

Notes for EECS 550: Information Theory

Yiwei Fu, Instructor: David Neuhoff

FA 2022

Contents

1	Source Codes	1
1.1	Lossless Coding	1
1.1.1	Fixed-length to Fixed-length Block Codes (FFB)	2
1.2	Shannon-McMillian Theorem	7
1.3	Fixed Length to Variable Length (FVB) Lossless Source codes	9
1.4	Huffman's Code Design	12
1.5	Buffering	13
1.6	Universal Source Coding	14
2	Entropy	15
2.1	Entropy	15
2.2	Basic Properties of Entropy	18
2.3	Conditional Entropy	20
2.4	Convexity	21
3	Information	22
3.1	Information	22
3.2	Conditional Information	24
3.3	Cryptography From Information Perspective	24
3.4	Continuous Random Variables and Information	25
3.5	Differential Entropy	26
3.6	Properties of Differential Entropy	27
4	Estimation, Decision	28
4.1	Estimation Theory	28
4.2	Decision Theory	28

Office hours:

Chapter 1

Source Codes

1.1 Lossless Coding

Lossless coding is a type of data compression.

GOAL to encode data into bits so that

1. bits can be decoded perfectly or with very high accuracy back into original data;
2. we use as few bits as possible.

We need to model for data, a measure of decoding accuracy, a measure of compactness.

MODEL FOR DATA

Definition 1.1.1. A *source* is a sequence of i.i.d (discrete) random variables U_1, U_2, \dots

We would like to assume a known alphabet $A = \{a_1, a_2, \dots, a_Q\}$ and known probability distribution either through probability mass functions $p_U(u) = \Pr[U = u]$.

Definition 1.1.2. Source coding

PERFORMANCE MEASURES A measure of compactness (efficiency)

Definition 1.1.3. *Encoding rate*, also called *rate*, is the average number of encoded bits per data symbol.

There are two versions of average rate:

1. Empirical average rate

$$\langle r \rangle := \lim_{N \rightarrow \infty} \frac{1}{N} \sum_{k=1}^N L_k(U_1, \dots, U_k),$$

2. Statistical average rate

$$\bar{r} := \lim_{N \rightarrow \infty} \frac{1}{N} \sum_{k=1}^N \mathbb{E}[L_k(U_1, \dots, U_k)]$$

where L_K is the number of bits out of the encoder after U_k and before U_{k+1} .

Definition 1.1.4. The per-letter frequency of error is defined as

$$\langle F_{LE} \rangle := \lim_{N \rightarrow \infty} \frac{1}{N} \sum_{k=1}^N I(\hat{U}_k = U_k)$$

and per-letter error probability is defined as

$$p_{LE} := \lim_{N \rightarrow \infty} \frac{1}{N} \sum_{k=1}^N \mathbb{E}[I(\hat{U}_k = U_k)] = \lim_{N \rightarrow \infty} \frac{1}{N} \sum_{k=1}^N \Pr(\hat{U}_k = U_k)$$

1.1.1 Fixed-length to Fixed-length Block Codes (FFB)

characteristics

Definition 1.1.5. A code is *perfectly lossless* (PL) if the $\beta(\alpha(\underline{u})) = \underline{u}$ for all $\underline{u} \in A_U^k$ (the set of all sequences u_1, \dots, u_k).

In order to be perfectly loss, α must be one-to-one. Encode must assign a distinct code-word (L bits) to each data sequences. rate $= L/K$. We seek $R_{PL}^*(k)$ the smallest rate of any PL code.

Number of sequences of size $k = Q^k$, and number binary sequence of size $L = 2^L$. We need $2^L \gg Q^k$.

$$\bar{r} = \frac{L}{k} \geq \frac{k \log_2 Q}{k} = \log_2 Q$$

Choose $\lceil k \log_2 Q \rceil$, then we have

$$R_{PL}^*(k) = \frac{\lceil k \log_2 Q \rceil}{k} \leq \frac{k \log_2 Q + 1}{k} = \log_2 Q + \frac{1}{k}.$$

$$\log_2 Q \leq R_{PL}^*(k) \leq \log_2 Q + \frac{1}{k}$$

Let R_{PL}^* be the least rate of any PL FFB code with any k . $R_{PL}^*(k) \rightarrow \log_2 Q$ as $k \rightarrow \infty$.

$$R_{PL}^* = \inf_k R_{PL}^*(k)$$

Now we want rate less and $\log_2 Q$ almost lossless codes.

$$R_{AL}^* = \inf\{r, \text{there is an FFB code with } \bar{r} \leq r \text{ and arbitrarily small } P_{LE}\}$$

$$= \inf\{r, \text{there is an FFB code with } \bar{r} \leq r \text{ and } P_{LE} < \delta \text{ for all } \delta > 0\}$$

Instead of per-letter probability P_{LE} , we focus on block error probability $P_{BE} = \Pr(\hat{U} \neq U)$

Lemma 1.1.1. $P_{BE} \geq P_{LE} \geq \frac{P_{BE}}{k}$

Proof. See homework. ■

To analyze, we focus on the set of correctly encoded sequences. $G = \{\underline{u} : \beta(\alpha(\underline{u})) = \underline{u}\}$

Then we have

$$P_{BE} = 1 - \Pr[U \in G], |G| \leq 2^k, L \geq \lceil \log_2 |G| \rceil.$$

QUESTION How large is the smallest set of sequences with length k form A_U with probability ≈ 1 ?

We need to use weak law of large numbers (WLLN).

Theorem 1.1.1. Suppose $A_x = \{1, 2, \dots, Q\}$ with probability p_1, \dots, p_Q . Given $\underline{u} = (u_1, \dots, u_k) \in A_U^k$.

$$n_q(\underline{u}) := \# \text{times } a_q \text{ occurs in } \underline{u}, \quad f_q(\underline{u}) = \frac{n_q(\underline{u})}{k} = \text{frequency}$$

Fix any $\varepsilon > 0$,

$$\Pr[f_q(\underline{u}) \doteq p_q \pm \varepsilon] \rightarrow 1 \text{ as } k \rightarrow \infty.$$

Moreover,

$$\Pr[f_q(\underline{u}) \doteq p_q \pm \varepsilon, q = 1, \dots, Q] \rightarrow 1 \text{ as } k \rightarrow \infty.$$

NOTATION $a \doteq b \pm \varepsilon \iff |a - b| \leq \varepsilon$

Consider subset of A_U^k that corresponds to this event x .

$$T_k = \{\underline{u} : f_q(\underline{u}) \doteq p_q \pm \varepsilon, q = 1, \dots, Q\}.$$

$$\Pr[\underline{U} = \underline{u}] = p(u_1)p(u_2) \dots p(u_k).$$

By WLLN, $\Pr(T_k) \rightarrow 1$ as $k \rightarrow \infty$.

KEY FACT all sequences in T_k have approximately the same probability.

For $\underline{u} \in T_k$,

$$\begin{aligned} p(\underline{u}) &= p(u_1)p(u_2) \dots p(u_k) \\ &= p_1^{n_1(u)} p_2^{n_2(u)} \dots p_k^{n_k(u)} \\ &= p_1^{kf_1(u)} p_2^{kf_2(u)} \dots p_k^{kf_k(u)} \\ &\approx \tilde{p}^k \text{ where } \tilde{p} = p_1^{p_1} p_2^{p_2} \dots p_Q^{p_Q}. \end{aligned}$$

So we have $|T_k| \approx \frac{1}{\tilde{p}^k}$.

