left(\frac{P(x)}{Q(x)} \right) + \sum_{x \in \mathcal{X}} P(x) \log_2 \left(1 + \frac{Q(x)}{P(x)} \right) \\
&= D(P \parallel Q) + \sum_{x \in \mathcal{X}} P(x) \log_2 \left(1 + \frac{Q(x)}{P(x)} \right) \\
&\leq D(P \parallel Q) + \sum_{x \in \mathcal{X}} P(x) \log_2 (e^{Q(x)/P(x)}) && (1 + x \leq e^x \text{ for all } x \geq 0) \\
&= D(P \parallel Q) + \sum_{x \in \mathcal{X}} P(x) \frac{Q(x)}{P(x)} \log_2 e \\
&= D(P \parallel Q) + \log_2 e \sum_{x \in \mathcal{X}} Q(x) \\
&= D(P \parallel Q) + \log_2 e
\end{aligned}$$

Therefore,

$$\mathbb{E}[\log_2 i^*] \leq A + B \leq D(P \parallel Q) + 2 \log_2 e$$

completing the proof. \square

Intuition: for any $x \in \mathcal{X}$, if $\alpha_i(x) \leq Q(x) \lll P(x)$, then you need an expected amount of $\frac{P(x)}{Q(x)}$ steps to succeed, because you just won't roll x that often.

Also, if $\alpha_{i+1}(x) > 0$ (any round prior to termination), $(1 - p_{i-1}^*(x))Q(x) \leq \alpha_i(x)$.

Proposition 2.2.9

For any rejection sampler, let i^* be the index where it returns. Then,

$$\mathbb{E}[\ell(i^*)] \geq D(P \parallel Q)$$

Proof. For convenience, redefine $\alpha_i(x) := \Pr[i^* = i \wedge \mathbf{X}_i = x]$.

First, observe that for any $x \in \mathcal{X}$, a rejection sampler must have

$$\alpha_i(x) \leq Q(x)$$

because we only have a $Q(x)$ chance of rolling x to accept it in round i .

Now, fix $x \in \mathcal{X}$. Consider the random variable $i^*|_{\mathbf{X}_{i^*}=x}$. Then, by [Kraft's inequality](#),

$$\begin{aligned}
 \mathbb{E}[\ell(i^*) \mid \mathbf{X}_{i^*} = x] &\geq H(i^* \mid \mathbf{X}_{i^*} = x) \\
 &= \sum_{i=1}^{\infty} \Pr[i^* = i \mid \mathbf{X}_{i^*} = x] \log_2 \frac{1}{\Pr[i^* = i \mid \mathbf{X}_{i^*} = x]} \\
 &= \sum_{i=1}^{\infty} \frac{\alpha_i(x)}{P(x)} \log_2 \frac{P(x)}{\alpha_i(x)} \\
 &\geq \sum_{i=1}^{\infty} \frac{\alpha_i(x)}{P(x)} \log_2 \frac{P(x)}{Q(x)} \\
 &= \log_2 \frac{P(x)}{Q(x)} \cdot \sum_{i=1}^{\infty} \frac{\alpha_i(x)}{P(x)} \\
 &= \log_2 \frac{P(x)}{Q(x)}
 \end{aligned}$$

because $\sum_{i=1}^{\infty} \alpha_i(x) = P(x)$. Apply the law of total probability:

$$\begin{aligned}
 \mathbb{E}[\ell(i^*)] &= \sum_{x \in \mathcal{X}} \Pr[\mathbf{X}_{i^*} = x] \mathbb{E}[\ell(i^*) \mid \mathbf{X}_{i^*} = x] \\
 &= \sum_{x \in \mathcal{X}} P(x) \mathbb{E}[\ell(i^*) \mid \mathbf{X}_{i^*} = x] \\
 &\geq \sum_{x \in \mathcal{X}} P(x) \log_2 \frac{P(x)}{Q(x)} \\
 &= D(P \parallel Q)
 \end{aligned}$$

as desired. □

Chapter 3

Mutual information

3.1 Definition and chain rules

Notation. Given two jointly distributed random variables (X, Y) over sample space $\mathcal{X} \times \mathcal{Y}$, write p_{xy} for $\Pr[X = x, Y = y]$.

*Lecture 9
June 3*

Definition 3.1.1

Given two jointly distributed random variables (X, Y) over sample space $\mathcal{X} \times \mathcal{Y}$, define the mutual information $I(X : Y)$ by

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

where the conditional entropy $H(X | Y)$ is

$$\sum_{y \in \mathcal{Y}} p_y \cdot H((X|_{Y=y}))$$

This is entirely analogous to saying that $|A \cap B| = |A| + |B| - |A \cup B| = |A| - |A \setminus B|$.

