

Math 5251 Probability (Some of §1.4, 1.5)

We want to talk about the **average length** of code words when we imagine the source words

$$W = \{w_1, w_2, \dots, w_m\}$$

being emitted randomly with certain probabilities $\{p_1, p_2, \dots, p_m\}$

from a **memoryless source**, meaning the previous words don't affect the probability p_i that the next word is w_i .

Not a reasonable model for most messages locally, but not so unreasonable for long messages from a source with known word frequencies.

EXAMPLE For one of our earlier encodings

$$W = \{A, B, C, D, E\} \quad \text{with } \Sigma = \{0, 1, 2\}$$

$f \downarrow$

$$\Sigma^* \{0, 1, 20, 21, 22\} = \mathcal{C}$$

if we assume source words appear with these probabilities

$$P(A) = \frac{1}{2} p_1$$

$$P(B) = P(C) = P(D) = P(E) = \frac{1}{8} p_2, p_3, p_4, p_5$$

then what is the **average length** of a codeword?

DEFIN:

average
codeword
length

$$:= p_1 l(w_1) + \dots + p_m l(w_m)$$

$$= \frac{1}{2} \cdot \underbrace{1}_{l(0)} + \frac{1}{8} \cdot \underbrace{1}_{l(1)} + \frac{1}{8} \cdot \underbrace{2}_{l(20)} + \frac{1}{8} \cdot \underbrace{2}_{l(21)} + \frac{1}{8} \cdot \underbrace{2}_{l(22)}$$

$$= \frac{1}{2} + \frac{1}{8} + \frac{6}{8} = \frac{11}{8} = 1.375$$

This is an example of the **expected value** of a **random variable** on a **probability space**...

DEF'N: A finite probability space is a finite set $\Omega = \{\omega_1, \omega_2, \dots, \omega_m\}$ (like $W = \{\omega_1, \dots, \omega_m\}$) with probabilities $P(\omega_i) = p_i$ assigned to each ω_i

"the probability that sampling from Ω produces ω_i is p_i "

such that
$$\begin{cases} p_i \in [0, 1] \\ p_1 + p_2 + \dots + p_m = 1. \end{cases}$$

DEF'N: A **random variable** X on Ω is a function $X: \Omega \rightarrow \mathbb{R}$
 $\omega_i \mapsto X(\omega_i)$

and its **expected value**

$$\mathbb{E}X := \sum_{i=1}^m p_i X(\omega_i)$$

EXAMPLES

$$(1) \quad \Omega = W = \{A, B, C, D, E\}$$

with

$$(p_1, p_2, p_3, p_4, p_5)$$
$$= \left(\frac{1}{2}, \frac{1}{8}, \frac{1}{8}, \frac{1}{8}, \frac{1}{8}\right)$$

is a probability space, and we have
a random variable $X: \Omega \rightarrow \mathbb{R}$

codeword
length

$$w_i \mapsto l(f(w_i))$$

where f :

A	\mapsto	0
B	\mapsto	1
C	\mapsto	20
D	\mapsto	21
E	\mapsto	22

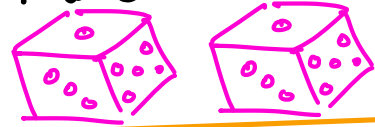
whose expected value

$$E[X] = \text{average length} = \sum_{i=1}^5 p_i l(f(w_i))$$

$$= \frac{1}{8} \cdot 1 + \frac{1}{2} \cdot 9 + \frac{1}{2} \cdot 2 + \frac{1}{2} \cdot 2 + \frac{1}{2} \cdot 2$$

$$= \frac{11}{8} = 1.375 \text{ from before}$$

(2) What is the expected value for the roll of one fair 6-sided die? Total of two fair dice?



$$\Omega_1 = \{ \overset{\omega_1}{1}, \overset{\omega_2}{2}, \overset{\omega_3}{3}, \overset{\omega_4}{4}, \overset{\omega_5}{5}, \overset{\omega_6}{6} \} = \text{outcomes for one die}$$

$$X \downarrow P(\omega_i) = \frac{1}{6} \quad \forall i \quad \leftarrow \text{called the uniform probability space on } \Omega$$

$$\mathbb{R} \quad X(i) = i \quad \text{when } P(\omega_i) = \frac{1}{|\Omega|} \forall i$$

$$EX = \frac{1}{6} \cdot 1 + \frac{1}{6} \cdot 2 + \dots + \frac{1}{6} \cdot 6$$

$$= \frac{1}{6} (1+2+3+4+5+6) = \frac{1}{6} \cdot 21 = \frac{7}{2} = 3.5$$

$$\Omega_2 = \Omega_1 \times \Omega_1 = \left\{ \begin{array}{l} (1,1), (1,2), \dots, (1,6), \\ (2,1), (2,2), \dots, (2,6), \\ \vdots \\ (6,1), (6,2), \dots, (6,6) \end{array} \right\} \quad P((i,j)) = \frac{1}{36}$$

$$= \frac{1}{|\Omega_2|}$$

↑ the uniform probability space again

$$X \downarrow \mathbb{R} \quad X(i,j) = i+j = \text{total of the dice}$$

$$EX = \frac{1}{6}(1+1) + \frac{1}{6}(1+2) + \dots + \frac{1}{6}(6+6)$$

$$= \frac{1}{36} \cdot 2 + \frac{2}{36} \cdot 3 + \frac{3}{36} \cdot 4 + \frac{4}{36} \cdot 5 + \frac{5}{36} \cdot 6 + \frac{6}{36} \cdot 7 = \frac{252}{36} = 7$$

$$\frac{1}{36} \cdot 12 + \frac{2}{36} \cdot 11 + \frac{3}{36} \cdot 10 + \frac{4}{36} \cdot 9 + \frac{5}{36} \cdot 8$$

