# Part 1: Theory
# L04: Compression (Algorithmic Derivation of Entropy via Compression)

Javed Aslam, Wolfgang Gatterbauer

cs7840 Foundations and Applications of Information Theory (fa24)

https://northeastern-datalab.github.io/cs7840/fa24/

9/16/2024

## Last time

- Expectation
- Variance
- Markov Chains

---

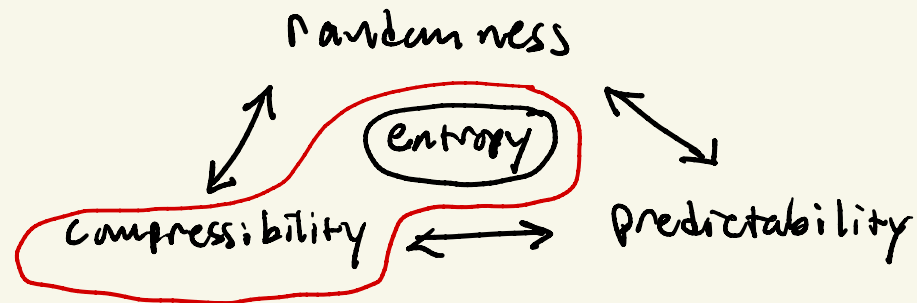- Intuitive Derivation of Entropy

Hartley $\implies$ Shannon

$$H(x) = H(\vec{p}) = \sum_i p_i \lg \frac{1}{p_i}$$

## Today

- Algorithmic Derivation of Entropy via Compression

## Next time

- Fundamental concepts in Information Theory



randomness

entropy

compressibility ⟷ predictability

# Today: Motivate Entropy via Compression

- Consider codes with codewords of length $l_1, l_2, l_3, \ldots$

- Kraft's Inequality: $\sum_i 2^{-l_i} \leq 1$   or generally   $\sum_i D^{-l_i} \leq 1$

   (binary codes)                              (D-ary codes)

Instantaneous (prefix-free) code  ①$\Longleftrightarrow$  Kraft's Inequality  ③$\Longleftarrow$  uniquely decodable codes

②$\Downarrow$

$$l_i^* = \lg \frac{1}{p_i}$$

Note:

$$\lg \triangleq \log_2$$

$$E[L] = \sum_i p_i \cdot l_i \geq \sum_i p_i \cdot l_i^* = \sum_i p_i \lg \frac{1}{p_i} = H(X)$$

Compression setup:

- A source code $C$ for r.v. $X$ is a mapping from $\mathcal{X}$, the range of $X$, to $\mathcal{D}^*$, the set of strings over encoding alphabet $\mathcal{D}$.

- The expected length $L(C) = \sum_{x \in \mathcal{X}} p(x) \cdot \ell(x)$

- A code is <u>non-singular</u> if every element of $\mathcal{X}$ maps to a <u>unique</u> string in $\mathcal{D}^*$, i.e.,
$$x_i \neq x_j \implies C(x_i) \neq C(x_j)$$

- The <u>extension</u> $C^*$ of code $C$ is a mapping from finite length strings from $\mathcal{X}$ to finite length strings from $\mathcal{D}$.