Theorem 3.1.2 (chain rule for entropy)

Given two jointly distributed random variables (X, Y) over a discrete sample space $\mathcal{X} \times \mathcal{Y}$,

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

Proof. Do a bunch of algebra:

$$\begin{aligned}
 H(X) + H(Y | X) &= \sum_{x \in \mathcal{X}} p_x \log \frac{1}{p_x} + \sum_{x \in \mathcal{X}} p_x \sum_{y \in \mathcal{Y}} \Pr[Y = y | X = x] \log \frac{1}{\Pr[Y = y | X = x]} \\
 &= \sum_{x \in \mathcal{X}} p_x \log \frac{1}{p_x} + \sum_{x \in \mathcal{X}} p_x \sum_{y \in \mathcal{Y}} \frac{p_{xy}}{p_x} \log \frac{p_x}{p_{xy}} \\
 &= \sum_{\substack{x \in \mathcal{X} \\ y \in \mathcal{Y}}} p_{xy} \log \frac{1}{p_x} + \sum_{\substack{x \in \mathcal{X} \\ y \in \mathcal{Y}}} p_{xy} \log \frac{p_x}{p_{xy}} \\
 &= \sum_{\substack{x \in \mathcal{X} \\ y \in \mathcal{Y}}} p_{xy} \left(\log \frac{1}{p_x} + \log \frac{p_x}{p_{xy}} \right) \\
 &= \sum_{\substack{x \in \mathcal{X} \\ y \in \mathcal{Y}}} p_{xy} \log \frac{1}{p_{xy}} \\
 &= H((X, Y))
 \end{aligned}$$

□

Corollary 3.1.3. For two independent variables, since $(Y | X) = Y$, we have $H((X, Y)) = H(X) + H(Y)$ as expected.

Corollary 3.1.4. $H((X_1, X_2, X_3)) = H(X_1) + H(X_2 | X_1) + H(X_3 | (X_1, X_2))$

Proof. Consider $(X_1, X_2, X_3) = ((X_1, X_2), X_3)$. Then, by the [chain rule for entropy](#),

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

and then by another application,

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

as desired.

□

Theorem 3.1.5 (general chain rule for entropy)

For k random variables X_1, \dots, X_k ,

$$H((X_1, \dots, X_k)) = \sum_{i=1}^k H(X_i | (X_1, \dots, X_{i-1}))$$

Proof. By induction on the [chain rule for entropy](#).

□

Notation. Although relative entropy is defined only on *distributions*, write $D(X \parallel Y)$ to be $D(f_X \parallel f_Y)$.

Theorem 3.1.6 (chain rule for relative entropy)

Let p and $q : \mathcal{X} \times \mathcal{Y} \rightarrow [0, 1]$ be distributions. Let $p(x) := \sum_{y \in \mathcal{Y}} p(x, y)$ denote marginals of p and $p(y|x) := \frac{p(x, y)}{p(x)}$ denote conditionals of p . Then,

$$\begin{aligned} D(p(x, y) \parallel q(x, y)) &= D(p(x) \parallel q(x)) + D(p(y|x) \parallel q(y|x)) \\ &= D(p(x) \parallel q(x)) + \sum_{x \in \mathcal{X}} p(x) \cdot D((p(y|x))_{y \in \mathcal{Y}} \parallel (q(y|x))_{y \in \mathcal{Y}}) \end{aligned}$$

where $D(p(y|x) \parallel q(y|x))$ is the conditional relative entropy.

Equivalently, let (X_1, Y_1) and (X_2, Y_2) be two joint random variables. Then,

$$D((X_1, Y_1) \parallel (X_2, Y_2)) = D(X_1 \parallel X_2) + \sum_{x \in \mathcal{X}} \Pr[X_1 = x] \cdot D(Y_1|_{X_1=x} \parallel Y_2|_{X_2=x})$$

