

Today : DATA COMPRESSION

SOME LIMITS TO COMPRESSIBILITY.

---

First : A Review of the last two lectures.

first few lectures : Studied properties of

- entropy, conditional entropy, mutual information, relative entropy
- One consequence : if variable  $X$  has large conditional entropy given  $Y$  then any attempt to determine  $X$  from  $Y$  has large error (Fano's lemma)

Lecture 4 : restrict to variables  $\bar{X}, \bar{Y}$

that are sequence of i.i.d. variables

distributed according to  $X$

when  $\bar{X} = (X_1, \dots, X_n)$  then

w.p.  $\geq 1 - \delta$

$$(AEP) \quad \mathbb{P}_{\bar{X}} \left[ p(\bar{X}) \in \left[ 2^{-(H(x)+\epsilon)n}, 2^{-(H(x)-\epsilon)n} \right] \right] \geq 1 - \delta$$

$$(\text{Typical Set}) \quad A_\epsilon^{(n)} = \left\{ \bar{X} \in \mathcal{X}^n \mid p(\bar{X}) \in \left[ \downarrow \right] \right\}$$

$A_\epsilon^{(n)}$  has size  $\doteq 2^{H(x) \cdot n}$

with each el't having prob.  $2^{-H(x) \cdot n}$ .

(flat disk. on small set)

[didn't say but  $\delta \approx \exp(-\epsilon^2 n)$ ]

Leads to formal proof that

-  $X$  can be compressed into  $\sim H(x) \cdot n$

bits in

expectation

-  $X$  can't be compressed to much less than

$H(x) \cdot n$  bits

## LECTURE 5: $\bar{X}$ = stochastic process

(esp. time-invariant,  
irreducible,  
aperiodic, Markov chain)

Stochastic  $\supseteq$  Stationary  $\supseteq$  Time Inv. Markov  
Chain in Stat.  
distribution

Entropy Rate of  $(\bar{X})$ :  $H(X_2|X_1)$

$$= \sum_{i,j} \mu_i P_{ij} \log \frac{1}{P_{ij}}$$

where  $P$  = transition probability matrix

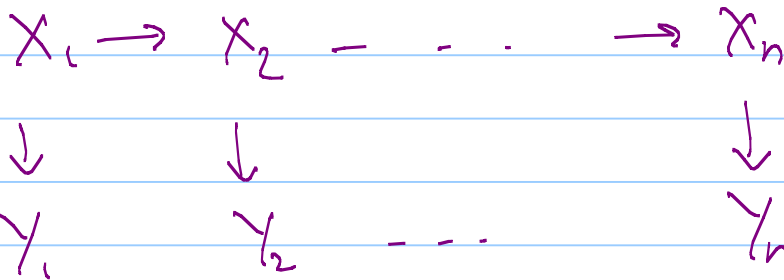
$\mu$  = stationary dist.

(AEP):

$$P_r \left[ P(\bar{X}) \in \left[ \frac{-(H(x) + \epsilon)n}{2}, \frac{-(H(x) + \epsilon)n}{2} \right] \right] \geq 1 - \delta$$
$$\delta \approx \exp(-\epsilon^2 n)$$

Finally

Hidden Markov chains



Entropy rate

$$H(\mathcal{Y}) = \lim_{n \rightarrow \infty} H(Y_n | Y_1, \dots, Y_{n-1})$$

Exists, but not easy to compute

Now to today's lecture

Goal : to "compress"  $X$  into string  
of letters over  $D$ .

but ....

- ① What is compression?
- ② What is the objective?

E.g.  $X \in \{1 \dots N\}$

$X=i$  w.p.  $p_i$

Is  $C(i) = 0 \quad \forall i$  "compression"?

Trade off between "compression" & loss of  
information ....

We wish to compress without losing  
(too much) information.

## Terminology

$\Omega$  = space of  $X$

$D$  = alphabet for compression

$D^*$  =  $\bigcup_{n \geq 0} D^n$  = set of finite strings  
over  $D$

for  $w = (w_1 \dots w_n) \in D^n$ ,  $|w| = n$

## Compression algorithm

$C: \Omega \rightarrow D^*$

accompanied by Decompressor

$Dec: D^* \rightarrow \Omega$

Lossy compression:  $\Pr_x [Dec(C(x)) \neq x] = \text{small}$

Lossless compression:  $\forall x \quad Dec(C(x)) = x$

if  $C$  st.  $\exists$  Dec s.t.