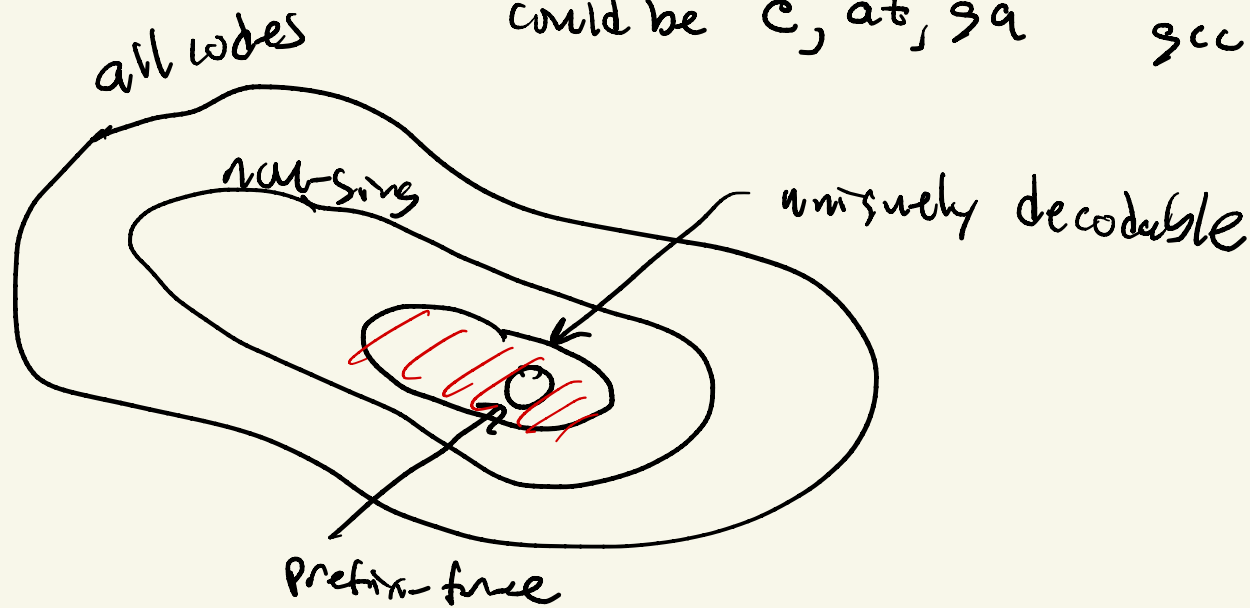- A code is <u>uniquely</u> <u>decodable</u> if its extension is <u>non-singular</u>
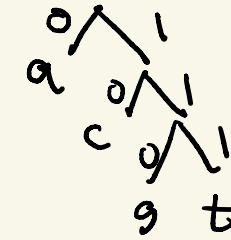
# Classes of codes

| X | Singular | Non-singular but not uniquely decodable | uniquely decodable but not instantaneous | instantaneous (prefix-free) |
|---|---|---|---|---|
| a | 0 | 0 | 10 | 0 |
| c | 0 | 010 | 00 | 10 |
| g | 0 | 01 | 11 | 110 |
| t | 0 | 10 | 110 | 111 |

↓

e.g. 010

could be c, at, ga

↓

110000

gcc



all codes

non-sing

uniquely decodable

Prefix-free

Claim: Instantaneous (prefix-free) code $\iff$ Kraft's Inequality

Pf ($\Rightarrow$): Let $l_{max}$ be longest code word

0
10
110
011
111



Prefix-free property:
internal node code
makes unavailable all
possible codes in
subtree below.

- Code of length $l_i$ wipes out how many leaves? $2^{l_{max}-l_i}$
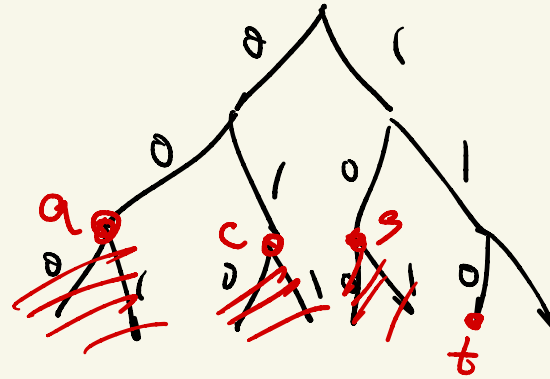
- tree only has $2^{l_{max}}$ leaves

$$\sum_i 2^{l_{max}-l_i} \leq 2^{l_{max}}$$

$\rightarrow$ dividing both sides by $2^{l_{max}}$

$\Rightarrow \sum_i 2^{-l_i} \leq 1$

$(\Leftarrow)$

| | | $l_i$ |
|---|---|---|
| a | 10 | 2 |
| c | 00 | 2 |
| g | 11 | 2 |
| t | 110 | 3 |



- Sort by length

$$l_1 \leq l_2 \leq \ldots \leq l_n$$

- assign source symbol associated w/ $l_1$ to first code lexicographically available of length $l_1$

- remove all children as possible codes

- repeat for $l_2, l_3, \ldots, l_n$

| a | 00 |
|---|---|
| c | 01 |
| g | 10 |
| t | 110 |

Proof sketch:

- subtrees assigned contiguously left-to-right
- If code lengths satisfy Kraft's inequality, never run out of subtrees

