

CO781 / QIC890 Fall 2016 Lec 6.

Topic 2: Entropy & data compression.

Reference: NC Sec 12.2, Prskill Sec 10.1.1, 10.3, 10.4. Cover & Thomas

Goal: part 1 of quantifying info & redundancy.

- X : random variable
 - Ω : sample space, $|\Omega| = m = \# \text{ outcomes of } X$
 - p : prob distribution of X
- $p: \Omega \rightarrow [0, 1]$
 $x \mapsto p(x)$
- capital = rv
lowercase = outcomes

$$\text{s.t. } \sum_{x \in \Omega} p(x) = 1.$$

eg. biased coin

$$\Omega = \{0, 1\}, \quad p(0) = 0.1, \quad p(1) = 0.9.$$

- A discrete "information source" is a sequence of r.v.'s X_1, X_2, X_3, \dots with a common sample space Ω (source alphabets).
(n draws have m^n outcomes)

eg. can toss the above biased coin as many times as wished

eg. weather every day, $\Omega = \{\text{sun, cloud, rain, snow}\}$

In general X_i need not be independent, or identically distributed.

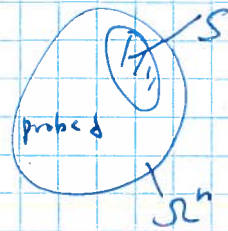
- If X_i 's are independent and identically distributed, we call X_1, X_2, \dots an "iid source".

(*)

Typicality & the asymptotic equipartition theorem:

Idea: Consider $X^n = X_1 X_2 \dots X_n$.

For large n , \exists large probability subset S
in Ω^n with low cardinality.



Furthermore, elements of S have similar probabilities.

Why?

Consider an arbitrary $x^n = x_1 x_2 \dots x_n$

$$p(x^n) = p(x_1) p(x_2) \dots p(x_n) \quad (\text{by independence})$$

$$= 2^{\log p(x_1)} 2^{\log p(x_2)} \dots 2^{\log p(x_n)} \quad (\log \text{ base } 2)$$

$$= 2^{-n \left[\frac{1}{n} \sum_{i=1}^n (-\log p(x_i)) \right]}$$

Empirical average of $-\log p(x)$

\downarrow LLN

$$2^{-n \left[\sum_{x \in \mathcal{X}} p(x) (-\log p(x)) \right]} \quad \left\langle \mathbb{E}_{\mathcal{P}} (-\log p(x)) =: H(X) \right.$$

So as $n \rightarrow \infty$, $p(x^n) \rightarrow 2^{-n H(X)}$. "typical"

If $p(x^n) \approx 2^{-n H(X)}$, then $x^n \in S$.

Def: [Shannon Entropy]

$$H(x) \text{ or } H(p) := - \sum_{x \in \Omega} p(x) \log p(x).$$

eg. for the biased coin,

$$H(x) = -0.1 \log 0.1 - 0.9 \log 0.9 = 0.469$$

Def: [typical sequence]

$$x^n \text{ is } \delta\text{-typical if } \left| \frac{1}{n} \log p(x^n) - H(x) \right| \leq \delta \quad \left(p(x^n) \approx 2^{-nH(x)} \right)$$

Def: [typical set]

$$T_{n,\delta} = \{x^n : x^n \text{ is } \delta\text{-typical}\}.$$

eg. For the biased coin, $n=100$, $\delta=0.1$

If there are t 0's & $n-t$ 1's,

$$\text{then } \frac{1}{n} \log p(x^n) = -\frac{t}{n} \log 0.1 - \frac{n-t}{n} \log 0.9$$

$$\text{and } -\frac{1}{n} \log p(x^n) \in [0.369, 0.569] \text{ for } 7 \leq t \leq 13.$$

So $T_{100,0.1} =$ all 100-bit strings with 7 to 13 0's.

Idea: $T_{n,\delta}$ is a large prob set with low cardinality,
because $-\frac{1}{n} \log p(x^n) \rightarrow H(x)$ as $n \rightarrow \infty$.

$$\text{eg. Prob}(T_{100,0.1}) = 0.75897$$

$$\left. \begin{array}{l} |T_{100,0.1}| = 8.3 \times 10^{15} \\ |\Omega^{100}| = 1.3 \times 10^{30} \end{array} \right\} \frac{|T_{100,0.1}|}{|\Omega^{100}|} \ll 1.$$

Asymptotic Equipartition Thm (AEP):

$\forall \epsilon > 0, \forall \delta > 0, \exists n_0$ s.t. $\forall n \geq n_0$

① $p(\text{Typ. S}) \geq 1 - \epsilon$

② $(1 - \epsilon) 2^{n(H(x) - \delta)} \leq |\text{Typ. S}| \leq 2^{n(H(x) + \delta)}$

③ $\forall A \subseteq \Omega^n, p(A) \geq 1 - \epsilon \Rightarrow |A| \geq (1 - 2\epsilon) 2^{n(H(x) - \delta)}$

NB: ① = large prob, ② small size, and how small ③ large prob set can't be smaller

Bonus: Typ. S = equiprobable by definition

See Preskill for full motivating example for the biased coin.