$(C, Dec)$  lead to lossless coding

then  $C$  is called Non-Singular

We will talk only about Non-Singular

(Lossless Coding) but many results

extend to Lossy Coding as well

(er... we mean limitations result of course)



Goal / Measure = ?

Expected Lengths :  $\mathbb{E}_x [ |C(x)| ]$

↑  
Minimize

x

Example  $\Omega = \left\{ \begin{array}{ll} 1 & \text{w.p. } \frac{1}{2} \\ 2 & \text{w.p. } \frac{1}{4} \\ 3 & \text{w.p. } \frac{1}{8} \\ 4 & \text{w.p. } \frac{1}{8} \end{array} \right\}$

$D = \{0, 1\}$

Code 1:

$$C_1(1) = 00$$

$$C_1(2) = 01$$

$$C_1(3) = 10$$

$$C_1(4) = 11$$

$$E[|C_1(x)|] = 2$$

Code 2:

$$C_2(1) = 0 \qquad C_2(3) = 00$$

$$C_2(2) = 1 \qquad C_2(4) = 10$$



$$E[|C_2(x)|] = \frac{3}{4} \cdot 1 + \frac{1}{4} \cdot 2 = \frac{5}{4}$$

- But something unsatisfactory about second code?

Second code is not extendible.

Real goal to get a code for

sequence  $(x_1, \dots, x_n)$  where  $x_i \sim X$   
(say i.i.d.)

Code for  $X$  suggests extension

code of  $X^n$

$$C(x_1, \dots, x_n) = C(x_1), C(x_2), \dots, C(x_n)$$

Code  $C_2$  above is non-singular

but extension is not!

$$\text{(E.g. } C_2(\uparrow 1) = C_2(3)$$

$$\text{ \& no } C_2(113) = C_2(311) )$$

Motivates another definition

---

$C$  is uniquely decodable

if  $\forall n$ ,  $n$ -extension of  $C$  is  
non-singular.

Example of uniquely-decodable code

$$C_3(1) = 0$$

1 w.p.  $\frac{1}{2}$

$$C_3(2) = 10$$

2 w.p.  $\frac{1}{4}$

$$C_3(3) = 110$$

3 "  $\frac{1}{8}$

$$C_3(4) = 111$$

4 "  $\frac{1}{8}$

$$E[\text{length}] = 1.75$$

Another example

$$C_4(1) = 0$$

$$C_4(2) = 01$$

$$C_4(3) = 011$$

$$C_4(4) = 111$$

## Proof of unique decodability of $C_2$ .

### Property of $C_2$

- $\forall x, y \in \Sigma$   $C(x)$  is not a prefix of  $C(y)$

Motivates another definition

$C$  is a prefix code if  $\forall$   
 $x, y \in \Sigma$   $C(x)$  not prefix of  $C(y)$

Claim: Prefix code is uniquely decodable

Proof: Let  $C(x_1 \dots x_n) = w_1 \dots w_m$   
 $= C(y_1 \dots y_\ell)$

then  $x_1 = y_1$  (since  $C(x_i)$  must

be a prefix of  $C(y_i)$  or vice versa).

But now by induction we have

$$C(x_2 \dots x_n) = w_{k+1} \dots w_m = C(y_2 \dots y_l)$$

$$\text{Where } C(x_1) = w_1 \dots w_k = C(y_1)$$

$$\Delta \text{ as } \underbrace{x_2 \dots x_n = y_2 \dots y_l}_{\gamma} \quad [\Delta n=l].$$

Proof of unique decodability of  $C_4$

(well  $C_4 = C_3$  in reverse)

Prefix codes ( $C_3$ ) are nicer than

non-prefix codes ( $C_4$ ) since they can be decoded online.