Setup: Find $l_i$ where min $L(C) = \sum_i p_i \cdot l_i$

s.t. $\sum_i 2^{-l_i} \le 1$

① $f(x) = x^2$     $\min_x f(x)$ ?     <span style="color:red">univariate optimization</span>

$$\frac{df}{dx} = 2x = 0 \implies x = 0 \quad f(0) = 0 \checkmark$$

② $f(x,y) = x^2 + y^2$     $\min_{x,y} f(x,y)$ ?     <span style="color:red">multivariate optimization</span>
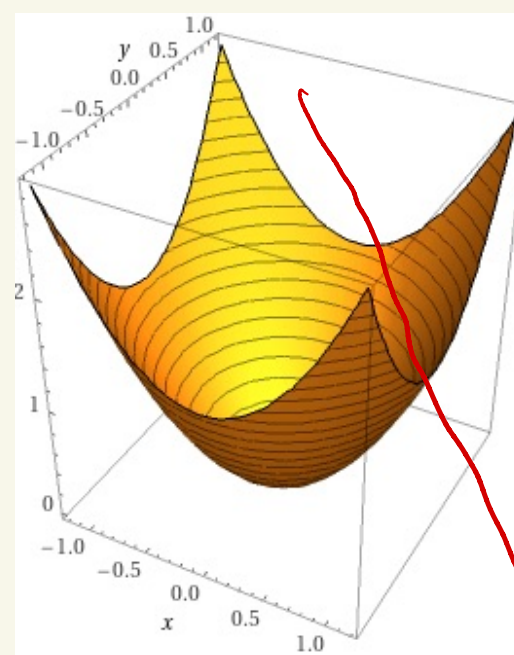
$$\frac{\partial f}{\partial x} = 2x = 0 \quad x = 0$$

$$\frac{\partial f}{\partial y} = 2y = 0 \quad y = 0$$

$$f(0,0) = 0$$

$$\min \quad f(x,y) = x^2 + y^2$$

Constrained optimization via Lagrange multipliers

$$\text{s.t} \quad \boxed{x+y=1}$$

$$J(x,y,\lambda) = x^2 + y^2 + \lambda(x+y-1)$$

$$\frac{\partial J}{\partial x} = 2x + \lambda = 0$$

$$\frac{\partial J}{\partial y} = 2y + \lambda = 0$$

subtract

$$2x - 2y = 0$$
$$x - y = 0$$

$$\frac{\partial J}{\partial \lambda} = x+y-1 = 0 \longrightarrow x+y=1$$

add

$$2x = 1$$
$$x = \tfrac{1}{2}$$
$$y = \tfrac{1}{2}$$

Find $l_i$ where $\min L(C) = \sum_i p_i \cdot l_i$

s.t. $\sum_j 2^{-l_j} \leq 1$

$$J(\vec{l}, \lambda) = \sum_i p_i l_i + \lambda \left( \sum_j 2^{-l_j} - 1 \right)$$

$2^{-l_i} = e^{-l_i \cdot \ln 2}$

$$\forall_i \frac{\partial J}{\partial l_i} = p_i + \lambda \cdot 2^{-l_i} \cdot (-\ln 2) = 0$$

$$\frac{\partial J}{\partial \lambda} = \sum_j 2^{-l_j} - 1 = 0 \Rightarrow \sum_j 2^{-l_j} = 1$$

$$\sum_i \left( p_i + \lambda \cdot 2^{-l_i} (-\ln 2) \right) = 0$$

$$\Rightarrow \sum_i p_i - (\ln 2) \cdot \lambda \sum_i 2^{-l_i} = 0$$

$$\sum_i p_i - (\ln 2) \cdot \lambda \cdot 1 = 0$$

$$1 - (\ln 2) \cdot \lambda = 0 \Rightarrow \boxed{\lambda = \frac{1}{\ln 2}}$$

$p_i + \frac{1}{\ln 2} 2^{-l_i} (-\ln 2) = 0$

$p_i - 2^{-l_i} = 0$

$2^{-l_i} = p_i$

$$\Rightarrow l_i = \lg \frac{1}{p_i}$$

So what is $\min L(c) = \sum_i p_i l_i$ ?

