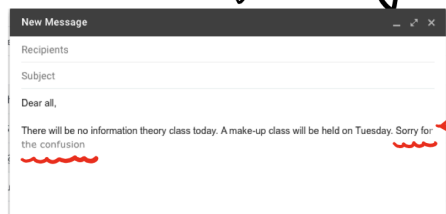


## 2: Uncertainty, compression, and entropy ①

### [A] Source Model

With probability as our means for modeling uncertainty, we can now return to our goal of answering "how much information is revealed when a question is answered."

Recall the following example:



for the confusion

for the inconvenience

for the confusion

Gmail server can simply assume that this unknown phrase was generated by a probability distribution.

We call the generator of this phrase (the user)

a source which is captured by a pmf  $P$  on  $\mathcal{X}$ .

- We may treat this entire phrase as only symbol

$$X \sim P \text{ on } \mathcal{X}.$$

- Or, we can treat each word as a symbol and the source generates a sequence of symbols.

\* Memoryless Source: The sequence  $\{X_t\}_{t=1}^{\infty}$  is output of

a memoryless source if it is independent & identically distributed

\* Markov source of order m: Given  $(X_{t-m}, \dots, X_{t-1})$ , the random variable  $X_t$  is independent of  $\{X_{t-j}\}_{j>m}$ .

(2)

This implies 
$$p(x_1, \dots, x_n) = \prod_{t=1}^n p(x_t | x_1, \dots, x_{t-1})$$
$$= \prod_{t=1}^n p^{(m)}(x_t | x_{t-m}, \dots, x_{t-1})$$

which is a fixed conditional prob called the transition matrix

→ Now that we know how to model uncertainty, how do we measure it?

Thesis: Uncertainty of a source equals the number of bits used to represent its output

avg. no. of bits or "likely" no. of bits

(both are roughly the same by estimates we saw in the previous lecture)

### B An Example

For our running example, suppose the gmail server goes over all my composed emails and finds that the first phrase follows similar expression  $\frac{1}{2}$  the times, the next  $\frac{1}{4}$  times, the next  $\frac{1}{8}$ ,  $\frac{1}{16}$ ,  $\frac{1}{32}$  times. Overall gmail found 256 possible phrases.

all the remaining possibilities are equally likely

$$P = \left( \frac{1}{2}, \frac{1}{4}, \frac{1}{8}, \frac{1}{16}, \frac{1}{32}, \dots \right)$$

Then, how do we store these phrases?

(3)

\* Make a look-up table and just store the entry number:

Record #	Phrase
1	phrase 1
2	phrase 2
3	phrase 3

- Uses 1 byte to represent the phrase.

\* Why just one table? Why not two?

→ First use 1 bit to represent which table and then represent the entry in the table:

$$\underbrace{\text{Table 1} \leftarrow \left\{ \frac{1}{2}, \frac{1}{4} \right\}}_{1 \text{ bit}}; \quad \underbrace{\text{Table 2} \leftarrow \left\{ \frac{1}{8}, \frac{1}{16}, \frac{1}{32}, \dots \right\}}_{8 \text{ bits}}$$

Worst-case requirement is now \_\_\_\_\_

$$\begin{aligned} \text{But the average \# of bits} &= 1 + \frac{3}{4} \cdot 1 + \frac{1}{4} \cdot 8 \\ &= \underline{\underline{3.75}} \text{ bits} \end{aligned}$$

- Maybe we can try 2 bits for the first table:

$$\text{Table 1} \leftarrow \left\{ \frac{1}{2}, \frac{1}{4}, \frac{1}{8}, \frac{1}{16} \right\}$$

$$\Rightarrow \text{avg. no. of bits} = 1 + \frac{15}{16} \cdot 2 + \frac{1}{16} \cdot 8 = \frac{39}{16} \approx \underline{\underline{2.44}} \text{ bits}$$

- Maybe 3 bits for table 1: avg. no. of bits  $\approx \underline{\underline{4.15}}$  bits

(we can convince ourselves that 2.44 bits is the best)

\* Why just two tables?

(4)

Let's make multiple tables. We should make at least 4 tables.