Pf ① We apply LLN on $Y = \log p(X)$, an induced r.v. s.t. $\forall x \in \Omega, Y = \log p(x)$ w.p. $p(x)$.

$$\text{So } \mathbb{E}Y = \sum_x p(x) \log p(x) = -H(X).$$

Since X^n is i.i.d., so is $Y^n = Y_1 \dots Y_n$.

Let $x^n = x_1, x_2, \dots, x_n \in \Omega^n, y_i = \log p(x_i)$

$$\text{Then } x^n \notin T_{n,\delta} \Leftrightarrow \left| \frac{1}{n} \sum_{i=1}^n y_i - \mathbb{E}Y \right| > \delta \quad (*)$$

To bound the prob of $(*)$, recall

Chebyshev's ineq for a r.v. Z :

$$\Pr \left\{ \underbrace{\left| Z - \mathbb{E}Z \right|}_{\text{outcome r.v.}} \geq K \sqrt{\text{Var}Z} \right\} \leq \frac{1}{K^2}.$$

$$\text{Choose } Z = \frac{1}{n} \sum_{i=1}^n Y_i.$$

$$\bullet n_0 = \frac{\text{Var}Y}{\delta^2 \epsilon} \quad \text{so if } n \geq n_0, \text{Var}Z = \frac{1}{n} \text{Var}Y \leq \frac{1}{n_0} \text{Var}Y = \delta^2 \epsilon$$

$$\bullet K = \frac{\delta}{\sqrt{\text{Var}Z}} \quad \text{so } K \sqrt{\text{Var}Z} = \delta.$$

So Chebyshev's ineq gives:

$$\Pr \left\{ \underbrace{\left| \frac{1}{n} \sum_{i=1}^n y_i - \mathbb{E}Y \right|}_{\text{outcome r.v.}} \geq \delta \right\} \leq \frac{\text{Var}Z}{\delta^2} \leq \epsilon \quad (**)$$

$$\therefore \text{Prob}(x \notin T_{n,\delta}) \stackrel{(*)}{=} \Pr \left\{ \left| \frac{1}{n} \sum_{i=1}^n y_i - \mathbb{E}Y \right| > \delta \right\} \stackrel{(**)}{\leq} \epsilon$$

dropping "=" decreases P

$$\therefore \Pr(T_{n,\delta}) > 1 - \epsilon.$$

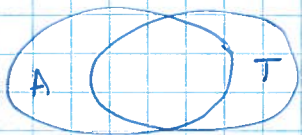
Pf(2) $1 - \epsilon \stackrel{\textcircled{1}}{\leq} p(T_{n,\delta}) \leq 1$

$$|T_{n,\delta}| \cdot \frac{1}{\min p(x^n)} \cdot 2^{-n(H(x)+\delta)} \leq \sum_{x^n \in T_{n,\delta}} p(x^n) \leq |T_{n,\delta}| \cdot \frac{1}{\max p(x^n) \text{ for } x^n \in T_{n,\delta}} \cdot 2^{-n(H(x)-\delta)}$$

$$\therefore |T_{n,\delta}| \leq 2^{n(H(x)+\delta)} \quad \& \quad (1-\epsilon) \geq 2^{-n(H(x)-\delta)} \leq |T_{n,\delta}|$$

NB: $\frac{|T_{n,\delta}|}{|\Omega^n|} \leq 2^{-n[\log_2 |\Omega| - H(x) - \delta]}$ \leftarrow exp decaying in n .
+ve for most x .

Pf(3) Let $A \subseteq \Omega^n$ have $p(A) > 1 - \epsilon$, $T = T_{n,\delta}$.



$$\begin{aligned} P(A \cap T) &= P(A) - P(A \setminus T) \\ &\geq P(A) - P(\Omega \setminus T) \\ &\geq 1 - \epsilon - \epsilon \\ &= 1 - 2\epsilon. \end{aligned}$$