Hence also called instantaneous

In the rest of the lecture we will study

Limitations on compressibility

- Kraft lower bound

- McMillan lower bound

- Entropy lower bound

Next lecture will match goals

Main target for this lecture

$$\bullet \quad \mathbb{E}_x [ |C(x)| ] \geq \frac{H(x)}{\log D}$$

for any prefix-free or uniquely decodable code  $C$ .

• Actually yields  
if

$$E_x [ |C(x)| ] = K$$

then  $K + O(\log K) \geq \frac{H(x)}{\log D}$

for any non-singular code  $\left[ \begin{array}{l} \text{won't show this} \\ \text{but will show} \\ K + 2\sqrt{K+2} \geq \dots \end{array} \right]$

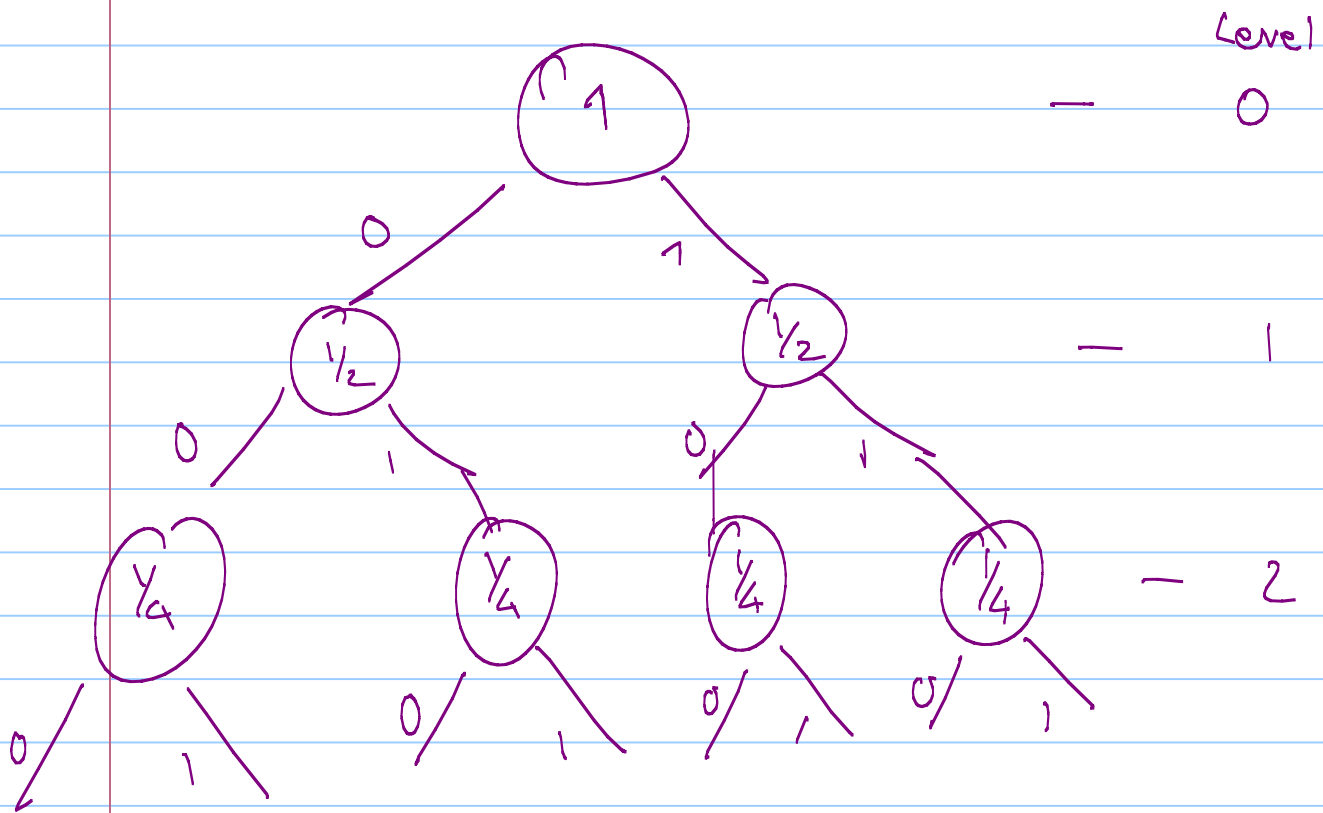
Main tool towards Entropy bound

↳ Kraft's Inequality

Theorem (Kraft): if  $C(i)$  has length  $l_i$   
for  $i = 1 \dots N$  over  $D$ -ary alphabet

then  $\sum D^{-l_i} \leq 1$ , if  $C$  prefix code.