Proof (for distributions). Do algebra:

$$\begin{aligned} &D(p(x) \parallel q(x)) + D(p(y|x) \parallel q(y|x)) \\ &= \sum_{x \in \mathcal{X}} p_x \log \frac{p_x}{q_x} + \sum_{x \in \mathcal{X}} p_x \sum_{y \in \mathcal{Y}} p(y|x) \log \frac{p(y|x)}{q(y|x)} \\ &= \sum_{x \in \mathcal{X}} p_x \log \frac{p_x}{q_x} + \sum_{x \in \mathcal{X}} p_x \sum_{y \in \mathcal{Y}} \frac{p_{xy}}{p_x} \log \frac{p_{xy} q_x}{q_{xy} p_x} \\ &= \sum_{\substack{x \in \mathcal{X} \\ y \in \mathcal{Y}}} p_{xy} \log \frac{p_x}{q_x} + \sum_{\substack{x \in \mathcal{X} \\ y \in \mathcal{Y}}} p_{xy} \log \frac{p_{xy} q_x}{q_{xy} p_x} \\ &= \sum_{\substack{x \in \mathcal{X} \\ y \in \mathcal{Y}}} p_{xy} \left(\log \frac{p_x}{q_x} + \log \frac{p_{xy} q_x}{q_{xy} p_x} \right) \\ &= \sum_{\substack{x \in \mathcal{X} \\ y \in \mathcal{Y}}} p_{xy} \log \frac{p_{xy}}{q_{xy}} \\ &= D(p(x, y) \parallel q(x, y)) \end{aligned}$$

as in the proof of [chain rule for entropy](#). □

Fact 3.1.7.

$$I[X : Y] = \mathbb{E}_{x \leftarrow X} [D(Y|_{X=x} \parallel Y)] = \sum_{x \in \mathcal{X}} p_x D(Y|_{X=x} \parallel Y)$$

Proof. First, claim that

$$I[X : Y] = D((X, Y) \parallel \tilde{X} \otimes \tilde{Y}) \quad (3.1)$$

where $\tilde{X} \otimes \tilde{Y}$ denotes a random variable consisting of \tilde{X} (resp. \tilde{Y}) independently sampled according

to the distribution of \mathbf{X} (resp. \mathbf{Y}) so that $\Pr[\tilde{\mathbf{X}} = x, \tilde{\mathbf{Y}} = y] = p_x p_y$. Expand the left-hand side:

$$\begin{aligned}
 I[\mathbf{X} : \mathbf{Y}] &= \sum_{x \in \mathcal{X}} p_x \log \frac{1}{p_x} + \sum_{y \in \mathcal{Y}} p_y \log \frac{1}{p_y} - \sum_{\substack{x \in \mathcal{X} \\ y \in \mathcal{Y}}} p_{xy} \log \frac{1}{p_{xy}} \\
 &= \sum_{x \in \mathcal{X}} \sum_{y \in \mathcal{Y}} p_{xy} \log \frac{1}{p_x} + \sum_{y \in \mathcal{Y}} \sum_{x \in \mathcal{X}} p_{xy} \log \frac{1}{p_y} - \sum_{\substack{x \in \mathcal{X} \\ y \in \mathcal{Y}}} p_{xy} \log \frac{1}{p_{xy}} \\
 &= \sum_{x \in \mathcal{X}} \sum_{y \in \mathcal{Y}} p_{xy} \left(\log \frac{1}{p_x} + \log \frac{1}{p_y} - \log \frac{1}{p_{xy}} \right) \\
 &= \sum_{\substack{x \in \mathcal{X} \\ y \in \mathcal{Y}}} p_{xy} \log \frac{p_{xy}}{p_x p_y} \\
 &= D((\mathbf{X}, \mathbf{Y}) \parallel \tilde{\mathbf{X}} \otimes \tilde{\mathbf{Y}})
 \end{aligned}$$

Now, apply the [chain rule for relative entropy](#):

$$\begin{aligned}
 D((\mathbf{X}, \mathbf{Y}) \parallel \tilde{\mathbf{X}} \otimes \tilde{\mathbf{Y}}) &= D(\mathbf{X} \parallel \tilde{\mathbf{X}}) + D((\mathbf{X}, \mathbf{Y}) \mid (\mathbf{X}, \tilde{\mathbf{X}}) \parallel (\tilde{\mathbf{X}} \otimes \tilde{\mathbf{Y}}) \mid (\mathbf{X}, \tilde{\mathbf{X}})) \\
 &= 0 + \sum_x p_x D(\mathbf{Y} \mid \mathbf{X} = x \parallel \mathbf{Y}) \\
 &= \mathbb{E}_{x \leftarrow \mathbf{X}} D(\mathbf{Y} \mid \mathbf{X} = x \parallel \mathbf{Y})
 \end{aligned}$$

□

Lecture 11
June 10

Theorem 3.1.8 (chain rule for mutual information)