Then we have

$$\bar{r} = \frac{\log_2 |T_k|}{k} = -\frac{k \log_2 \tilde{p}}{k} = -\log_2 \tilde{p}.$$

Is that rate good? Can we do better? Can we have a set S with probability ≈ 1 and significantly smaller?

Since $\Pr(\underline{U} \in A_U^k \setminus T_k) \approx 0 \implies \Pr(\underline{U} \in S) \approx \Pr(\underline{U} \in S \cap T_k) \approx \frac{|S|}{|T_k|}$. So when k is large, T_k is the smallest set with large probability. And $R_{AL}^* \approx -\log \tilde{p}$.

How to express \tilde{p} .

$$\begin{aligned} -\log \tilde{p} &= -\log \prod_{i=1}^Q p_i^{p_i} \\ &= -\sum_{i=1}^Q p_i \log p_i =: \text{entropy} = H. \end{aligned}$$

Some properties of H :

1. its unit is bits
2. $H \geq 0$.
3. $H = 0 \iff p_q = 1$ for some q .
4. $H \leq \log_2 Q$.
5. $H = \log_2 Q \iff p_q = \frac{1}{Q}$ for all q .

Identify the set that WLLN says has probability $\rightarrow 1$. Suppose X_1, X_2, \dots i.i.d. real-valued variables.

$$T = \{\underbrace{x_1 \dots x_n}_{\underline{x}} \in A_X^N : \frac{1}{N} \sum_{i=1}^N x_i \doteq \bar{x} \pm \varepsilon\}$$

is called a typical set. $\Pr(\underline{X} \in T) \approx 1$ when N is large.

Now suppose X_1, X_2, \dots i.i.d. A_x -valued random variables, function $g : A_x \rightarrow \mathbb{R}$. Con-

sider Y_1, Y_2, \dots with $Y_i = g(X_i)$. Y_i 's are i.i.d. random variables.

If $\mathbb{E}[g(X)]$ is finite then we can apply WLLN that

$$\Pr\left(\frac{1}{N} \sum_{i=1}^N Y_i \doteq \mathbb{E}[Y] \pm \varepsilon\right) \rightarrow 1 \quad \text{as } N \rightarrow \infty.$$

Typical sequences wrt g :

$$T_{x, p_\lambda, g, \varepsilon}^N = \left\{ \underline{x} : \frac{1}{N} \sum_{i=1}^N g(x_i) \doteq \overline{g(X)} \pm \varepsilon \right\}.$$

If $\mathbb{E}[g(X)]$ is finite then by WLLN we have

$$\Pr(\underline{X} \in T_g) \rightarrow 1 \quad \text{as } N \rightarrow \infty.$$

Example 1.1.1 (Indicator function). Suppose $F \subset A_X$, and $g(x) = \begin{cases} 1 & x \in F \\ 0 & x \notin F \end{cases}$. Then

$\frac{1}{N} \sum_{i=1}^N g(x_i) = f_F(x)$. Now

$$T_g = \{\underline{x} : f_F(x) \doteq \Pr(X \in F) \pm \varepsilon\}.$$

By WLLN,

$$\Pr(\underline{X} \in T_g) \rightarrow 1 \quad \text{as } N \rightarrow \infty, \implies \Pr(n_F(\underline{X}) \doteq \mathbb{E}[x] \pm \varepsilon) \rightarrow 1.$$

Example 1.1.2. $A_x = \mathbb{R}$, $g(x) = x^2$. $T_g = \{\underline{x} : \}$

Theorem 1.1.2. Now suppose M functions g_1, g_2, \dots, g_M . Fix ε . Then

$$T_{g_1, g_2, \dots, g_M} = \bigcap_{i=1}^M T_{g_i}.$$

$$\Pr(\underline{X} \in T_{g_1, g_2, \dots, g_M}) \rightarrow 1 \quad \text{as } N \rightarrow \infty.$$

Proof.

$$\begin{aligned}
\Pr(\underline{X} \notin T_{g_1, g_2, \dots, g_M}) &= \Pr\left(\underline{X} \in \left(\bigcap_{i=1}^M T_{g_i}\right)^c\right) \\
&= \Pr\left(\underline{X} \in \left(\bigcup_{i=1}^M T_{g_i}^c\right)\right) \\
&\leq \sum_{i=1}^M \Pr(\underline{X} \in T_{g_i}^c) \rightarrow 0 \quad \text{as } N \rightarrow \infty. \quad \blacksquare
\end{aligned}$$

IMPORTANT APPLICATION

Suppose $A_x = \{a_1, \dots, a_Q\}$ a finite alphabet with probability p_1, \dots, p_Q . The $g_q(x)$ be the indicator of a_q . $T_q = \{\underline{x} : f_q(\underline{x}) \doteq p_q \pm \varepsilon\}$. And $\tilde{T} = \bigcap_{i=1}^Q T_i = \{\underline{x} : \forall q, f_q(x) \doteq p_q \pm \varepsilon\}$. $\tilde{T}_{X, p_X, \varepsilon}^N$ very typical sequence. We have

$$\Pr(\underline{X} \in \tilde{T}) \rightarrow 1 \quad \text{as } N \rightarrow \infty.$$

If $\underline{x} \in \tilde{T}$, then $\underline{x} \in \tilde{T}_g$ for any other g . Consider any real-valued g . If $\underline{x} \in \tilde{T}_\varepsilon$ then $\underline{x} \in T_{g, \varepsilon c}$ for some c .

$$\frac{1}{N} \sum_{i=1}^N g(x_i) = \sum_{q=1}^Q \frac{n_q(x)}{N} g(Q_q) = \sum_{q=1}^Q (p_q \pm \varepsilon) q(a_q) = \mathbb{E}[g(X)] + \varepsilon \sum_{q=1}^Q g(Q_q)$$

$$\Pr(\underline{X} \in \tilde{T}) \rightarrow 1 \text{ as } N \rightarrow \infty.$$

If $\underline{x} \in \tilde{T}$,

$$\begin{aligned}
p(\underline{x}) &= p(x_1)p(x_2) \dots p(x_N) \\
&= p_1^{n_1(\underline{x})} \dots \\
&= p_1^{f_1(\underline{x})N} \dots \\
&\doteq p_1^{(p_1 \pm \varepsilon)N} \dots \\
&\doteq 2^{N(\sum_{q=1}^Q p_q \log p_q \pm \varepsilon \sum_{q=1}^Q \log p_q)} \quad \doteq 2^{-NH \pm N\varepsilon c}
\end{aligned}$$

Theorem 1.1.3 (Shannon-McMillian Theorem). Suppose X_1, X_2, \dots i.i.d, $A_x = \{a_1, \dots, a_Q\}$ with probability p_1, \dots, p_Q . Then

1.

$$\Pr(\tilde{X} \in \tilde{T}_\varepsilon^N) \rightarrow 1 \text{ as } N \rightarrow \infty.$$

2. If $\underline{x} \in \tilde{T}_\varepsilon^N$, $p(\underline{x}) \doteq 2^{-NH \pm N\varepsilon c}$.
3. $|\tilde{T}_\varepsilon^N| \doteq \Pr(\underline{X} \in \tilde{T}_\varepsilon^N) 2^{N(H \pm \varepsilon c)}$.

Proof. ■

1.2 Shannon-McMillian Theorem

Is \tilde{T} essentially the smallest set with probability ≈ 1 ?