$$\sum_i p_i \cdot l_i^* = \sum_i p_i \cdot \lg \frac{1}{p_i} = H(x)$$

# Part 1: Theory
# L06: Compression
# (uniquely decodable codes)

Javed Aslam, Wolfgang Gatterbauer

cs7840 Foundations and Applications of Information Theory (fa24)

https://northeastern-datalab.github.io/cs7840/fa24/

9/23/2024

Last time
- Basic results in Information Theory

Today
- Finish basic results
- Continue Compression

Next time
- Continue Compression

Instantaneous (prefix-free) code $\overset{①}{\Longleftrightarrow}$ Kraft's Inequality $\overset{③}{\Longleftarrow}$ Uniquely decodable codes

$②$

Note: $\lg \triangleq \log_2$

$$l_i^* = \lg \frac{1}{p_i}$$

$$E[L] = \sum_i p_i \cdot l_i \geq \sum_i p_i \cdot l_i^* = \sum_i p_i \lg \frac{1}{p_i} = H(X)$$

Claim:   Instantaneous (prefix-free) code $\iff$ Kraft's Inequality

Pf ($\Rightarrow$):   Let $l_{max}$ be largest codeword

$$0$$
$$10$$
$$110$$
$$111$$



- Code of length $l_i$ wipes out how many leaves? $2^{l_{max} - l_i}$

- tree only has $2^{l_{max}}$ leaves

$$\sum_i 2^{l_{max} - l_i} \leq 2^{l_{max}}$$

Prefix-free property:
internal node code
makes unavailable all $\rightarrow$ dividing both sides by $2^{l_{max}}$
possible codes in
subtree below.    $\Rightarrow$ $\boxed{\sum_i 2^{-l_i} \leq 1}$

- What is $\left(\sum\limits_{x \in \mathcal{X}} D^{-\ell(x)}\right)^k$ ?

Now ③

$x_2 \in \mathcal{X}$

| | a | c | g | t |
|---|---|---|---|---|
| a | $D^{-\ell(a)} D^{-\ell(a)}$  $D^{-\ell(a)} D^{-\ell(c)}$ | | $\cdots$ | |
| c | | | | |
| g | | | | |
| t | | | | |

$x_1 \in \mathcal{X}$

- Consider $\mathcal{X} = \{a, c, g, t\}$
  $\mathcal{D} = \{0,1\}$ and $k = 2$

  $a \to 0$
  $c \to 10$
  $g \to 110$
  $t \to 111$

- $\left(\sum\limits_{x \in \mathcal{X}} D^{-\ell(x)}\right)^2$

$= \sum\limits_{x_1 \in \mathcal{X}} \sum\limits_{x_2 \in \mathcal{X}} D^{-\ell(x_1)} \cdot D^{-\ell(x_2)}$

$= \sum\limits_{x_1 x_2 \in \mathcal{X}^2} D^{-\ell(x_1)} \cdot D^{-\ell(x_2)} = \sum\limits_{x_1 x_2 \in \mathcal{X}^2} D^{-(\ell(x_1) + \ell(x_2))}$

$\mathcal{X}^2 = \mathcal{X} \times \mathcal{X}$

$= \sum\limits_{x^2 \in \mathcal{X}^2} D^{-\ell(x^2)}$

$x^2 \in \mathcal{X}^2 = \mathcal{X} \times \mathcal{X}$

$\ell(aa) = 2$
$\ell(ac) = 3$
$\vdots$
$\ell(tt) = 6$

$= \sum\limits_{m=1}^{2 \cdot \ell_{max}} a(m) \cdot D^{-m} \leq \sum\limits_{m=1}^{2 \cdot \ell_{max}} D^m \cdot D^{-m} = 2 \cdot \ell_{max}$

$a(m) = \# x^2 \in \mathcal{X}^2$ where $\ell(x^2) = m$

(✗✗✗) unique decodability $\Rightarrow a(m) \leq D^m$ !

