

Brief Summary on Information Theory Basics

Petros Elia

EURECOM Fall 2024

November 19, 2024

- ① Introduction
- ② Definitions and Inequalities
- ③ Fundamental Limits of Compression
- ④ Communication and Channel Capacity
- ⑤ Shannon's Channel Coding Theorem
- ⑥ Information Theory of Continuous Data
- ⑦ The Gaussian Channel

- ① Introduction
- ② Definitions and Inequalities
- ③ Fundamental Limits of Compression
- ④ Communication and Channel Capacity
- ⑤ Shannon's Channel Coding Theorem
- ⑥ Information Theory of Continuous Data
- ⑦ The Gaussian Channel

1 Big impact of Information Theory

| 3

- ▶ Two main impactful applications
 - > Source compression
 - > Communications
- ▶ Also applies in
 - > Distributed computing
 - > Learning
 - > Investment theory
 - > Statistical physics

For compression:



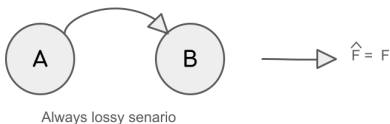
- ▶ Lossless or Lossy compression
- ▶ Q: What is the minimum amount of space needed?

1 Applications of Information Theory

| 4

For communications:

- ▶ Q: How fast can we transmit such that $\hat{F} = F$?



- ▶ Depending on the physics, the fundamental limits change.
- ▶ Basic channel, BSC is defined by the probability of transition p .
- ▶ Multi-user scenario: both in compression and in communications

① Introduction

② Definitions and Inequalities

- Definitions

- Convexity and Jensen's inequality

- Bound on Entropy

- Data processing inequality

- Fano's inequality

③ Fundamental Limits of Compression

④ Communication and Channel Capacity

⑤ Shannon's Channel Coding Theorem

2 Definitions - Concept of Information in Bits

| 6

- ▶ Random variable (R.V.) $X = \text{"Volleyball Winner at Paris Olympics"}$
- ▶ Alphabet $\mathcal{X} = \{\text{Argentina, France, Italy, Brazil}\}$
 - > with probability p_1, p_2, p_3, p_4 .
 - > Random event $A \rightarrow p$ probability
- ▶ What is the Information content of an event (properties to be considered)

$$\begin{array}{ll} I(p = 1) = 0 & I(p \downarrow) \uparrow \\ I(p) \geq 0 & I(p \rightarrow 0) \rightarrow \infty \end{array}$$

- ▶ Also, if A, B independent; we should have $\rightarrow I(A, B) = I(A) + I(B)$
 - > Also nicely reflecting $p(A, B) = p_A p_B$
- ▶ Shannon information:
 - > $I(p) = -\log_2(p)$ in bits (assume this)
 - > $I(p) = -\log_e(p)$ in nats
- ▶ Ex: Rain = 99%, Sunny = 1%.
 - > $I(\text{'Today is raining'}) = -\log(0.99) = 0.014 \text{ bits}$
 - > $I(\text{'Today is sunny'}) = -\log(0.01) = 6.64 \text{ bits}$

► Definition 1: **Entropy**

$$H(X) = - \sum_{i=1}^{|\mathcal{X}|} p_i \log(p_i) \geq 0 \quad (\text{average information / } \underline{\text{UNCERTAINTY}})$$

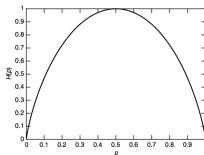
> Alphabet cardinality $|\mathcal{X}|$

► **Binary Entropy** for binary X (e.g. Rain vs. Sunny)

> $\mathcal{X} = \{0, 1\}$: $p(X = 0) = p$, $p(X = 1) = 1 - p$:

$$H(X) = -p \log p - (1 - p) \log(1 - p)$$

> Maximized at $p = 0.5$.



2 Entropy - Example

| 8

Calculate: $H(X) = -\sum_i^{|X|} p_i \log(p_i) \geq 0$

- For $X = \{(X^{(1)}, \frac{1}{2}), (X^{(2)}, \frac{1}{4}), (X^{(3)}, \frac{1}{8}), (X^{(4)}, \frac{1}{8})\}$
 - > {Volleyball: Argentina, France, Italy, Brazil}

$$\begin{aligned} H(X) &= -\frac{1}{2} \log_2 \frac{1}{2} - \frac{1}{4} \log_2 \frac{1}{4} - \frac{1}{8} \log_2 \frac{1}{8} - \frac{1}{8} \log_2 \frac{1}{8} \\ &= \frac{1}{2} \log_2 2 + \frac{1}{4} \log_2 2^2 + \frac{2}{8} \log_2 2^3 \\ &= \frac{1}{2} + \frac{2}{4} + \frac{2}{8} \times 3 = \frac{7}{4} \text{ bits} \end{aligned}$$