Yes. Let $S \in A_x^N$. We have

$$\Pr(\underline{X} \in S) = \Pr(X \in S \cap \tilde{T}) + \Pr(X \in S \cap \tilde{T}^c) \doteq |S \cap \tilde{T}| 2^{-NH \pm 2N\varepsilon c} + \Pr(\tilde{T}^c) \rightarrow 0 \text{ as } N \rightarrow \infty.$$

Theorem 1.2.1. For every $\varepsilon > 0$, there is a sequence $b_{\varepsilon,1}, b_{\varepsilon,2}, \dots$ s.t. $b_{\varepsilon,N} \rightarrow 0$ as $N \rightarrow \infty$, $b_{\varepsilon,B} \geq 0$.

For any N and any $S \subset A_X^N$,

$$|S| \geq (\Pr(\underline{X} \in S) - b_{\varepsilon,N}) 2^{NH - N\varepsilon c}.$$

An in hindsight shortcut

Let us directly consider

$$\begin{aligned} T_{S,\varepsilon}^N &= \left\{ \underline{x} : p(\underline{x}) \doteq 2^{-N(H \pm \varepsilon)} \right\} \\ &= \left\{ \underline{x} : -\frac{1}{N} \log p(\underline{x}) \doteq H \pm \varepsilon \right\} \\ &= \left\{ \underline{x} : -\frac{1}{N} \sum_{i=1}^N \log p(x_i) \doteq H \pm \varepsilon \right\} \end{aligned}$$

compare \tilde{T}_ε^N and $T_{s,\varepsilon}^N$.

Claim: $\tilde{T}_\varepsilon^N \subset T_{s,\varepsilon}^N$ where $c = -$; $\sum_{q=1}^Q \log p_q$.

Suppose $\underline{x} \in \tilde{T}_\varepsilon^N$. Show if it is also in $T_{s,\varepsilon}^N$. Check the following $p(x) \doteq 2^{-NH \pm N\varepsilon c}$, $-\log p(x) \doteq NH \pm N\varepsilon c$.

$$\begin{aligned}
-\log p(\underline{x}) &= -\log \prod_{i=1}^N p(x_i) \\
&= -\log \prod_{q=1}^Q p_q^{n_q(\underline{x})} \\
&= -\log \prod_{q=1}^Q p_q^{N f_q(\underline{x})} \\
&\doteq -\log \prod_{q=1}^Q p_q^{N(p_q \pm \varepsilon)} \\
&\doteq -\sum_{q=1}^Q N(p_q \pm \varepsilon) \log p_q \\
&\doteq NH \pm N\varepsilon \sum_{q=1}^C \log p_k \\
&\doteq NH \pm N\varepsilon c.
\end{aligned}$$

Extreme example:

$$A_x = \{0, 1\}, p_0 = p_1 = \frac{1}{2}, H = 1.$$

$$p(\underline{x}) = 2^{-N}.$$

$$T_{s,\varepsilon}^N = \{\underline{x} : p(\underline{x}) = 2^{-N(H \pm \varepsilon)} = 2^{-N}\} = A_X^N.$$

$$\tilde{T}_\varepsilon^N = \{\underline{x} : n_1(\underline{x}) \doteq N(\frac{1}{2} + \varepsilon)\}.$$

$$|T_{s,\varepsilon}^N| \doteq 2^{N(H \pm \varepsilon)}, |\tilde{T}_\varepsilon^N| \doteq 2^{N(H \pm 2\varepsilon c)}.$$

T_s is called probability typical. \tilde{T} is called frequency typical.

$$\text{Example } A_x = \{0, 1\}, p_1 = \frac{1}{4}, p_0 = \frac{3}{4}, \tilde{T}_\varepsilon^N = \{\underline{x} : f_1(\underline{x}) \doteq \frac{1}{4} + \varepsilon\}.$$

$$T_{s,\varepsilon}^T = \left\{ \underline{x} : f_1(\underline{x}) = \frac{1}{4} \pm N\varepsilon \log \frac{1-p_1}{p_1} \right\}$$

Typical sequences for an infinite alphabet

There are two cases: A_x is countably infinite / random variables are continuous

In the first case, frequency typical approach doesn't work. Probabilistic typical approach works just as is. $H = -\sum_{q=1}^{\infty} p_q \log p_q$ can be infinite.

Let $S_{\delta,N}$ = size of the smallest set of N sequences from A_x with probability at least $1 - \delta$.

Then for any $0 < \delta < 1$ and any h , $\frac{S_{\delta,N}}{2^{Nh}} \rightarrow \infty$ as $N \rightarrow \infty$.

1.3 Fixed Length to Variable Length (FVB) Lossless Source codes

Recall that FFB perfectly lossless has $R_{PL}^* = \log_2 |A_x|$, and FFB almost lossless has $R_{AL}^* = H$.

FVB perfectly lossless $R_{VL}^* \leq \log_2 |A_x|$.

Suppose we have a source with $A_x = \{a, b, c, d\}$ with probability $\{\frac{1}{2}, \frac{1}{4}, \frac{1}{8}, \frac{1}{8}\}$.

$p(u)$	u	code1	code2	code3	code4	code5	code6
$\frac{1}{2}$	a	00	0	0	0	0	0
$\frac{1}{4}$	b	01	10	10	10	1	01
$\frac{1}{8}$	c	10	110	10	11	01	011
$\frac{1}{8}$	d	11	111	11	111	10	0111
Rate		2	1.75	1.5	1.625	1.25	1.875

We can see that code 3-5 are all bad.

Code 6 has an advantage that you know 0 represents the start of a codeword. We will see later why (Example 1.4.2).

FVB source code is characterized by

- source length k
- codebook of binary codewords $C = \{v_1, v_2, \dots, v_{Q^k}\}$, $Q = |A_U|$.
- encoding rule $\alpha : A_U^K \rightarrow C$
- decoding rule $\beta : C \rightarrow A_U^K$.

The encoder operates in block fashion. The decoder does not.

Distinguish codes that look like code2 and codes that look like code6.

Definition 1.3.1. A codebook C is *prefix-free* if no codeword is the prefix of another.

A prefix-free code is called a prefix code. We will stick to prefix codes until states otherwise. (instantaneously decodable)

We like to draw binary tree diagrams of code.

Code 1:

233

A prefix is perfectly lossless if and only if α is 1-to-1. The rate: $\bar{r}(c) = \frac{\bar{L}}{K} = \frac{1}{K} \sum_u p(u) L(u)$

(length of codeword assigned to \underline{u})

$$R_{VL}^*(k) = \min \{ \bar{r}(c) : c \text{ is perfectly lossless FVB with source length } k \}.$$

$$R_{VL}^* = \inf \{ \bar{r}(c) : c \text{ is PL FVB prefix code with any source length} \} = \inf_K R_{VL}^*(k).$$

How does one design a prefix code to have small or smallest rate?

Focus first $k = 1$. Shannon's idea: $L_q \approx -\log_2 p_q$.

$$\sum_{q=1}^Q p_q L_q \approx -\sum_{q=1}^Q p_q \log p_q = H.$$

Question. Is there a prefix code with $L_q \approx -\log p_q$ for $q = 1, 2, \dots, Q$? Could there be prefix codes with even smaller rate?

Theorem 1.3.1 (Kraft inequality theorem). *There is a binary prefix code with length L_1, L_2, \dots, L_Q iff the Kraft sum*

$$\sum_{q=1}^Q 2^{-L_q} \leq 1.$$

Proof. Suppose $\underline{v}_1, \dots, \underline{v}_Q$ is a prefix code with length L_1, \dots, L_Q . Let $L_{\max} = \max_q L_q$.