In general, $\left(\sum\limits_{x \in \mathcal{X}} D^{-\ell(x)}\right)^k \leq k \cdot \ell_{max}$

$$\left( \sum_{x \in \mathcal{X}} D^{-\ell(x)} \right)^k \leq k \cdot \ell_{max}$$

$$\Rightarrow \quad \sum_{x \in \mathcal{X}} D^{-\ell(x)} \leq \left( k \cdot \ell_{max} \right)^{1/k}$$

$$= e^{\ln\left[ \left( k \cdot \ell_{max} \right)^{1/k} \right]}$$

$$= e^{\frac{\ln(k \cdot \ell_{max})}{k}}$$

must hold for all $k$, e.g., $k \to \infty$

$$\lim_{k \to \infty} = ? \quad \Rightarrow \quad e^0 = 1$$

$$\Rightarrow \quad \sum_{x \in \mathcal{X}} D^{-\ell(x)} \leq 1$$

# Part 1: Theory
# L07: Compression
# (uniquely decodable codes continued)

Javed Aslam, Wolfgang Gatterbauer

cs7840 Foundations and Applications of Information Theory (fa24)

https://northeastern-datalab.github.io/cs7840/fa24/

9/25/2024

Last time

Basic results in Information theory

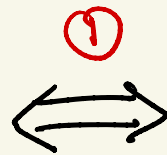Kraft's Inequality for uniquely decodable codes

Today

- Implications of Kraft's Inequality
- Shannon Codes
- Block coding
- Asymptotic Equipartition Property (AEP)

Next time

- Continue...

instantaneous (prefix-free) code $\overset{①}{\Longleftrightarrow}$ Kraft's Inequality $\overset{③}{\longleftarrow}$ uniquely decodable codes

$$② \Downarrow$$

Note:

$$\lg \triangleq \log_2$$

$$\ell_i^* = \lg \frac{1}{p_i}$$

$$E[L] = \sum_i p_i \cdot \ell_i \geq \sum_i p_i \cdot \ell_i^* = \sum_i p_i \lg \frac{1}{p_i} = H(X)$$

# Bounds on n! : Stirling's Approximation  <span style="color:red">(Mathematical Digression)</span>

$$n! = n \cdot (n-1) \cdot (n-2) \cdots 3 \cdot 2 \cdot 1$$

Approximations:

good: $\left(\frac{n}{e}\right)^n$

better: $\left(\frac{n}{e}\right)^n \sqrt{2\pi n}$

even better! $\left(\frac{n}{e}\right)^n \sqrt{2\pi n} \left(1 + O\left(\frac{1}{n}\right)\right)$

best: $\left(\frac{n}{e}\right)^n \sqrt{2\pi n} \; e^{\lambda_n}$

where $\dfrac{1}{12n+1} < \lambda_n < \dfrac{1}{12n}$

Example

| | lg |
|---|---|
| $10! = 3,628,800$ | $21.79$ |
| $\left(\frac{n}{e}\right)^n = 453,999$ | $18.79$ |
| $\left(\frac{n}{e}\right)^n \sqrt{2\pi n} = 3,598,696$ | $21.78$ |
| $\left(\frac{n}{e}\right)^n \sqrt{2\pi n} \; e^{\frac{1}{12n}} = 3,628,810$ | $21.79$ |

Kraft's Inequality: $\sum_i 2^{-l_i} \leq 1$

- consider $r_i = \dfrac{2^{-l_i}}{\sum_j 2^{-l_j}}$ — this is a distribution

- $L - H(x) = \sum_i p_i \cdot l_i - \left(-\sum_i p_i \lg p_i\right)$

$$= \sum_i p_i \cdot \lg \frac{1}{2^{-l_i}} + \sum_i p_i \lg p_i$$

$$= \sum_i p_i \cdot \lg \frac{1}{r_i \cdot \sum_j 2^{-l_j}} + \sum_i p_i \lg p_i$$

$$= \sum_i p_i \lg \frac{1}{r_i} + \sum_i p_i \cdot \lg \frac{1}{\sum_j 2^{-l_j}} + \sum_i p_i \lg p_i$$

$$= \sum_i p_i \lg \frac{p_i}{r_i} + \lg \frac{1}{\sum_j 2^{-l_j}}$$

$$= D(\vec{p} \| \vec{r}) + \lg \frac{1}{\sum_j 2^{-l_j}}$$

$\geq 0 \qquad\qquad \geq 0$

$= 0$ iff

$2^{-l_i} = p_i$ and $\sum_j 2^{-l_j} = 1$

$\iff l_i = \lg \frac{1}{p_i}$

To show:  $H(x) \leq L^* \leq H(x)+1$       $L^*$ is opt code

Shannon codes:  $l_i = \lceil \lg \frac{1}{p_i} \rceil$

Claim:  these $l_i$ satisfy Kraft's inequality

Pf:  $\sum_i 2^{-\lceil \lg \frac{1}{p_i} \rceil} \leq \sum_i 2^{-\lg \frac{1}{p_i}} = \sum_i 2^{\lg p_i} = \sum_i p_i = 1$