Proof : Consider  $D$ -ary tree with branches denoting various letters of  $D$ . Example with  $D=2$



for node at level  $i$  associate weight

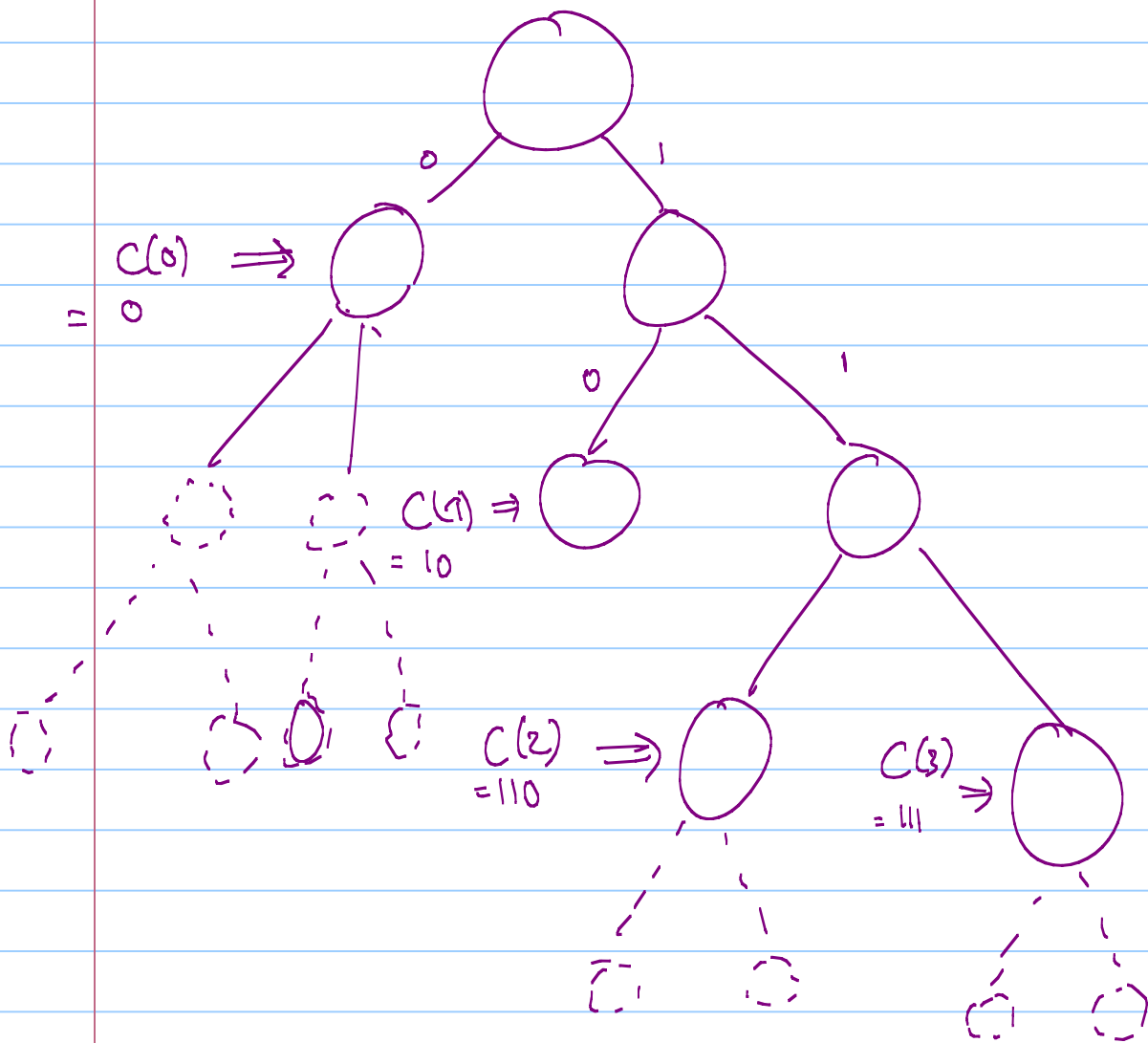
$D^{-i}$ . Note that the weight of the root =  $1$ ;

weight of every node = Sum of wts. of children

Now lets understand " $C(i)$ "'s.

for every  $i$  mark node  $C(i)$ .