- Vs. 100m Sprint:
 $X = \{(X^{(1)}, \frac{9999}{10000}), (X^{(2)}, \frac{1}{30000}), (X^{(3)}, \frac{1}{30000}), (X^{(4)}, \frac{1}{30000})\}$
 - > {Usain Bolt, Me, You, Obama}

$$H(X) = -\frac{9999}{10000} \log_2 \frac{9999}{10000} - \frac{3}{30000} \log_2 \frac{1}{30000} = 0.0016 \text{ bits}$$

Definition 2: **Joint Entropy** of two R.Vs. X, Y

- ▶ Let X, Y be two discrete R.V (discrete alphabets)
- ▶ $X \in \mathcal{X} = \{x^{(1)}, x^{(2)}, \dots\}$, $Y \in \mathcal{Y} = \{y^{(1)}, y^{(2)}, \dots\}$
- ▶ $p(X), p(Y), p(X, Y)$
- ▶ Consider $Z = (X, Y)$, $Z \in \mathcal{Z} = \mathcal{X} \times \mathcal{Y}$
- ▶ Thus

$$H(X, Y) = -E[\log p(X, Y)] = - \sum_{x \in \mathcal{X}} \sum_{y \in \mathcal{Y}} p(x, y) \log p(x, y)$$

- > Naturally $H(X, Y) \leq H(X) + H(Y)$

Definition 3: **Conditional Entropy:** $H(X|Y)$

$$\begin{aligned} H(X|Y) &\triangleq E_y[H(X|Y=y)] \\ &= - \sum_{y \in \mathcal{Y}} \sum_{x \in \mathcal{X}} p(x|Y=y)p(y) \log p(x|Y=y) \\ &\stackrel{(BR)}{=} - \sum_{y \in \mathcal{Y}} \sum_{x \in \mathcal{X}} p(x,y) \log \frac{p(x,y)}{p(y)} \end{aligned}$$

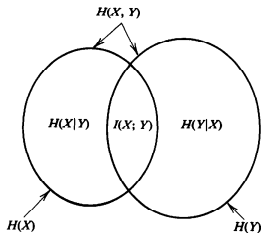
► Naturally **Conditioning Reduces Entropy**

$$H(X|Y) \leq H(X)$$

> with equality if X and Y are independent.

Definition 4: **Mutual information:**

$$\begin{aligned} I(X; Y) &= H(X) - H(X|Y) \\ &= H(Y) - H(Y|X) \\ &= H(X) + H(Y) - H(X, Y) \end{aligned}$$



- $I(X; Y) = 0$ when X and Y are independent.

Definition 5: **Chain Rule:**

- ▶ X, Y pair of discrete R.V
- ▶ Then

$$H(X, Y) = H(Y) + H(X|Y) = H(Y, X) = H(X) + H(Y|X)$$

- ▶ Example

$$H(X_1, X_2, \dots, X_n) = H(X_1) + H(X_2|X_1) + H(X_3|X_1, X_2) + \dots + H(X_n|X_1, X_2, \dots, X_{n-1})$$

2 Chain Rule and Conditional Mutual Information

| 13

- ▶ We have two chain rules
 - > Recall: Entropy chain rule

$$H(X, Y) = H(X) + H(Y|X) = H(Y) + H(X|Y)$$

- > Information chain rule

$$I((X, Y); Z) = I(X; Z) + \underbrace{I(Y; Z|X)}_*$$

- ▶ ***Conditional mutual information**

$$I(Y; Z|X) = H(Y|X) - H(Y|Z, X)$$

Definition 6: **Relative entropy** (Kullback-Leibler distance/divergence)

- ▶ Given RV X and two distributions $p(X)$ and $q(X)$, then what is the distance between them?

$$D(p||q) = \sum_{x \in \mathcal{X}} p(x) \log \frac{p(x)}{q(x)} \geq 0$$

- > Takes form of Informational Distance

$$D(p||q) = E_{p(X)}[\log p(X) - \log q(X)]$$

- ▶ Can measure the degree of independence between X and Y , i.e:

$$D(p(X, Y)||p(X)P(Y)) \text{ Thus most importantly}$$

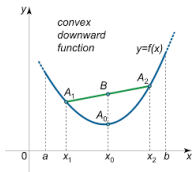
$$D(p(X, Y)||p(X)p(Y)) = I(X; Y)$$

2 Convexity and Jensen's inequality

| 15

- Convexity: For $x = px_1 + (1 - p)x_2, p \in \mathbb{R}^+, p \leq 1$, then

$$f(x) \leq pf(x_1) + (1 - p)f(x_2)$$



- X : discrete random variable
- f : convex function

$$E[f(X)] \geq f(E[X])$$

- f : concave function

$$E[f(X)] \leq f(E[X])$$

- ▶ X is a discrete RV $x \in \mathcal{X}$
 - > with $|\mathcal{X}|$ elements, and distribution p .
- ▶ Then

$$H(X) \leq \log |\mathcal{X}|$$

- > with equality if

$$p(x) = \frac{1}{|\mathcal{X}|}, \forall x \in \mathcal{X} : \quad (\text{uniform distribution})$$