$\Rightarrow$ valid code lengths; can easily be turned into prefix-free code by earlier results

Now:  $\lg \frac{1}{p_i} \leq l_i = \lceil \lg \frac{1}{p_i} \rceil < \lg \frac{1}{p_i} + 1$

$\sum_i p_i \lg \frac{1}{p_i} \leq \sum_i p_i \cdot l_i < \sum_i p_i (\lg \frac{1}{p_i} + 1)$

$H(x) \leq L < H(x) + 1$

Block Coding: encode blocks of length $n$ at a time

- induced distribution $p(x_1, x_2 \cdots x_n) = p(x_1) \cdot p(x_2) \cdots p(x_n)$

- $L_n = \frac{1}{n} \sum p(x_1 \cdots x_n) \cdot \ell(x_1 \cdots x_n)$

- From last slide

$$H(x_1 \cdots x_n) \le E\, \ell(x_1 \cdots x_n) \le H(x_1 \cdots x_n) + 1$$

$$\Rightarrow \quad n\, H(x) \le n \cdot L_n \le n\, H(x) + 1$$

$$\Rightarrow \quad H(x) \le L_n \le H(x) + \frac{1}{n}$$

$\therefore$ Can drive inefficiency down arbitrarily small by using larger and larger blocks.

- Issue: Code book grows exponentially with block length
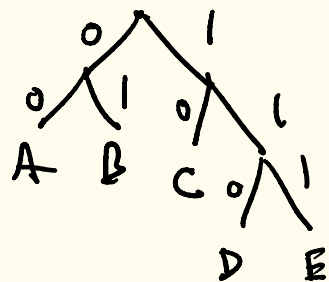
2 other codes ... Example:

$H(\vec{p}) = 2.23$

| | A | B | C | D | E | |
|---|---|---|---|---|---|---|
| $p_i$ | .35 | .17 | .17 | .16 | .15 | |
| $\lg \frac{1}{p_i}$ | 1.51 | 2.56 | 2.56 | 2.64 | 2.74 | |
| $\lceil \lg \frac{1}{p_i} \rceil$ | 2 | 3 | 3 | 3 | 3 | Shannon |

$L_s = 2.65 \text{ bpc}$

(bpc - bits per character on average)

**Fano:**

| A | B | C | D | E |
|---|---|---|---|---|
| .35 | .17 | .17 | .16 | .15 |

· Sort prob.
· Split as close to 50/50 as possible
· Recurse



A 00
B 01
C 10
D 110
E 111

2.31 bpc

$L_F = 2.31 \text{ bpc}$

**Huffman**

· Combine least probable events bottom-up, creating combined events along the way

| A | B | C | {DE} |
|---|---|---|------|
| .35 | .17 | .17 | .31 |

| A | {BC} | {DE} |
|---|------|------|
| .35 | .34 | .31 |



A 0
B 100
C 101
D 110
E 111

2.3 bpc

$L_H = 2.3 \text{ bpc}$

AEP : Asymptotic Equipartition Property

Consider biased coin $\begin{cases} Pr(H) = 1/3 \\ Pr\{T\} = 2/3 \end{cases}$

Flip coin $n$ times.

- what is most <u>probable</u> outcome?  TTT---T  $Pr = (2/3)^n$

- <u>typical</u> sequence has about $1/3 H$ $2/3 T$  $Pr = (1/3)^{\frac{n}{3}} \cdot (2/3)^{2n/3}$

- Let  $Pr[\text{"typical sequence"}] = x = (1/3)^{n/3} \cdot (2/3)^{2n/3}$

- what is $x$?

$$\lg x = \frac{n}{3} \lg (1/3) + \frac{2n}{3} \lg (2/3)$$

$$= n \cdot \left[ \frac{1}{3} \lg 1/3 + \frac{2}{3} \lg 2/3 \right]$$

$$= - n \cdot H(x)$$

$$\Rightarrow Pr[\text{"typical sequence"}] = x = 2^{-n H(x)}$$

How many strings of $\frac{1}{3}$ H & $\frac{2}{3}$ T ?