Let $\mathbf{X}_1, \mathbf{X}_2$, and \mathbf{Y} be random variables. Then,

$$I((\mathbf{X}_1, \mathbf{X}_2) : \mathbf{Y}) = I(\mathbf{X}_1 : \mathbf{Y}) + I(\mathbf{X}_2 : (\mathbf{Y} \mid \mathbf{X}_1))$$

and in general

$$I((\mathbf{X}_1, \dots, \mathbf{X}_n) : \mathbf{Y}) = \sum_{i=1}^n I(\mathbf{X}_i : (\mathbf{Y} \mid (\mathbf{X}_1, \dots, \mathbf{X}_{i-1})))$$

3.2 Markov chains, data processing, and sufficient statistics

Definition 3.2.1

The random variables \mathbf{X} , \mathbf{Y} , and \mathbf{Z} form a Markov chain if the conditional distribution of \mathbf{Z} depends only on \mathbf{Y} and is conditionally independent of \mathbf{X} . Equivalently,

$$\Pr[\mathbf{X} = x, \mathbf{Y} = y, \mathbf{Z} = z] = \Pr[\mathbf{X} = x] \cdot \Pr[\mathbf{Y} = y \mid \mathbf{X} = x] \cdot \Pr[\mathbf{Z} = z \mid \mathbf{Y} = y]$$

Then, we write $\mathbf{X} \rightarrow \mathbf{Y} \rightarrow \mathbf{Z}$.

Example 3.2.2 (Legend of the Drunken Master). In $\Omega = \mathbb{R}^2$, Jackie Chan is drunk and takes steps in random directions. He starts at $J_0 = (0,0)$. Then, $J_1 = J_0 + d_1$ where d_1 is an independent random unit vector in \mathbb{R}^2 , and $J_2 = J_1 + d_2$ and so on.

First, J_3 and J_1 are not independent. But if we fix $J_2 = j_2 \in \mathbb{R}^2$, then $J_1 \mid J_2 = j_2$ and $J_3 \mid J_2 = j_2$ are independent. In fact, they are uniformly distributed random points on the circle of radius 1 centred at j_2 .

Proposition 3.2.3

Let X , Y , and Z be random variables. TFAE:

1. $X \rightarrow Y \rightarrow Z$
2. X and Z are conditionally independent given Y . That is,

$$\Pr[X = x, Z = z \mid Y = y] = \Pr[X = x \mid Y = y] \cdot \Pr[Z = z \mid Y = y]$$

3. Z is distributed according to $f(Y, R)$ for some R independent of X and Y .

Exercise 3.2.4. Prove the definitions are equivalent.

Theorem 3.2.5 (data-processing inequality)

If $X \rightarrow Y \rightarrow Z$, then $I(X : Z) \leq I(X : Y)$.

Equality happens if and only if $X \rightarrow Z \rightarrow Y$.

Proof. By the chain rule,

$$I(X : (Y, Z)) = I(X : Y) + I(X : Z \mid Y) \overset{0}{=} I(X : Z) + I(X : Y \mid Z)$$

so that

$$I(X : Y) = I(X : Z) + I(X : Y \mid Z)$$

One may show that the mutual information is always non-negative, so we have $I(X : Y) \geq I(X : Z)$ as desired. \square

3.3 Communication complexity

Problem 3.3.1

Suppose there is a joint distribution (X, Y) that Alice and Bob wish to jointly compute. Alice and Bob have access to a shared random string $R = (R_i)$. Alice is given $x \in \mathcal{X}$ and wants to send Bob a prefix-free message of minimum length so that Bob can compute a sample from $Y \mid X = x$.

Lecture 10
June 5

Definition 3.3.2

A protocol Π is a pair of functions (M, y) where $M : \mathcal{X} \times \Omega_R \rightarrow \{0, 1\}^*$ is the message Alice sends to Bob and $y : \{0, 1\}^* \times \Omega_R \rightarrow \mathcal{Y}$ is Bob's output.

The performance of Π is $\mathbb{E}_{X, R} |M(X, R)|$

Suppose X and Y are independent. Then, Bob needs no information so we can use the trivial protocol $M(X, R) = \emptyset$ with performance 0.

Otherwise, we can use a strategy of prefix-free encoding x so that $\mathbb{E} |M(X, R)| \approx H(X)$.

Theorem 3.3.3

There exists a protocol $\Pi = (M, y)$ such that expected message length

$$\mathbb{E} |M(X, R)| \leq I(X : Y) + \mathcal{O}(\log I(X : Y))$$

For all other protocols $\Pi' = (M', y')$,

$$\mathbb{E} |M'(X, R)| \geq I(X : Y)$$

Proof. Let X be a random point on the hypercube $\{\pm 1\}^n$. Let Y be a random point on $\{\pm 1\}^n$ that is ε -correlated with X . That is, $Y_i = X_i$ with probability ε and is uniformly random otherwise.