Example





Retain only the part of the tree that lies on the path from some  $C(i)$  to the root.

Reassign weights, leaving weights of  $C(i)$  as they were & putting in for every node  $\text{weight} = \text{sum of weight of its children}$ .

- On the one hand weights of nodes don't go up

- On the other hand weight of root

$$= \sum_{i=1}^N D^{-l_i}$$

thus  $\sum_{i=1}^N D^{-l_i} \leq 1$

## McMillan's Bound

Theorem: if  $C$  is uniquely decodable

&  $|C(i)| = l_i$  then

$$\sum 2^{-l_i} \leq 1$$

Proof: (Optional. See Cover & Thomas)

## Entropy lower bound

Note: Kraft says nothing about probabilities

To relate to Expected decoding lengths...

$$E_x [ |C(x)| ] = \sum p_i l_i$$

But now lets write  $l_i = \frac{-\log D^{-l_i}}{\log D}$

$$E_x [ |C(x)| ] = - \sum_{i=1}^N p_i \frac{\log D^{-l_i}}{\log D}$$

But  $\sum D^{-l_i} \leq 1$  [Kraft]

$$\text{Let } q_i = D^{-l_i}$$

$$\text{so let } \sum D^{-l_i} = 1 - q_0$$

$$E_x [ |C(x)| ] = - \sum_{i=1}^N \frac{p_i \log q_i}{\log D}$$

$$= \frac{H(x)}{\log D} + \frac{D(p \parallel q)}{\log D} \geq H(x)$$

## Conclude

Theorem: if  $C$  is prefix free or uniquely decodable then for  $x \sim p$

$$E_x [ |C(x)| ] \geq H(x)$$

—————  $x$  —————

Addendum What about non uniquely-decodable codes?

Answer 1: Not a fair question. Really need unique decodability.

Answer 2: Doesn't make much difference.

§

Lemma: if  $C$  is non-singular code of <sup>exp.</sup> length  $K$ , then  $\exists C'$  prefix free of expected length  $K + 2\lceil\sqrt{K}\rceil$ .

(Can probably do better....)

Proof:

Given  $C$  first produce  $C_1$

s.t.  $\forall i \ |C_1(i)|$  is a multiple of  $\lceil\sqrt{K}\rceil$ . We have

$$\textcircled{*} \quad \mathbb{E}_x [ |C_1(x)| ] \leq \mathbb{E} [ |C(x)| ] + \lceil\sqrt{K}\rceil$$

(since no string extends by more than  $\lceil\sqrt{K}\rceil$ ).

Now produce  $C_2$  where  $|C_2(i)|$  is a multiple of  $\lceil\sqrt{K}\rceil + 1$

as follows.

$$\text{Suppose } C_1(i) = \boxed{w_1} \boxed{w_2} \dots \boxed{w_c}$$

$$\text{where } |w_j| = \sqrt{K}$$

$$\text{then } C_2(i) = \boxed{w_1} \boxed{0} \boxed{w_2} \boxed{0} \dots \boxed{w_c} \boxed{1}$$

↑      ↗      ↗      ↑  
all zeroes      ONE

Claim -  $C_2$  is prefix free (why?)

$$\geq E[|C_2(x)|]$$

$$\leq \frac{[\sqrt{K}] + 1}{\sqrt{K}} \cdot E[|C_1(x)|]$$

$$= (\sqrt{K} + 1)(\sqrt{K} + 1) \leq K + O(\sqrt{K}).$$

Conclude

Essentially for any reasonable  
encoding

$$E[|C(x)|] \geq \underbrace{H(x)}_{(\log D)}$$

But is this tight?

Will see in next lecture.