$\binom{n}{n/3}$ typical sequences

$$\binom{n}{n/3} = \frac{n!}{\left(\frac{n}{3}\right)! \left(\frac{2n}{3}\right)!} \sim \frac{\left(\frac{n}{e}\right)^n}{\left(\frac{n/3}{e}\right)^{n/3} \cdot \left(\frac{2n/3}{e}\right)^{2n/3}} = \frac{\left(\frac{n}{e}\right)^n}{\left(\frac{n}{e}\right)^{n/3} \cdot \left(\frac{1}{3}\right)^{n/3} \cdot \left(\frac{n}{e}\right)^{2n/3} \cdot \left(\frac{2}{3}\right)^{2n/3}}$$

$$= \frac{1}{\left(\frac{1}{3}\right)^{n/3} \cdot \left(\frac{2}{3}\right)^{2n/3}} = 2^{n H(x)}$$

- There are about $2^{n H(X)}$ typical sequences
- Each has about $2^{-n H(x)}$ probability

$\Rightarrow$ almost everything one is likely to see is <u>typical</u> and <u>equally likely</u>

# Part 1: Theory
# L08: Compression
# (AEP = Asymptotic Equipartition Property)

Javed Aslam, Wolfgang Gatterbauer

cs7840 Foundations and Applications of Information Theory (fa24)

https://northeastern-datalab.github.io/cs7840/fa24/

9/30/2024

## Last time

- Implications of Kraft's Inequality
- Shannon, Fano, Huffman Codes
- Block Coding
- Motivating the AEP