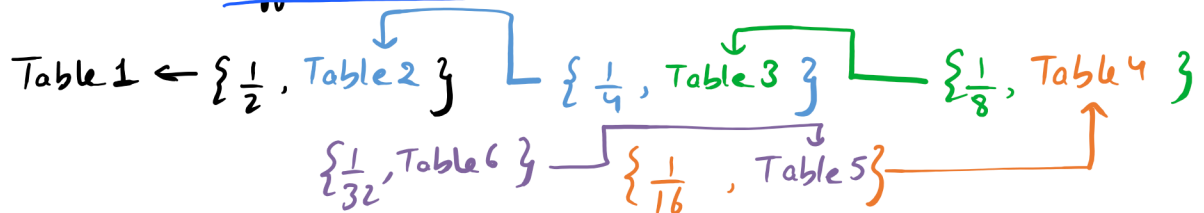
- 00 : Table 1  $\leftarrow \{1/2, 1/4\}$  1 bit
- 01 : Table 2  $\leftarrow \{1/8, 1/16\}$  1 bit  $\Rightarrow$  avg. no. of bits =  $2 + \frac{3}{4} \cdot 1 + \frac{3}{16} \cdot 1 + \frac{1}{32} \cdot 8$
- 10 : Table 3  $\leftarrow \{1/32\}$  0 bit
- 11 : Table 4  $\leftarrow \{251 \text{ symbols}\}$  8 bits  $\approx 3.25$

Easy to improve: put  $1/2$  in Table 1 alone.

$$\Rightarrow \text{avg. no. of bits} = 2 + \frac{1}{2} \cdot 0 + \frac{3}{8} \cdot 1 + \frac{3}{32} \cdot 1 + \frac{8}{32}$$

$$\approx 2.72$$

Another view: Tables are "super-symbols" with prob. equal to sum of prob. of all symbols inside it. This view leads to Huffman code



$$\Rightarrow \text{avg. no. of bits} = \frac{1}{2} \cdot 1 + \frac{1}{4} \cdot 2 + \frac{1}{8} \cdot 3 + \frac{1}{16} \cdot 4 + \frac{1}{32} \cdot 5 + \frac{1}{32} \cdot 8$$

$$\approx \underline{\underline{2.03}}$$

Yet another view:

?  $\leftarrow$  0, ?  $\leftarrow$  1, ?  $\leftarrow$  00, ?  $\leftarrow$  01, ?  $\leftarrow$  11,  $\leftarrow$  8 bit sequences remaining symbols  
 $1/2$        $1/4$        $1/8$        $1/16$        $1/32$

$$\Rightarrow \text{avg. no. of bits} = \frac{1}{2} + \frac{1}{4} + \frac{1}{8} \cdot 2 + \frac{1}{16} \cdot 2 + \frac{1}{32} \cdot 2 + \frac{1}{32} \cdot 8$$

$$\approx \underline{\underline{1.44 \text{ bits}}}$$

Where does this improvement end?

## A few lessons

(5)

- (a) The two-table example is instructive: In some applications we can treat the no. of bits used for the first table as the effective length of the final representation, if symbols in Table 2 appear with very small prob.
- (b) The final scheme we presented is not practical, but can serve as a tool for characterizing a measure of uncertainty.

### C A Basic Compression Problem

The first remark above motivates the following quantity:

$$L_{\varepsilon}(P) = \min \{ \lceil \log_2 |A| \rceil : A \subseteq \mathcal{X} \text{ s.t. } P(A) \geq 1 - \varepsilon \}$$

"fixed-length code"  $\equiv$  min no. of bits required to represent a set of probability  $\geq 1 - \varepsilon$

$\equiv$  min no. of bits required to represent Table 1

The second remark motivates the next definition:

$$\bar{L}(P) = \min \left\{ \sum_{x \in \mathcal{X}} p(x) |e(x)| : e: \mathcal{X} \rightarrow \{0, 1\}^* \text{ is 1-1} \right\}$$

"variable-length code"  $\equiv$  min avg length of a code which represents  $\mathcal{X}$  as bits

Note that both these quantities have been defined as mathematical entities, motivated by our search for a measure of uncertainty, and are devoid of any practical constraints. When we revisit compression later, we will consider practical schemes for compression.

Building on our scheme for the example we saw, we can identify the optimal schemes for both  $L_\varepsilon(P)$  and  $\bar{L}(P)$ . ⑥

→ Simply arrange the probabilities in a decreasing order and assign sequences of increasing length to the symbols.

Specifically, let  $P_{(1)} \geq P_{(2)} \geq \dots \geq P_{(k)}$  denote the probabilities of the pmf  $P$  arranged in a decreasing order, with  $(i)$  denoting a symbol with the  $i^{\text{th}}$  highest probability.