From the tree, the number of sequences of length L_{\max} prefixed by any codeword, is $\sum_{q=1}^Q 2^{L_{\max}-L_q} \leq 2^{L_{\max}} \implies \sum_{q=1}^Q 2^{-L_q} \leq 1$. So the Kraft inequality holds. ■

Now suppose

$$L_q = \lceil -\log_2 p_q \rceil, q = 1, \dots, Q. \quad (1.3.1)$$

Is there a code with these lengths? Check Kraft.

$$\begin{aligned} \sum_{q=1}^Q 2^{-L_q} &= \sum_{q=1}^Q 2^{-\lceil -\log p_q \rceil} \\ &\leq \sum_{q=1}^Q 2^{-(-\log p_q)} \\ &\leq \sum_{q=1}^Q 2^{\log p_q} \\ &\leq \sum_{q=1}^Q p_q = 1. \end{aligned}$$

So the Kraft inequality holds. \exists a prefix code with length L_1, \dots, L_Q given by (1.3.1), called Shannon-Fano code.

Now the question is how good is this Shannon-Fano Code?

For the Shannon-Fano code, the rate (average length) is

$$\bar{L}_{SF} = \sum_{q=1}^Q p_q \lceil -\log p_q \rceil.$$

We have the following bounds:

$$H = \sum_{q=1}^Q p_q (-\log p_q) \leq \bar{L}_{SF} < \sum_{q=1}^Q p_q (-\log p_q + 1) = H + 1.$$

Question. Can we do better now?

We will show that $\bar{L} \geq H$ for any prefix code.

Let C be a prefix code with length L_1, \dots, L_Q . Take the difference $\bar{L} - H = \sum_{q=1}^Q p_q L_q + \sum_{q=1}^Q p_q \log p_q$.

$$\begin{aligned} \bar{L} - H &= \sum_{q=1}^Q p_q L_q + \sum_{q=1}^Q p_q \log p_q \\ &= - \sum_q p_q \log \frac{2^{-L_q}}{p_q} \\ &= - \sum_q p_q \ln \frac{2^{-L_q}}{p_q} \frac{1}{\ln(2)} \\ &\geq - \sum_q p_q \left(\frac{2^{-L_q}}{p_q} - 1 \right) \frac{1}{\ln(2)} \\ &\geq - \frac{1}{\ln(2)} \sum_q 2^{-L_q} + \sum_q p_q \frac{1}{\ln(2)} = \frac{1}{\ln 2} (1 - 1) = 0. \end{aligned}$$

In homework we will that that \bar{L} can get very close to $H + 1$.

Now allow $k \geq 1$. WE have a $C = \{\underline{v}_1, \dots, \underline{v}_{Q^k}\}$ of length L_1, \dots, L_{Q^k} . We want small

$$\bar{r}(c) = \frac{\bar{L}}{K} = \frac{\sum_u p(\underline{u}) L(\underline{u})}{k}.$$

Shannon-Fano code achieve that

$$H^k \leq \bar{L}_{SF} < H^k + 1 \implies \frac{H^k}{k} \leq \bar{r}_{SF} = \frac{\bar{L}_{SF}}{k} < \frac{H^k}{k} + \frac{1}{k}.$$

Since $H^k = kH$ we have

$$H \leq \bar{r}_{SF} < H^k + \frac{1}{k}.$$

Similarly we have for any prefix code, we have

$$\bar{r} = \frac{\bar{L}}{K} \geq \frac{H^k}{k} = H.$$

This leads to a new coding theorem.

Theorem 1.3.2. *Given i.i.d. source U with alpha A_U and entropy H . We have*

1. *For every k , $R_{VL}^* \leq R_{VL}^*(k) < H + \frac{1}{k}$.*
2. *For every k , $R_{VL}^*(k) \geq R_{VL}^* \geq H$.*

Combined we have

$$\forall k \in \mathbb{Z}_{>0}, H \leq R_{VL}^*(k) < H + \frac{1}{k}$$

and

$$R_{VL}^* = H.$$

1.4 Huffman's Code Design

Given p_1, \dots, p_Q , it finds a prefix code with smallest \bar{L} .

Algorithm 1.4.1 Huffman Code

Input: Alphabet probability $\{p_i | i = 1, \dots, Q\}$, WLOG assume $p_1 \geq p_2 \geq \dots \geq p_Q$.

Output: FVB Codebook for alphabet $\{a_i | i = 1, \dots, Q\}$.

```

1: function HUFFMAN( $P_Q = \{p_i | i = 1, \dots, Q\}$ )
2:   if  $Q = 2$  then return  $\{0, 1\}$ 
3:   end if
4:    $p'_{Q-1} \leftarrow p_{Q-1} + p_Q$ 
5:    $P_{Q-1} \leftarrow (P_Q \setminus \{p_{Q-1} + p_Q\}) \cup \{p'_{Q-1}\}$ 
6:    $c_{Q-1} \leftarrow \text{HUFFMAN}(P_{Q-1}) =: \{\underline{v}_1, \dots, \underline{v}_{Q-1}\}$ 
7:    $c_Q \leftarrow \{\underline{v}_1, \dots, \underline{v}_{Q-2}, \underline{v}_{Q-1}0, \underline{v}_{Q-1}1\}$ 
8: end function
```

Proposition 1.4.1. *If c_{Q-1} is optimal for P_{Q-1} then c_Q is optimal for P_Q .*

Example 1.4.1.

We found that

But there is a tighter upper bound

$$\bar{L}^* \leq \begin{cases} H + p_{\max} & p_{\max} < \frac{1}{2} \\ H + p_{\max} + 0.086 & p_{\max} \geq \frac{1}{2}. \end{cases}$$

Hence

$$\mathcal{R}_{VL}^*(k) \leq \begin{cases} H + \frac{p_{\max}^k}{k} & (p_{\max})^k < \frac{1}{2} \\ H + \frac{p_{\max}^k}{k} + \frac{0.086}{k} & (p_{\max})^k \geq \frac{1}{2}. \end{cases}$$

Up till now we've only focused on i.i.d RV's. Now suppose RV's are dependent, then

$$\frac{H^k}{k} < \frac{kH}{k}.$$

For a stationary random process,

$$\frac{H^k}{k} \searrow H_{\infty}.$$

For example, English has

$$H^1 \approx 4.08, H_{\infty} \approx 1.$$

The bits produced by a good lossless source code ($\bar{r} \approx H$) are approximately i.i.d. equiprobable.

Synchronization and transmission entropy

Example 1.4.2. Suppose $\{01, 001, 101, 110\}$ for $\{a, b, c, d\}$.

$$\underline{u} = dddddddd \dots \implies \underline{z} = 110110110110110110 \dots$$

$$(\text{if one leading 1 is missing}) \underline{z}' = 101101101101 \dots \implies \hat{u} = ccccccc \dots$$

Now if $\{1, 01, 001, 000\}$ for $\{a, b, c, d\}$. Then the same problem will not happen.

1.5 Buffering

Suppose the source is outputting at R symbols per second. The encoder would have $R\bar{r}$ bits per second.

Buffer overflow happens when a long sequence of low probability symbols are encoded.

Buffer underflow happens when a long sequence of high probability symbols are en-

coded. Buffer will be empty. Include an additional codeword in codebook called a “flag”. Insert this codeword when buffer becomes empty.

We focused prefix codes. There are some non prefix codes that can be decoded losslessly.

Definition 1.5.1. A code is *separable* if any finite sequence of codewords is different from any other finite sequence of codewords.