- > uniform is scenario with highest uncertainty.
- ▶ Proof Sketch: (any p , q uniform). Use² that $D(p||q) \geq 0$.

$$\begin{aligned} D(p||q) &= \sum p(x) \log \frac{p(x)}{q(x)} = \sum p(x) \log(p(x)|\mathcal{X}|) \\ &= \sum p(x) \log p(x) + \sum p(x) \log |\mathcal{X}| = -H(X) + \log |\mathcal{X}| \geq 0 \end{aligned}$$

²From Log-Sum Ineq: $\sum_{i=1}^{|\mathcal{X}|} p_i \log \frac{p_i}{q_i} \geq (\sum_{i=1}^{|\mathcal{X}|} p_i) \log \left(\frac{(\sum_{i=1}^{|\mathcal{X}|} p_i)}{(\sum_{i=1}^{|\mathcal{X}|} q_i)} \right)$

2 Data Processing Inequality

| 17

- ▶ X is a R.V.



- ▶ Then

$$I(X; Z) \leq I(X; Y)$$

- ▶ You cannot increase the information you get for Z from X , by "massaging" X

2 Markov chains (MC): (X, Y, Z)

| 18

- ▶ (X, Y, Z) form a MC if and only if
 - > $p(Z|X, Y) \stackrel{(mc1)}{=} p(Z|Y)$
 - > $p(X, Y, Z) \stackrel{(mc2)}{=} p(X)p(Y|X)p(Z|Y)$
 - > (mc2): Apply chain rule and then (mc1)
- ▶ Meaning: dependency of z from x happens through y only
- ▶ (mc1) and (mc2) are equivalent
 - > Sketch of Alternate Proof that $(mc1) \rightarrow (mc2)$

$$\begin{aligned} P(X, Y, Z) &\stackrel{(br)}{=} P(Z|X, Y)P(X, Y) \\ &\stackrel{(mc1)}{=} P(Z|Y)P(X, Y) \stackrel{(br)}{=} P(Z|Y)P(Y|X)P(X) \rightarrow (mc2) \end{aligned}$$

2 Proof of DPI using Markov chains

| 19

- ▶ Want to prove that $I(X, Z) \leq I(X, Y)$
- ▶ Consider Information Chain Rule (icr)

$$\begin{aligned} I(X; Y, Z) &\stackrel{(icr)}{=} I(X; Y) + I(X; Z|Y) \\ &\stackrel{(icr)}{=} I(X; Z) + I(X; Y|Z) \\ &\Rightarrow I(X; Y) + I(X; Z|Y) = I(X; Z) + I(X; Y|Z) \quad (**) \end{aligned}$$

- ▶ Now note that:

$$P(X, Z|Y) \stackrel{(br)}{=} \frac{P(X, Y, Z)}{P(Y)} \stackrel{(br)}{=} \frac{P(Z|Y, X)P(X, Y)}{P(Y)} \stackrel{(mc+br)}{=} P(Z|Y)P(X|Y)$$

- ▶ Thus, given Y , then X and Z are independent
- ▶ Thus $I(X; Z|Y) = 0$
- ▶ Thus from (**), we have

$$I(X; Y) = I(X; Z) + I(X; Y|Z)$$

- ▶ Thus $I(X, Z) \leq I(X, Y)$. (DPI proved)

2 Fano's inequality

| 20

- ▶ X : unknown R.V.
- ▶ We observe Y , correlated with X via $P(Y|X)$
- ▶ We want to build an estimate of X using Y
 - > build $\hat{X} = g(Y)$
- ▶ Prob of error: $P_e = \text{Prob}(\hat{X} \neq X)$
- ▶ Fano bound: $P_e \geq ?$

2 Deriving Fano's Bound

| 21

► Let E be a binary R.V. $E(\text{error event}) = \begin{cases} 0 & \hat{X} = X \\ 1 & \hat{X} \neq X \end{cases} \begin{matrix} 1 - P_e \\ P_e \end{matrix}$

► $H(E, X|Y) \stackrel{(cr)}{=} H(X|Y) + H(E|X, Y) = H(X|Y) \quad (*)$

> Since $\rightarrow H(E|X, Y) = 0$, since E is fully determined by X and Y

► Rewrite same

$$H(E, X|Y) \stackrel{(cr)}{=} H(E|Y) + H(X|E, Y) \quad (**)$$

$$H(X|E, Y) = H(X|Y) - H(E|Y) \stackrel{(cre)}{\geq} H(X|Y) - H(E) = H(X|Y) - H(P_e) \quad (***)$$

$$H(X|E, Y) = (1 - P_e)H(X|Y, E = 0) + P_e H(X|Y, E = 1) \stackrel{(maxH)}{\leq} P_e \log(|\mathcal{X}| - 1)$$