Entropy of a sample space (§2.2)

In 1948, Claude Shannon tried to quantify how much **information** we acquire when we are told the outcome ω_i of a sampling from a probability space $\Omega = \{\omega_1, \omega_2, \dots, \omega_m\}$

having probabilities p_1, \dots, p_m so $P(\omega_i) = p_i$. We will eventually call this the **Shannon entropy**

$$H(\Omega) = H(p_1, p_2, \dots, p_m) \text{ of } \Omega$$

The idea is to first **define/normalize** by saying the **self-information** $I(\omega_i)$ for an outcome heads/tails of a fair coin flip

$$\Omega = \{\text{heads}, \text{tails}\}$$

$$\begin{array}{cc} \omega_1 & \omega_2 \\ P(\text{heads}) = \frac{1}{2} & P(\text{tails}) = \frac{1}{2} \\ = p_1 & = p_2 \end{array}$$

is $I(\text{heads}) := I(\text{tails}) := 1 \text{ bit}$

Then if one did **2 coin flips**, each outcome would have $P(\omega_i) = \frac{1}{2} \cdot \frac{1}{2} = \frac{1}{4} = P(\text{heads, heads}) = P(\text{heads, tails})$ and should have **twice the self-information**, that is, $I(\omega_i) = 2$ bits.

Similarly k coin flips have outcomes ω_i with all $P(\omega_i) = \underbrace{\frac{1}{2} \cdot \frac{1}{2} \cdots \frac{1}{2}}_{k \text{ times}} = \frac{1}{2^k}$

and should have $I(\omega_i) = k$ bits
 $= -\log_2\left(\frac{1}{2^k}\right) = -\log_2(p_i)$

This motivates the choice that ...

DEFIN: An outcome ω_i in Ω having $\Pr(\omega_i) = p_i$ has **self-information** $I(\omega_i) := -\log_2(p_i)$ and the (Shannon) **entropy/information** for Ω is the expected value EI of the self-information:

$$\begin{aligned} H(\Omega) &:= H(p_1, p_2, \dots, p_m) \\ &= -p_1 \log_2(p_1) - \dots - p_m \log_2(p_m) \\ &= -\sum_{i=1}^m p_i \log_2(p_i) \quad \text{in bits.} \end{aligned}$$

EXAMPLES

$$(1) \quad \Omega = \{A, B, C, D, E\} \text{ with} \\ (p_1, p_2, p_3, p_4, p_5) \\ = \left(\frac{1}{2}, \frac{1}{8}, \frac{1}{8}, \frac{1}{8}, \frac{1}{8}\right)$$

$$\text{has } H(\Omega) = \frac{1}{2} \cdot 1 + \frac{1}{8} \cdot 3 + \frac{1}{8} \cdot 3 + \frac{1}{8} \cdot 3 + \frac{1}{8} \cdot 3 \\ = \frac{1}{2} + \frac{12}{8} = \frac{1}{2} + \frac{3}{2} = 2$$

$$(2) \quad \Omega = \{\omega_1, \dots, \omega_m\} \text{ with uniform distribution} \\ (p_1, \dots, p_m) = \left(\frac{1}{m}, \dots, \frac{1}{m}\right)$$

$$\text{has } I(\omega_i) = -\log_2\left(\frac{1}{m}\right) = \log_2(m) \quad \forall i$$

$$\text{and } H(\Omega) = H\left(\frac{1}{m}, \dots, \frac{1}{m}\right) = \frac{1}{m} \log_2(m) + \dots + \frac{1}{m} \log_2(m) \\ = \log_2(m)$$

e.g.

$$H\left(\frac{1}{2}, \frac{1}{2}\right) < H\left(\frac{1}{3}, \frac{1}{3}, \frac{1}{3}\right) < H\left(\frac{1}{4}, \frac{1}{4}, \frac{1}{4}, \frac{1}{4}\right) < \dots \\ = \log_2(2) \quad = \log_2(3) \text{ bits} \quad = \log_2(4) \\ = 1 \text{ bit} \quad \quad \quad = 2 \text{ bits}$$

Q: Why is $H(\Omega)$ always nonnegative?

How well does $H(\Omega)$ capture the notion of the information conveyed by knowing the outcome ω from a sampling of Ω ? A supporting result...

THEOREM (Roman THM 1.1.1) Any function $H(p_1, \dots, p_m)$ defined for all sequences (p_1, \dots, p_m) with $p_i \in [0, 1]$ $\sum_{i=1}^m p_i = 1$

having these properties

- (i) H is **continuous** as a function of the p_i
- (ii) $H(\frac{1}{n}, \dots, \frac{1}{n}) < H(\frac{1}{n+1}, \dots, \frac{1}{n+1})$ for all $n=1, 2, \dots$

(iii) $H(p_1, \dots, p_r, q_1, \dots, q_s) =$
 $H(p, q) + p H(\frac{p_1}{p}, \dots, \frac{p_r}{p}) + q H(\frac{q_1}{q}, \dots, \frac{q_s}{q})$
let $p = p_1 + \dots + p_r$ let $q = q_1 + \dots + q_s$

must be of the form

$$H(p_1, \dots, p_m) = - \sum_{i=1}^m p_i \log_b(p_i)$$

for some choice of **base $b > 1$**

CONVENTION:
 $p \log(p) = 0$
 if $p=0$

NOTE:

This pins $H(\Omega)$ down up to a multiple:
 (Shannon picked $b=2$)

$$\log_b(p) = \frac{\log_2(p)}{\log_2(b)}$$