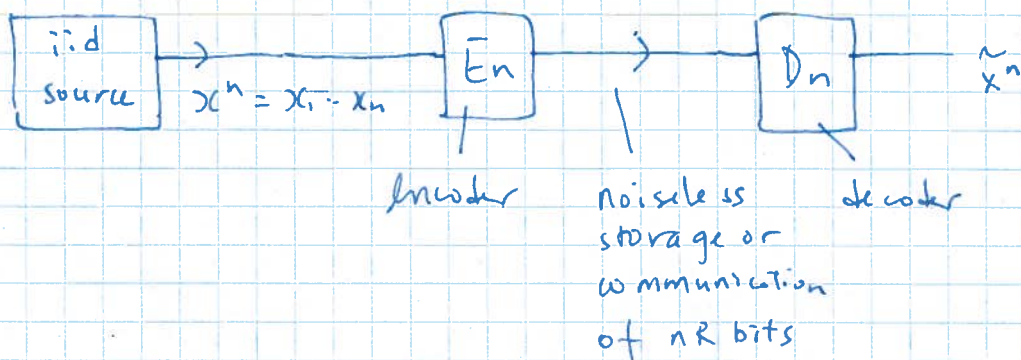
$$|A| \geq |A \cap T| \geq \frac{P(A \cap T)}{\max_{a \in A \cap T} p(a)} \geq \frac{1 - 2\epsilon}{2^{-n(H(x)-\delta)}}$$

use typicality

Application: data compression of iid sources.

Aka: Shannon's noiseless coding theorem.

Setting:



Goal: min R while keeping prob $(\hat{x}^n \neq x^n)$ negligible.

Idea: transmit only typical sequences & ignore the rest

Simplest "code book" or E_n :

Assign to each $x^n \in T_{n,\delta}$ a unique label of $n(H(X) + \delta)$ bits.

$$E_n = \begin{cases} x^n \mapsto b(x^n) & \text{if } x^n \in T_{n,\delta} \\ x^n \mapsto \text{ERR} & \text{otherwise} \end{cases}$$

The decoder D_n simply invert b if the received message $\neq \text{ERR}$

b = bijection agreed upon by Alice & Bob.

$$\Pr(\hat{x}^n \neq x^n) = \Pr(x^n \notin T_{n,\delta}) \leq \epsilon \quad \forall n \geq n_0 = \frac{\text{Var}[\log p(x)]}{\delta^2 \epsilon}$$

Shannon's noiseless coding theorem:

Let X_1, X_2, \dots be an iid source.

$$\textcircled{1} \quad \forall \epsilon > 0 \quad \forall R > H(X)$$

$$\exists n_0 \text{ s.t. } \forall n \geq n_0 \quad \exists E_n, D_n$$

$$\text{s.t. } \Pr(D_n \circ E_n(x^n) \neq x^n) < \epsilon.$$

$$\textcircled{2} \quad \forall R < H(X)$$

$$\exists n_0 \text{ s.t. } \forall n \geq n_0 \quad \forall E_n, D_n$$

$$\Pr(D_n \circ E_n(x^n) = x^n) \leq \epsilon + 2^{-n \left[\frac{H(X) - R}{2} \right]}$$

Summary: compression rate R

$$= \begin{cases} H(X) + o(1) & \text{achievable} \\ H(X) - o(1) & \text{not} \end{cases}$$

Pf $\textcircled{1}$: The E_n, D_n described in "Idea" work.

Pf $\textcircled{2}$: WLOG, E_n is deterministic.

(Otherwise, E_n succeeds w/ high prob over the random coins and \exists coin values r for which E_n succeeds whp.)

Replace E_n by its restriction to r .)

So 2^{nR} symbols are encoded, the rest turned to ZRR.

Let $A =$ these 2^{nR} symbols, $\delta = (H(X) - R)/2 > 0$, $T = T_{n,\delta}$

$$\Pr(A) = \Pr(A|T) + \Pr(A|\bar{T})$$

$$\leq \epsilon + |A| \max_{x^n \notin T} P(x^n)$$

$$\leq \epsilon + 2^{nR} \cdot 2^{-n(H(X) - \delta)}$$

$$= \epsilon + 2^{-n[(H(X) - R) - \delta]}$$

$$= \epsilon + 2^{-n \left[\frac{H(X) - R}{2} \right]}$$

Comments:

- Allowing an arbitrarily small error ϵ reduces the compression rate from $n \log |\Sigma|$ to $n H(x)$ bits.
- Note w.p. $1 - \epsilon$, the ENTIRE 2^n correct.
- Data compression gives $H(x)$ an operational meaning.
 - How much space needed to rep each symbol ASYMPTOTICALLY.
 - How uncertain a symbol is.

↑
larger limit
- We considered "block codes" where n is fixed.
We were not concerned about complexity of E_n, D_n .
See Cover & Thomas for other codes.
Ex: Huffman code is exactly but variable length
with brute force E_n / D_n . Same EXPECTED rate.
- There are universal compression algorithm that requires
only an upper bound on $H(x)$ without knowing what x is!
(Generalizations use f in Ω setting.)