Observe that, individually, X and Y have the same distribution. In particular, in the ε case, then $Y_i = X_i$ is $\text{Uniform}\{\pm 1\}$. In the $1 - \varepsilon$ case, $Y_i \sim \text{Uniform}\{\pm 1\}$ by definition.

We can calculate $H(X) = H(Y) = n$.

Also, $H(Y | X) = \sum_x p_x H(Y | X = x) \approx (1 - \varepsilon)n$. One can show that $Y | X = x$ is approximately uniformly distributed over the vectors of length n that agree on εn coordinates with x . This sample space has size $2^{(1-\varepsilon)n}$.

Therefore, $I(X : Y) = H(Y) - H(Y | X) \approx \varepsilon n$.

By prop. 2.2.8, there exists a rejection sampler such that $\mathbb{E}[\ell(i^*)] \leq D(P \parallel Q) + \mathcal{O}(\log D(P \parallel Q))$.

Recall from STAT 230 that we can transform R into any distribution with the change of variable bullshit. In particular, transform R_i to IID $Y_i \sim Y$ and the biased coins.

Alice will run $\text{REJECTION_SAMPLER}(Y|_{X=x}, Y)$ to find a random index i^* such that Y_{i^*} has distribution $Y|_{X=x}$.

Alice sends a prefix-free encoding of i^* . Bob outputs Y_{i^*} . The performance is:

$$\begin{aligned} \mathbb{E}_{X, R} |M(X, R)| &= \sum_{x \in \mathcal{X}} p_x \mathbb{E}_{i^*, Y_1, Y_2, \dots} [\ell(i^*)] \\ &\leq \sum_{x \in \mathcal{X}} p_x (D(Y|_{X=x} \parallel Y) + \mathcal{O}(\log D(Y|_{X=x} \parallel Y))) \\ &= I(X : Y) + \sum_{x \in \mathcal{X}} p_x \mathcal{O}(\log D(Y|_{X=x} \parallel Y)) \\ &\leq I(X : Y) + \mathcal{O}(\log I(X : Y)) \end{aligned}$$

where the last step is by Jensen's inequality.

Now, let Π be any protocol. We will apply the [data-processing inequality](#).

Lecture 11

June 10

cont.

Notice that $X \rightarrow (M(X, R), R) \rightarrow Y$ if and only if Π is a valid protocol. If we sample $x \sim X$ and Alice sends $M(x, R)$, then Bob outputs something distributed according to $Y \mid X = x$, i.e., just Y since x was arbitrary. Then,

$$\begin{aligned}
 I(X : Y) &\leq I(X : (M(X, R), R)) && \text{(data processing inequality)} \\
 &= I(X : R) + \sum_{r \in \Omega_R} p_r I(X|_{R=r} : M(X, R)|_{R=r}) && \text{(chain rule)} \\
 &= 0 + I(X : M(X, R) \mid R) && \text{(independence)} \\
 &\leq H(M(X, R) \mid R) && (I(A : B) \leq \min\{H(A), H(B)\}) \\
 &\leq H(M(X, R)) && (H(A \mid B) \leq H(A)) \\
 &\leq \mathbb{E} |M(X, R)| && \text{(Kraft inequality)}
 \end{aligned}$$

completing the proof. \square

3.4 Parameter estimation

List of Named Results

1.1.4	Theorem (Jensen's inequality)	3
1.3.5	Theorem (Kraft's inequality)	5
2.0.7	Theorem (Sterling's approximation)	11
3.1.2	Theorem (chain rule for entropy)	21
3.1.5	Theorem (general chain rule for entropy)	22
3.1.6	Theorem (chain rule for relative entropy)	23
3.1.8	Theorem (chain rule for mutual information)	24
3.2.5	Theorem (data-processing inequality)	25

Index of Defined Terms

boolean k -slice, [11](#)

code

prefix-free, [5](#), [9](#)

Shannon–Fano, [6](#)

uniquely decodable, [5](#),
[9](#)

concentration of measure,
[12](#)

entropy, [2](#)

conditional, [21](#)

relative, [8](#)

conditional, [23](#)

Hamming k -slice, [11](#)

KL divergence, [8](#)

Markov chain, [24](#)

mutual information, [21](#)

rejection sampler, [15](#)