Remark. Prefix codes are separable. And determining if a non prefix code is separable is not easy.

Could prefix codes have smaller I ?

McMillian’s theorem says that Kraft inequality holds for separable codes. If you have a separable code whose codeword satisfy Kraft, then there is a prefix code with same lengths.

Lossless coding for source with infinite alphabet.

Suppose $A_n = \{1, 2, 3, \dots\}$.

1. FFB codes can’t have finite rate if perfectly lossless
2. AL FFB then SM theorem
3. FVB. Current approach is based on Kraft inequality. It still holds for infinite case. (See Appendix). But Huffman’s optimal design does not apply.

Other forms of variable length lossless source codes.

1. Run-length coding
2. Dictionary coding

1.6 Universal Source Coding

Suppose you are to encode 10^6 symbols from alphabet $A = \{a, b, c, d\}$. We can calculate $n_a(\underline{u}), n_b(\underline{u}), n_c(\underline{u}), n_d(\underline{u})$ and similarly, frequencies. Then we can apply Huffman or Shannon-Fano code.

Chapter 2

Entropy

The star of this chapter is

$$\ln x \leq x - 1$$

2.1 Entropy

Entropy

$$H := - \sum_x p(x) \log p(x)$$

is a measure of randomness or uncertainty.

$$H_q := - \sum_x p(x) \log_q p(x), \quad H_q = H_r \frac{1}{\log_r q}.$$

$$H(X) = \mathbb{E}[\log p_X(X)]$$

Remark. 1. $H(X) \geq 0$ and $H(X) = 0 \iff p(x) = 1$ for some x .

2. $H(X) = \infty$ if X is continuous or has continuous component.

3. $H(X, Y) \geq H(X)$

4. $H(X, Y) = H(X) + H(Y)$ if X and Y are independent.

5. $H(X_1, X_2, \dots, X_n) = H(X_{\sigma(1)}, \dots, X_{\sigma(n)})$ for any $\sigma \in S_n$.

Divergence is a measure of dissimilarity of two probability distribution.

Definition 2.1.1. Suppose p and q are probability mass functions. The *divergence* from p

to q is

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

Remark. 1. $p = q \implies D(p\|q) = 0$

2. It is *not symmetric*. $D(p\|q) \neq D(q\|p)$. You can make symmetric by taking the sum, but then it is not nicely related to information theory.

What if $p(x) = 0$ for some x ? We take $0 \log \frac{0}{q(x)} = 0$. So if $\exists x$ s.t. $q(x) = 0$ and $p(x) \neq 0$. Then $D(p\|q) = \infty$.

When alphabet A_x is infinite, $D(p\|q)$ can be ∞ even when $p(x) > 0$ and $q(x) > 0$ for all $x \in A_x$.

Is $\sum_x p(x) \log \frac{p(x)}{q(x)}$ always defined? Write

$$\sum_x p(x) \log \frac{p(x)}{q(x)} = \sum_{x, p(x) > q(x)} \log \frac{p(x)}{q(x)} + \sum_{x, p(x) < q(x)} \log \frac{p(x)}{q(x)}$$

We will show later that the second term is never $-\infty$, so it is always well-defined.

Proposition 2.1.1 (Divergence inequality). *For any p, q ,*

$$D(p\|q) \geq 0, D(p\|q) = 0 \iff p = q.$$

Proof.

$$\begin{aligned} D(p\|q) &= \sum_x p(x) \log \frac{p(x)}{q(x)} \\ &= \sum_x p(x) \ln \frac{p(x)}{q(x)} \frac{1}{\ln 2} \\ &= - \sum_x p(x) \ln \frac{q(x)}{p(x)} \frac{1}{\ln 2} \\ &\leq - \sum_x p(x) \frac{q(x) - p(x)}{p(x)} \frac{1}{\ln 2} \\ &\leq - \left(\sum_x p(x) - q(x) \right) \frac{1}{\ln 2} = 0 \end{aligned}$$

For first equality, \Leftarrow is clear. Now suppose $D(p\|q) = 0$. Then

$$\ln \frac{q(x)}{p(x)} = \frac{q(x)}{p(x)} - 1 \implies p(x) = q(x) \text{ for all } x \text{ with } p(x) > 0.$$

■

Let's rewrite the divergence inequality a little bit.

$$0 \leq D(p||q) = \sum_x p(x) \log \frac{p(x)}{q(x)} = -H(x) - \sum_x p(x) \log(q_x)$$

$$H(x) \leq - \sum_x p(x) \log(q_x)$$

with $= \iff p = q$. $-\sum_x p(x) \log(q_x)$ is called cross entropy.

Definition 2.1.2. The *cross entropy* of p with respect to q is

$$H_c(p, q) := - \sum_x p(x) \log q(x).$$

Cross-entropy inequality: for any p, q

$$H_p(X) \leq H_c(p, q)$$

with $= \iff p = q$.

Remark.

$$D(p||q) = H_c(p, q) - H_p(X) \iff H_c(p, q) = H_p(X) + D(p||q).$$

Definition 2.1.3. Variation distance

$$V(p, q) = \sum_x |p(x) - q(x)|$$

How does $D(p||q)$ compare to $V(p, q)$?

Proposition 2.1.2 (Pinskev's inequality).

$$V(p, q) \leq \sqrt{(2 \ln 2) D(p||q)}.$$

So small $D(p||q) \implies$ small $V(p, q)$ even when $|A_x| = \infty$. On the other hand the converse is not true.

Lemma 2.1.1. If there exists $0 < \delta < 1$ such that $\frac{|p(x)-q(x)|}{p(x)} \leq \delta$ for all x s.t. $p(x) > 0$ then

$$D(p||q) \leq \frac{\delta}{1-\delta} \frac{1}{\ln 2}$$

If $D(p||q) \approx 0$ then $p \approx q$, meaning $V(p, q) \approx 0$.

If p, q are percentage wise close, then $D(p||q) \approx 0$.

Log-sum inequality

Suppose $u_1, u_2, \dots, u_n, v_1, \dots, v_n$ nonnegative. Then

$$\sum_i u_i \log \frac{u_i}{v_i} \geq \left(\sum_i u_i \right) \log \frac{\sum_i u_i}{\sum_i v_i}.$$

This is a generalization of divergence inequality.

2.2 Basic Properties of Entropy

Proposition 2.2.1.

$$H(X^N) \leq \sum_{i=1}^N H(X_i)$$

with $= \iff X_1, \dots, X_N$ are independent.

Proof.

$$H(X^N) \leq H_L(p, q)$$

where p is probability mass function of X_1, \dots, X_N . Choose

$$q(x_1, \dots, x_n) = p(x_1)p(x_2) \dots p(x_n).$$

■

Since we are dealing with discrete random variables, it is useful to think about probability mass function as a set of probabilities $\{p_1, p_2, \dots\}$. Write

$$H(p_1, p_2, \dots) = - \sum_i p_i \log p_i$$

Let $p'_i = p_i + p_j$, replace p_i, p_j with p'_i and leave all others the same. We have

$$-p_i \log p_i - p_j \log p_j \geq -(p_i + p_j) \log(p_i + p_j)$$

i.e. entropy decreases when two probabilities are merged.

Proposition 2.2.2. *If X is Q -ary with $Q < \infty$ then*

$$H(X) \leq \log_2 Q$$

Proof. Let $q(x) = \frac{1}{Q}$.

$$\begin{aligned} H(X) &\leq H_c(p, q) = - \sum_x p(x) \log p(x) \\ &= \sum_x p(x) \log_2 Q = \log_2 Q. \end{aligned}$$

■

Proposition 2.2.3. *Suppose $Y = g(X)$. Then*

$$H(Y) = H(g(X)) \leq H(X)$$

and $= \iff g$ is one-to-one (probabilistically).