Theorem  
(a)  $L_\varepsilon(P) = \lceil \log N_\varepsilon(P) \rceil$

where

$$N_\varepsilon(P) = \min \left\{ \ell : \sum_{i=1}^{\ell} P_{(i)} \geq 1 - \varepsilon \right\}.$$

$$(b) \quad \bar{L}(P) = \sum_{i=1}^k P_{(i)} \left\lceil \log \left( \frac{i+1}{2} \right) \right\rceil$$

$= \min \left\{ j \in \mathbb{N} : \sum_{\ell=1}^j 2^{-\ell} \geq i \right\}$

Proof: (a) Clearly,  $L_\varepsilon(P) \leq \lceil \log N_\varepsilon(P) \rceil$ . For the other direction,

consider a set  $A \subset \mathcal{X}$  with  $P(A) \geq 1 - \varepsilon$ . Note that

$$1 - \varepsilon \leq \sum_{i=1}^{|A|} p_i \leq \sum_{i=1}^{|A|} P_{(i)} \Rightarrow |A| \geq N_\varepsilon(P),$$

which implies  $L_\varepsilon(P) \geq \lceil \log N_\varepsilon(P) \rceil$ .

(b) Again,  $\bar{L}(P) \leq \sum_{i=1}^k P_{(i)} \lceil \log \left( \frac{i+1}{2} \right) \rceil$  follows by using the scheme which simply assigns the shortest available binary sequence to symbols in the decreasing order of their probabilities.

For the other direction, for any other scheme that assigns different binary sequences to each symbol, let  $l(i)$  denote the length of the sequence assigned to  $p(i)$ . We show that

$$\sum_{i=1}^t p(i) l(i) \geq \sum_{i=1}^t p(i) \lceil \log(\frac{i}{2} + 1) \rceil \text{ for all } 1 \leq t \leq k.$$

Clearly, the claim holds for  $t=1$ . Suppose it holds for a  $t \in \mathbb{N}$ .

$$\begin{aligned} \text{Then, } \sum_{i=1}^{t+1} p(i) l(i) &= \sum_{i=1}^t p(i) l(i) + p_{(t+1)} l_{(t+1)} \\ &\geq \sum_{i=1}^t p(i) \lceil \log(\frac{i}{2} + 1) \rceil + p_{(t+1)} l_{(t+1)} \end{aligned}$$

[ by induction hypothesis ]

Namely, the induction hypothesis claims that it is optimal to have the first  $t$  symbols represented by  $\{0, 1, 00, 01, 11, 10, \dots\}$ . We restrict to such schemes now. But then, if we represent  $(t+1)$  with the  $(t+1)^{th}$  shortest binary sequence, we will only improve the performance. Therefore,

$$\sum_{i=1}^{t+1} p(i) l(i) \geq \sum_{i=1}^{t+1} p(i) \lceil \log(\frac{i}{2} + 1) \rceil.$$

(Alternatively, we can use induction to show that the optimal scheme with represent the symbols  $(1), \dots, (t)$  using the  $t$ -shortest binary sequences. □

Exercise: Show that

$$\frac{\bar{L}(P) - \varepsilon \lceil \log |X| \rceil}{1 - \varepsilon} \leq L_\varepsilon(P) \leq \frac{\bar{L}(P)}{\varepsilon} + 1.$$

## D) Enter Shannon Entropy

(8)

The bounds we develop in the exercise above suggest (roughly) that the difference between  $L_\varepsilon(P)$  and  $\bar{L}(P)$  depends on  $\varepsilon$ , but both are fundamentally the same quantity.

So, we will simply focus on  $L_\varepsilon(P)$  and treat it as our measure of uncertainty.

The expression we have identified for  $L_\varepsilon(P)$  doesn't look very useful. For instance, it has an  $\varepsilon$  in-built in it but we seek a fundamental notion of uncertainty that does not depend on such external parameters.

We now find an approximate upper bound that is not as accurate as the previous formula, but brings out a new feature of  $P$  that measures its uncertainty.