## Today

- AEP
  - formal defs. & proofs
  - Consequences

~~~~~~~~~~

- Finish fundamentals

## Next time

- Continue Fundamentals

# (weak) Law of Large Numbers <span style="color:red">(Mathematical Digression)</span>

Roughly: Sample mean converges to true mean (in probability)

Technically:  Let  $\overline{X}_n = \frac{1}{n}\sum_{i=1}^{n} X_i$  ;  $E[X] = \mu$

- $\overline{X}_n \xrightarrow{P} \mu$  as  $n \to \infty$

- $\lim_{n \to \infty} Pr(|\overline{X}_n - \mu| < \varepsilon) = 1$  for any  $\varepsilon > 0$

- $(\forall \varepsilon > 0)(\forall \delta > 0)(\exists n_0)\ Pr(|\overline{X}_n - \mu| < \varepsilon) > 1 - \delta\ \ \forall n \geq n_0$

Thm: (AEP) If $x_1, x_2, \dots x_n$ are i.i.d. $\sim p(x)$, then

$$-\frac{1}{n} \lg p(x_1 x_2 - x_n) \longrightarrow H(X) \text{ in probability.}$$

Intuition: Consider an actual sequence $x_1 x_2 - x_n$

$$-\frac{1}{n} \lg p(x_1 x_2 - x_n) = \frac{1}{n} \boxed{\lg \frac{1}{p(x_1 x_2 - x_n)}} \sim \text{length of Shannon code for block}$$

avg. length per message

$\longrightarrow$ converges to $H(X)$ as $n \to \infty$ by block coding

Formal Proof: Since $x_i$ are i.i.d., then so are r.v. $\lg p(x_i)$ & $\lg \frac{1}{p(x_i)}$

By LLN:
$$-\frac{1}{n} \lg p(x_1 x_2 - x_n) = -\frac{1}{n} \sum_i \lg p(x_i)$$

Informal Relationship
$$-\frac{1}{n} \lg p(x_1 x_2 - x_n) \sim H(X)$$

$$\Leftrightarrow p(x_1 x_2 - x_n) \sim 2^{-n H(X)}$$

$$= \frac{1}{n} \sum_i \lg \frac{1}{p(x_i)}$$

$$\longrightarrow E\left[\lg \frac{1}{p(x)}\right] \text{ in probability}$$

$$= H(X)$$

Def: The __typical set__ $A_\epsilon^{(n)}$ with respect to $p(x)$ is the set of all sequences $(x_1 x_2 \ldots x_n) \in \mathcal{X}^n$ such that

$$2^{-n(H(X)+\epsilon)} \leq p(x_1 x_2 \ldots x_n) \leq 2^{-n(H(X)-\epsilon)}$$

<span style="color:red">← empirical probability of sequence</span>

thm: ① If $(x_1 x_2 \ldots x_n) \in A_\epsilon^{(n)}$, then $H(X)-\epsilon \leq -\frac{1}{n} \lg p(x_1 \ldots x_n) \leq H(X)+\epsilon$

② $\Pr[A_\epsilon^{(n)}] > 1-\epsilon$ for $n$ sufficiently large

③ $|A_\epsilon^{(n)}| \leq 2^{n(H(X)+\epsilon)}$

④ $|A_\epsilon^{(n)}| \geq (1-\epsilon) 2^{n(H(X)-\epsilon)}$ for $n$ sufficiently large

Pf: ① Immediate from definition of typical set

② By LLN $(\forall \epsilon > 0)(\forall \delta > 0)(\exists n_0)$

$$\Pr\left[ \underbrace{\left| -\frac{1}{n} \lg p(x_1 x_2 \ldots x_n) - H(X) \right| < \epsilon}_{\color{red}{\text{typical set by ①}}} \right] > 1-\delta \qquad \forall n > n_0$$

<span style="color:red">↑ choose $\delta = \epsilon$</span>

③ $|A_\varepsilon^{(n)}| \leq 2^{n(H(x)+\varepsilon)}$

④ $|A_\varepsilon^{(n)}| \geq (1-\varepsilon) 2^{n(H(x)-\varepsilon)}$ for $n$ sufficiently large

Let $\vec{x} = (x_1, x_2 - x_n)$

Pf ③:

$$1 = \sum_{\vec{x} \in \chi^n} P(\vec{x})$$

$$\geq \sum_{\vec{x} \in A_\varepsilon^{(n)}} P(\vec{x})$$

$$\geq \sum_{\vec{x} \in A_\varepsilon^{(n)}} 2^{-n(H(X)+\varepsilon)}$$

$$= 2^{-n(H(X)+\varepsilon)} \cdot |A_\varepsilon^{(n)}|$$

$$\Rightarrow |A_\varepsilon^{(n)}| \leq 2^{n(H(X)+\varepsilon)}$$

Pf ④: For $n$ sufficiently large, $\Pr[A_\varepsilon^{(n)}] > 1-\varepsilon$ so...

$$1-\varepsilon < \Pr[A_\varepsilon^{(n)}]$$

$$= \sum_{\vec{x} \in A_\varepsilon^{(n)}} P(\vec{x})$$

$$\leq \sum_{\vec{x} \in A_\varepsilon^{(n)}} 2^{-n(H(X)-\varepsilon)}$$

$$= 2^{-n(H(X)-\varepsilon)} \cdot |A_\varepsilon^{(n)}|$$

$$\Rightarrow |A_\varepsilon^{(n)}| \geq (1-\varepsilon) \cdot 2^{n(H(X)-\varepsilon)}$$

**Upshot:**

- $2^{-n(H(x)+\varepsilon)} \leq p(x_1, x_2 \dots x_n) \leq 2^{-n(H(x)-\varepsilon)}$

  *a typical sequence has empirical probability $\sim 2^{-n H(x)}$*

- $|A_\varepsilon^{(n)}| \leq 2^{n(H(x)+\varepsilon)}$
- $|A_\varepsilon^{(n)}| \geq (1-\varepsilon)\, 2^{n(H(x)-\varepsilon)}$

  *there are $\sim 2^{n H(x)}$ typical sequences*

- $Pr[A_\varepsilon^{(n)}] > 1-\varepsilon$   for $n$ sufficiently large

  *typical sequences contain almost all probability*

# Immediate Consequence for compression



$x^n$

typical set $\sim 2^{n\,H(x)}$ in size

$\Rightarrow$ high prob events are typical sequences which use $n\,H(x)+1$ bits or $H(x) + 1/n$ per encoded message, on average.

Block coding compression method (roughly):

- If $\vec{x}$ is <u>typical</u>, start with a 0 and encode the exact typical sequence in straight binary using
$$\lg\left(2^{n\,H(x)}\right) = n\,H(x)\text{ additional bits}$$

- If $\vec{x}$ is <u>not typical</u>, start with a 1 and encode the exact atypical sequence in straight binary using
$$\lg\left(|x|^n\right) = n\,\lg|x|\text{ additional bits}$$

(see text for more careful treatment taking into account $\epsilon$, etc.