Proof.

$$H(Y) = \sum_y p(y) \log p(y)$$

where $p(y) = \sum_{x, g(x)=y} p(x)$.

■

$p_i = q^i(1 - q), i = 0, 1, 2, \dots$, then

$$\begin{aligned}
H(X) &= - \sum_{i=0}^{\infty} p_i \log p_i \\
&= - \sum_{i=0}^{\infty} q^i(1 - q) \log q^i(1 - q) \\
&= - \sum_{i=0}^{\infty} q^i(1 - q) (i \log q + \log(1 - q)) \\
&= - \log q \sum_{i=0}^{\infty} q^i(1 - q)i - \sum_{i=0}^{\infty} q^i(1 - q) \log(1 - q) \\
&= - \log q \cdot q(1 - q) \frac{d}{dq} \sum_{i=0}^{\infty} q^i - \log(1 - q) \\
&= - \log q \cdot q(1 - q) \frac{1}{(1 - q)^2} - \log(1 - q) \\
&= - \log q \cdot \frac{q}{1 - q} - \log(1 - q) = \frac{-q \log(q) - (1 - q) \log(1 - q)}{1 - q} \\
&= H(Q) \frac{1}{1 - q} < \infty.
\end{aligned}$$

$p_i = \frac{\alpha}{i(\ln i)^2}, i = 2, 3, \dots$, then

$$\begin{aligned}
H(X) &= - \sum_{i=2}^{\infty} \frac{\alpha}{i(\ln i)^2} \log \left(\frac{\alpha}{i(\ln i)^2} \right) \\
&=
\end{aligned}$$

2.3 Conditional Entropy

$$H(X | Y) = \sum_{x,y} p(x, y) \log p(x|y) \geq 0$$

with equality iff X is a function of Y .

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

with equality iff X, Y are independent.

chain rule:

$$H(X, Y) = H(X) + H(Y | X) \implies H(X, Y) \geq H(X)$$

Conditional lossless source coding

2.4 Convexity

Goal: entropy is a concave (convex \cap).

Extended definition of entropy

Definition 2.4.1.

$$\overline{H}(x) = \sup_{\text{finite quantizers } Q} H(Q(x))$$

where finite quantizer is a function $Q : A \rightarrow B, |\{Q(x) : x \in A\}| < \infty$.

This gives normal definition for discrete random variables and ∞ for continuous and mixed random variables

Chapter 3

Information

3.1 Information

Not a good question: How much information is there in X ? Better questions: The information in X about random variable Y ?

Definition 3.1.1. The (mutual) information given by X about Y is defined as

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

Definition 3.1.2. Y = outcome of a fair 6-sided die. X = oddity of the outcome. We have

$$H(Y) = \log_2 6, H(Y | X) = \log_2 3 \implies I(X; Y) = \log_2 6 - \log_2 3 = \log_2 2 = 1.$$

Lemma 3.1.1. Suppose X, Y are discrete random variables. Then

1. $I(X; Y) \geq 0, = 0$ iff X, Y independent.

2. We have alternate formulas

$$\begin{aligned}
 I(X; Y) &= - \sum_y p(y) \log p(y) + \sum_{x,y} p(x, y) \log p(y|x) \\
 &= - \sum_{x,y} p(x, y) \log p(y) + \sum_{x,y} p(x, y) \log p(y|x) \\
 &= \sum_{x,y} p(x, y) \log \frac{p(y|x)}{p(y)} \\
 &= \sum_{x,y} p(x, y) \log \frac{p(x, y)}{p(x)p(y)}
 \end{aligned}$$

This shows that I is symmetric: $I(X; Y) = I(Y; X)$.

3. From above,

$$\begin{aligned}
 I(X; Y) &= - \sum p(x, y) \log p(x) - \sum p(x, y) \log p(y) + \sum p(x, y) \log p(x, y) \\
 &= H(X) + H(Y) - H(X, Y).
 \end{aligned}$$

4. We can view information as an expectation:

$$I(X; Y) = \mathbb{E} \left[\log \frac{p(Y | X)}{p(Y)} \right] = \mathbb{E} \left[\log \frac{p(X | Y)}{p(X)} \right] = \mathbb{E} \left[\log \frac{p(X, Y)}{p(X)p(Y)} \right].$$

5. We can view information with respect to divergence:

$$I(X, Y) = D(p_{XY} \| p_X p_Y).$$

Remark. Alternate notation $I_\phi(X, Y), I_{X;Y}(p), I(p)$.

What happens if $p(x, y) = 0$, or $p(x) = 0$, or $p(y) = 0$? $p(x) = 0$ or $p(y) = 0 \implies p(x, y) = 0$.

Remark. $I(X; Y)$ is possible. Suppose $H(Y) = \infty$ and Y is a function of X . Then $I(X; Y) = H(Y) - H(Y | X) = \infty - 0 = \infty$.

Information for more variables

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

Relations between information entropy

$$1. I(X; Y) = H(X) - H(X | Y) = H(Y) - H(Y | X).$$

2. $I(X; Y) \leq H(X)$, = iff X is a function of Y .
3. $I(X; X) = H(X)$.
4. $I(X; g(X)) = H(g(X))$.

3.2 Conditional Information

Recall that we have two concepts of conditional entropy.

$$H(X | Y = y), H(X | Y) = \sum_y p(y) H(X | Y = y)$$

We are going to use the same approach for conditional information.

Definition 3.2.1. Suppose X, Y, Z are two discrete random variables,

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

$$I(X; Y | Z) = \sum_z p(z) I(X; Y | Z = z).$$

Lemma 3.2.1. 1. $I(X; Y | Z = z) \geq 0$

2. $I(X; Y | Z) \geq 0$, = 0 iff X, Y are conditionally independent given Z .
3. $I(X; Y | Z) = H(X | Z) - H(X | Y, Z)$.
4. $I(X; Y, Z) = I(X; Z) + I(X; Y | Z)$

chain rule of conditional information

$$I(X; YZ | U) = I(X; Z | U) + I(X; Y | Z, U).$$

3.3 Cryptography From Information Perspective

$A_X = A_K$ and $|A_X| = |A_K| = 2^N$. $p_K(k) = 2^{-N}$, $k \in A_K$.

If $|K| < |A_X|$ then the crypto system is not perfect.

Fix $x \in |A_X|$. For each y we have $P(Y = y | X = x) = P(Y = y) > 0$. Therefore for each y , there must be some key $k \in K$ such that $y = e_K(x)$. It follows that $|K| \geq |Y|$. The encryption is injective giving $|Y| \geq |A_X|$.

3.4 Continuous Random Variables and Information

Definition 3.4.1. A random variable X is continuous if $\Pr(X = x) = 0$ for all x .