→ Instead of selecting the most likely symbols, we simply retain the symbols that are more likely than a threshold.

$$\begin{aligned} \text{Let } A_\lambda &= \{x: P(x) \geq 2^{-\lambda}\} \\ &= \{x: -\log p(x) \leq \lambda\} \end{aligned}$$

another interpretation: keep symbols that are not too rare, with rarity measured as  $-\log p(x)$



Exercise:  $|A_\lambda| \leq \lambda$

(9)

Theorem. If for  $\lambda > 0$  we have

$$P(\{x: -\log p(x) \leq \lambda\}) \geq 1 - \varepsilon.$$

Then,

$$L_\varepsilon(P) \leq \lambda$$

In words: A  $(1-\varepsilon)$ -prob. upper bound on the random variable

$Z = -\log p(x)$  is an upper bound for  $L_\varepsilon(P)$

And what is our estimate for a large prob. upper bound for a random variable? Markov ineq. gives an estimate of  $\mathbb{E}[Z]$  since  $P(Z \leq \frac{\mathbb{E}[Z]}{\varepsilon}) \geq 1 - \varepsilon.$

Thus, by the Theorem above,

$$L_\varepsilon(P) \leq \frac{\mathbb{E}[-\log p(x)]}{\varepsilon}$$

We have discovered our fundamental measure of uncertainty!

$$H(P) := \sum_x p(x) \log \frac{1}{p(x)} : \text{Shannon Entropy}$$

Exercise (i) Compute entropy for  $P$  given by

(a)  $(\frac{1}{1024}, \dots, \frac{1}{1024})$  (b)  $(\frac{1}{2}, \frac{1}{4}, \frac{1}{1024}, \dots, \frac{1}{1024})$  (c)  $(\frac{3}{4}, \frac{1}{1024}, \dots, \frac{1}{1024})$

(ii) Plot entropy of  $\text{Ber}(p)$  as a function of  $p$ .

## E Random hash

(10)

Up to now, we have seen that using  $P$  we can form a list of size  $\approx 2^{H(P)}$  which contains the unknown  $X \sim P$  with high probability.

In our question/answer analogy, the person knowing the answer can reveal the answer to the guesser using  $\approx H(P)$  bits, if he knows the list. We used this observation to justify that the information revealed by the answerer is  $\approx H(P)$ . This may not be satisfactory since there seems to be some hidden information (in form of the guesslist) shared between the guesser and the answerer.

We now show that the same result holds even when the answerer does not know the list. Specifically, we show:

The answerer can send  $l \approx H(P)$  bits about the answer, without knowing  $P$ , which will allow the guesser to get the answer with high probability.

→ The tool we use is a random hash.

Definition ( $l$ -bit random hash) Consider a random function  $F: \mathcal{X} \rightarrow \{0, 1\}^l$  chosen uniformly over the set of all functions

from  $\mathcal{X}$  to  $\{0, 1\}^l$ . The output  $F(x)$  is called the  $l$ -bit random hash of the input  $x$ . (11)

Lemma (Collision-resistance property) For a random function  $F$  as above and every distinct  $x, x' \in \mathcal{X}$ , we have.

$$\mathbb{P}(F(x) = F(x')) = 2^{-l}.$$

Proof. A simple exercise.

Thus, a random hash is an "almost" one-one mapping.

Suppose now that the guesser has a list of  $A$  such that  $x \in A$ . More generally, assume

$$\mathbb{P}(x \in A) \geq 1 - \epsilon.$$

The answerer can simply send  $F(x)$  and the guesser can declare the answer  $\hat{x}$  as that  $x \in \mathcal{L}$  for which  $F(x) = F(x)$ .

(If there are multiple, use any one)

Then,

$$\begin{aligned} \mathbb{P}(\hat{X} \neq X) &= \mathbb{P}(\hat{X} \neq X, X \in A) + \mathbb{P}(\hat{X} \neq X, X \in A^c) \\ &\leq \mathbb{P}(\hat{X} \neq X | X \in A) + \mathbb{P}(X \in A^c) \\ &\leq \mathbb{P}(\hat{X} \neq X | X \in A) + \epsilon \end{aligned}$$

$$\begin{aligned} \text{Also, } \mathbb{P}(\hat{X} \neq X | X \in A) &= \sum_{x \in A} \frac{p(x)}{P(A)} \cdot \sum_{x' \in A: x' \neq x} \mathbb{P}(F(x) = F(x')) \\ &\leq |A| \cdot 2^{-l}. \end{aligned}$$

Thus,

$$P(\hat{X} \neq X) \leq |A| 2^{-l} + \varepsilon$$

If we choose  $s.t.$   $|A| \leq 2^{L_\varepsilon(P)}$ , namely the small  $A$   $s.t.$   $P(A) \geq 1 - \varepsilon$ , then

$$P(\hat{X} \neq X) \leq 2^{L_\varepsilon(P) - l} + \varepsilon \leq \eta + \varepsilon$$

$$\text{if } l \geq L_\varepsilon(P) + \frac{1}{\eta}.$$

In particular, we can choose  $l \geq H(P)$  to get a small prob. of error in view of the results of the last section.