We assume alphabet is \mathbb{R} , X is absolutely continuous

Example 3.4.1. 1. Uniform

$$p(x) = \begin{cases} \frac{1}{b-a} & a \leq x \leq b \\ 0 & \text{otherwise} \end{cases}$$

2. Gaussian

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

3. Laplacian

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

4. Exponential

$$p(x) = \begin{cases} \frac{\sqrt{2}}{\sigma} e^{-\frac{\sqrt{2}}{\sigma}|x|} & x \geq 0 \\ 0 & x < 0 \end{cases}$$

Definition 3.4.2. The *support* of a random variable X or of its probability distribution is defined by

$$S := \{x : \Pr(X \in x \pm \varepsilon) > 0, \forall \varepsilon > 0\}.$$

Definition 3.4.3. Conditional probability

$$\Pr(F | X = x) = \frac{\Pr(F, X = x)}{\Pr(X = x)}.$$

When X is continuous,

$$\Pr(F | X = x) = \lim_{\delta \rightarrow 0} \Pr(F | X \in x \pm \delta)$$

Suppose $Y = 3X$

$$1 = \Pr(Y = 3 | X = 1) = \lim_{\delta \rightarrow 0} \frac{\Pr(Y = 3 | X \in 1 \pm \delta)}{\Pr(X \in 1 \pm \delta)} = 0, \text{ contradiction}$$

So we use

$$\Pr(Y \in F | X = x) = \lim_{\varepsilon \rightarrow 0} \lim_{\delta \rightarrow 0} \Pr(Y \in F_\delta | X \in x \pm \varepsilon)$$

where

$$F_\delta = \{y \mid \|y - y'\| \text{ for some } y' \in F\}.$$

Generalized sum

3.5 Differential Entropy

Recall that for continuous random variable, the entropy is ∞ , which is not very interesting. We recall the definition of information

$$I(X; Y) = \int \int p(x, y) \frac{\log p(x|y)}{\log p(x)} dx dy.$$

We can rewrite it as

$$\begin{aligned} I(X; Y) &= - \int \int p(x, y) \log p(x) dx dy + \int \int p(x, y) \log p(x|y) dx dy \\ &= - \int p(x) \log p(x) dx + \int \int p(x, y) \log p(x|y) dx dy \end{aligned}$$

The first term is defined as differential entropy, denoted by $H_d(X)$, and the second term is conditional differential entropy, denoted by $H_d(X|Y)$.

We have

$$I(X; Y) = H_d(X) - H_d(X|Y).$$

Example 3.5.1. Suppose $p_X(x) = \begin{cases} \frac{1}{\Delta} & 0 \leq x \leq \Delta \\ 0 & \text{otherwise} \end{cases}$, then $H_d(X) = \log \Delta$.

Example 3.5.2. Suppose X is Gaussian with variance σ^2 , then

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

But $H_d(X)$ can be negative and even $-\infty$.

$H_d(X) \geq H_d(X|Y)$, $= 0$ iff X, Y independent.

When $Y = \alpha X$, $\alpha > 0$, then $H_d(Y) = H_d(X) + \log \alpha$.

Differential entropy is not a measure of true randomness and uncertainty. We'll see it is a relative measure.

Back to typical sequences.

Suppose X_1, X_2, \dots IID continuous rv's with pdf p , then

$$T_{s,\varepsilon}^N = \left\{ \underline{x} : \frac{1}{N} \sum_{i=1}^N s(x_i) \doteq \mathbb{E}[s(X)] = H_d(X) \right\}$$

And by LLN,

$$\Pr(\underline{X} \in T_{s,\varepsilon}^N) \rightarrow 1 \text{ as } N \rightarrow \infty, \quad p(\underline{x}) \doteq 2^{-NH_d(X) \pm N\varepsilon}, \text{ for } x \in T_{s,\varepsilon}^N.$$

$$\begin{aligned} \Pr(\underline{X} \in T_{s,\varepsilon}^N) &= \int_{T_{s,\varepsilon}^N} p(\underline{x}) \, d\mathbf{x} \\ &\doteq 2^{-NH_d(X) \pm N\varepsilon} \int_{T_{s,\varepsilon}^N} 1 \, d\mathbf{x} \\ &\doteq 2^{-NH_d(X) \pm N\varepsilon} \text{vol}(T_{s,\varepsilon}^N) \end{aligned}$$

We also have

$$|T_{s,\varepsilon}^N| \doteq \Pr(\underline{X} \in T_{s,\varepsilon}^N) 2^{NH_d(X) \pm 2N\varepsilon} \doteq 2^{NH_d(X) \pm 2N\varepsilon}$$

So $|T_{s,\varepsilon}^N| \cong 2^{NH_d(X)}$ is the size(vol) of a smallest set with probability ≈ 1 .

3.6 Properties of Differential Entropy

Differential entropy decreases as probabilities aggregate

Suppose X is supported on S , then $H_d(X) \leq \log |S|$. $H_d(aX + b) = H_d(X) + \log a$.

$H_{d,X}(p)$ is convex \cap .

Chapter 4

Estimation, Decision

4.1 Estimation Theory

4.2 Decision Theory

Bayes rule

$$\arg \max_x \frac{p_X(x)p_{Y|X}(y|x)}{p_Y(y)}$$

Example 4.2.1. $Y = X + N$ with X, N independent Gaussian variables with mean 0.

$$r^*(y) = \frac{\sigma_X^2}{\sigma_X^2 + \sigma_N^2} y$$

$$\bar{d} = \sigma_X^2(1 - \rho^2), \rho := \frac{\mathbb{E}[XY]}{\sigma_x \sigma_y}.$$

The MMSE estimate is linear.

In many situations, the MMSE estimation $r^*(y) = \mathbb{E}[X|Y = y]$ is too complicated or unknown. In such cases, we often seek the best linear (technically, affine) estimate

$$r_L^*(y) = ay + b$$

choose a, b to minimize $\mathbb{E}[(X - r_L(Y))^2]$.

The best linear estimate

$$r_L^*(y) = \rho \frac{\sigma_X}{\sigma_Y} (y - \mathbb{E}[Y] + \mathbb{E}[X]), \quad d_L^* = \sigma_X^2 (1 - \rho^2)$$

We don't need to know $p(x|y), p(y|x), p(x)$.

Bayes Estimation: assume $p_{XY}, d(x, \hat{x})$.

Much of statistics considers estimation when only $p_{Y|X}$ is known.

Maximum Likelihood Rule

$$r_{ML}(y) = \arg \min_x p_{Y|X}(y|x)$$

We do not have an average distortion.

Lemma 4.2.1 (Fano's Lemma, decision theory). $|A_X| < \infty$ finite alphabet. For any decision rule r for deciding from Y we have

$$H(X|Y) \leq \mathcal{H}(p_E) + p_E \log(Q - 1) =: \mathcal{F}(p_E).$$

$\mathcal{F}(p_E)$ (called Fano's function) peaks at $p_E = \frac{Q-1}{Q}$ with maximum $\log_2 Q$.

Consider its inverse $\mathcal{F}_Q^{-1}(p_E)$ on $[0, \frac{Q-1}{Q}]$. Since $H(X|Y) = H(X) - I(X; Y)$, then as $I(X; Y)$ affects the lower bound of error probability.

Make an observation Y , want to know X , $\hat{X} = r(Y)$ estimate decision.

$$\begin{aligned} H(X|Y = y) &\leq \mathcal{H}(p_{X|Y}(a|y)) + (1 - p_{X|Y}(a|y)) \log(Q - 1) \\ H(X|Y = y) &\leq \mathcal{H}(1 - p_{X|Y}(r(y)|y)) + (1 - p_{X|Y}(r(y)|y)) \log(Q - 1) \\ &= \mathcal{H}(\Pr(X \neq r(Y)|Y = y)) + \Pr(X \neq r(Y)|Y = y) \log(Q - 1) \\ &= \mathcal{F}_Q(\Pr(X \neq r(Y)|Y = y)) \end{aligned}$$

Meanwhile,

$$\begin{aligned} H(X|Y) &= \sum_y p(y) H(X|Y = y) \\ &\leq \sum_y \mathcal{F}_Q(\Pr(X \neq r(Y)|Y = y)) \\ &\leq \mathcal{F}_Q \left(\sum_y \Pr(X \neq r(Y)|Y = y) \right) \\ &= \mathcal{F}_Q(\Pr(X \neq r(Y))) = \mathcal{F}_Q(p_E). \end{aligned}$$

Fano lower bound to MSE

If X is real-valued and $r(Y)$ is an estimate of X then

$$H_d(X|Y) \leq \frac{1}{2} \log(2\pi e \bar{d}), \quad \bar{d} = \mathbb{E}[(X - r(Y))^2]$$

$$\bar{d} \geq \frac{1}{2\pi e} 2^{2H_d(X|Y)} = \frac{1}{2\pi e} 2^{2(H_d(X) - I(X;Y))}.$$

Proof. Recall that

$$H_d(X) \leq \frac{1}{2} \log(2\pi e \sigma_X^2).$$

Fix a value $Y = y$, then

$$\begin{aligned} H_d(X|Y = y) &\leq \frac{1}{2} \log 2\pi e \mathbb{E}[(X - \mathbb{E}[X|Y = y])^2 | Y = y] \\ &\leq \frac{1}{2} \log 2\pi e \mathbb{E}[(X - r(Y))^2 | Y = y] \end{aligned}$$

Averaging over y we have

$$\begin{aligned} H_d(X|Y) &= \sum_y p_Y(y) H_d(X|Y = y) \\ &\leq \sum_y p_Y(y) \frac{1}{2} \log \left(2\pi e \mathbb{E}[(X - r(Y))^2 | Y = y] \right) \\ &\leq \frac{1}{2} \log \left(2\pi e \sum_y p_Y(y) \mathbb{E}[(X - r(Y))^2 | Y = y] \right) = \frac{1}{2} \log(2\pi e \bar{d}) \quad \blacksquare \end{aligned}$$

Lower bound to per-letter error probability

Suppose we have a rule for deciding $X^N = (X_1, \dots, X_N)$ from Y with $|A_X| = Q$. Distortion = per letter probability

$$\bar{d}_L = \frac{1}{N} \sum_{i=1}^N \Pr(X_i \neq r_i(Y))$$

Fano lower bound for p_{LE}

$$\frac{1}{N} H(X^N|Y) \leq \frac{1}{N} \sum_{i=1}^N H(X_i|Y) \leq \mathcal{F}_Q(p_{LE}), \quad p_{LE} \geq \mathcal{F}_Q^{-1} \left(\frac{1}{N} \sum_{i=1}^N H(X_i|Y) \right).$$

Fano lower bound to per-letter MMSE If R is a rule for estimating X_1, \dots, X_N (real-

valued) from Y , then

$$\bar{d}_L = \frac{1}{N} \sum_{i=1}^N \mathbb{E}(X_i - r_i(Y))^2 \geq \frac{1}{2\pi e} 2^{2\frac{1}{N} \sum_{i=1}^N H_d(X_i|Y)}$$

Block converse

When an FFB block source code with source length K and code length L rate $\bar{r} = \frac{L}{K} \leq H$, then

$$p_{BE} \geq \mathcal{F}_{Q^K}^{-1}(K(H\bar{r}))$$

U_1, \dots, U_K , encoder, Z_1, \dots, Z_L , decoder, $\hat{U}_1, \dots, \hat{U}_K$.

$$H(U^K|Z^L) = H(U^K) - I(U^K; Z^L) \geq KH - K\bar{r},$$

By Fano's lemma,

$$p_{BE} \geq \mathcal{F}_{Q^K}^{-1}(H(U^K|Z^L)) \geq \mathcal{F}_{Q^K}^{-1}(K(H - \bar{r}))$$

per-letter converse

When an FFB block source code with source length K , code length L , rate $\bar{r} = \frac{L}{K} < H$, then

$$p_{LE} \geq \mathcal{F}_Q^{-1}(H - \bar{r})$$

Subsection 9.9.2: If a lossless source code has rate close to H then the bits it produces are approximately IID equiprobable $D(p_{Z_1, \dots, Z_n} \| p_{IID}) \approx 0$.

upper bound to error probability of an optimal decision rule

$$H(X) \geq -\log p_{\max}$$

$$H(X|Y = y) \geq -\log \max_x p_{X|Y}(x|y) = -\log p_{X|Y}(r^*(y)|Y = y)$$

Average over y ,

$$\begin{aligned}
 H(X|Y) &= \sum_y p_Y(y) H(X|Y=y) \\
 &\geq - \sum_y p_Y(y) \log p_{X|Y}(r(y)|Y=y) \\
 &\geq - \log \left(\sum_y p_Y(y) p_{X|Y}(r(y)|y) \right) \\
 &\geq - \log \Pr(r(Y) = X) = - \log(1 - p_E).
 \end{aligned}$$

$$H(X|Y) \geq - \log(1 - p_E) \implies p_E \leq 1 - 2^{-H(X|Y)}$$

Upper bound to MMSE for estimating X from Y .

Special case: $Y = X + V$, X, V independent.

We have

$$I(X; Y) = H_d(Y) - H_d(Y|X).$$

$$H_d(Y) \leq \frac{1}{2} \log 2\pi e \sigma_Y^2 = \frac{1}{2} \log 2\pi e (\sigma_X^2 + \sigma_V^2)$$

$$H_d(Y|X) = H_d(X + V|X) = H_d(V|X) = H_d(V) = \frac{1}{2} \log 2\pi e \sigma_V^2$$

channel code scenario

data \rightarrow encoder \rightarrow noisy channel \rightarrow decoder \rightarrow data reproductions

Data: bits from binary symmetric source (BSS) Z_1, Z_2, \dots, Z_i 's are independent, identical, binary $\{0, 1\}$, equiprobable.

Channel: Discrete time system, with input alphabet A_X , output alphabet A_Y , and a stochastic input/output characterized by transition distribution q .

Example: (BSC).

Example 4.2.2 (Binary symmetric channel (BSC)).

$$A_X = A_Y = \{0, 1\}, \quad q(y|x) = \begin{cases} 1 - \epsilon & y = x, \quad 0 < \epsilon < 0.5 \\ \epsilon & y \neq x. \end{cases}$$

Transition diagram:

Additive model: $Y = X \oplus V$, X, V independent, and $p_V(1) = \varepsilon, p_V(0) = 1 - \varepsilon$.

Example 4.2.3 (Additive Gaussian channel).

$$A_X = A_Y = \mathbb{R}, q(y|x) = \frac{1}{\sqrt{2\pi}\sigma_V} e^{-\frac{(y-x)^2}{\sigma_V^2}} \iff Y = X + V, V \sim \eta(0, \sigma_V^2).$$

Stationary memoryless,

Definition 4.2.1 (Memoryless). Given X_i , a memoryless channel Y_i is independent of X_j 's and Y_j 's.

If the input X_1, \dots, X_n , the probability distribution of output is

$$p(y^N|x^N) = q(y_1|x_1)q(y_2|x_2) \dots q(y_N|x_N).$$

From now on we always assume this condition.

Performance: rate: # z bits per channel symbol. large is desired. accuracy:

$$p_{LE} = \lim_{K \rightarrow \infty} \frac{1}{K} \sum_{i=1}^K \Pr(\hat{Z}_k \neq Z_n)$$

Question: what is an achievable value for the (rate, accuracy) pair?