- ① Introduction
- ② Definitions and Inequalities
- ③ Fundamental Limits of Compression
(Weak) Law of Large Numbers
Asymptotic Equipartition Property (AEP)
- ④ Communication and Channel Capacity
- ⑤ Shannon's Channel Coding Theorem
- ⑥ Information Theory of Continuous Data

- ▶ Definition: **Random process** $\{X_i\} : X_1, X_2, \dots, X_n, \quad n \in \mathbb{Z}$
- ▶ Definition: i.i.d. random process
 - > X_i, X_j are independent $\forall i \neq j$, and $p(X_i) = p(X_j)$
 - > i.i.d. is worst case for compression.

3 (Weak) Law of Large Numbers

| 24

- ▶ Consider an i.i.d. R.P. $\{Z_i\}$
 - > $E[Z_i] = \mu$, $E[(Z_i - \mu)^2] = \sigma^2$.
- ▶ Define sample mean: $\bar{Z}_n = \frac{1}{n} \sum_{i=1}^n Z_i$
- ▶ LLN: \bar{Z}_n converges “in probability” to μ .

$$Prob(|\bar{Z}_n - \mu| \geq \epsilon) \xrightarrow{n \rightarrow \infty} 0, \quad \forall \epsilon > 0$$

- ▶ We say that $\bar{Z}_n \xrightarrow{\text{probability}} \mu$
- ▶ But LLN says even more ...

$$Prob(|\bar{Z}_n - \mu| \geq \epsilon) \leq \frac{\sigma^2}{n\epsilon} \xrightarrow{n \rightarrow \infty} 0$$

- ▶ This convergence is slower for correlated data
 - > more data needed to explore the entire space...!

3 Asymptotic equipartition

| 25

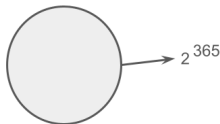
- ▶ Example: $n = 365$, $X_i \in \mathcal{X} = \{r, s\}$
 - > (weather data simplified to i.i.d.)
- ▶ $X^n = (X_1, X_2, \dots, X_n)$, $Prob(r) = 10\%$, $Prob(s) = 90\%$
- ▶ $X^n = (r, r, r, \dots, r)$, $X^n = (s, s, s, \dots, s)$ are atypical sequences:
 - > possible mathematically, but really low probability.
- ▶ $X^n = (s, s, r, s, s, \dots, r, s)$ more typical sequence: what we expect! about 10% of “r” and 90% “s”.

3 Asymptotic equipartition

| 26

- ▶ Take $\{X_i\}$ on i.i.d. R.P. $p(X_i)$ is the distribution of $X_i \sim X$ ($X_i \in \mathbb{R}$)
- ▶ Let $Z_i \triangleq \log p(X_i)$. $\{Z_i\}$ also i.i.d. process. ($Z_i \in \mathbb{R}^-$)
 - > i.e., Draw: $X^n = [X_1 \ X_2 \cdots X_n]$
 - > create instance $Z^n = [\log p(X_1) \ \log p(X_2) \ \dots \ \log p(X_n)] = [Z_1 \dots Z_n]$
- ▶ $\frac{1}{n} \sum_{i=1}^n Z_i \rightarrow E(Z_i) = E[\log p(X_i)] = \sum_{x \in \mathcal{X}} p(x) \log p(x) = -H(X)$
$$\frac{1}{n \log p(X_1, \dots, X_n) = \frac{1}{n} \log p(X^n)} \xrightarrow{n \rightarrow \infty} -H(X)$$

$$\Rightarrow P(X^n) \approx 2^{-nH(X)}, \ n \rightarrow \infty \quad (\text{simply intuition})$$



- ▶ $P(X^n) \approx \frac{1}{2^{nH(X)}}$ and not $\frac{1}{2^{365}}$
 - > Many sequences are negligible.
 - > Key for Compression: place focus on "typical" sequences

3 Asymptotic equipartition

| 27

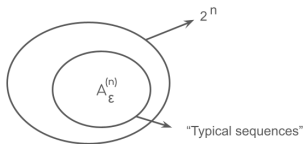
- ▶ LLN:

$$\text{Prob}(|\frac{1}{n} \log P(X^n) + H(X)| > \epsilon) \rightarrow 0$$

- ▶ Define typical set: $\mathcal{A}_\epsilon^{(n)}$

$$\mathcal{A}_\epsilon^{(n)} = \{X^n \in \mathcal{X}^n \text{ s.t. } 2^{-n(H(X)+\epsilon)} \leq P(X^n) \leq 2^{-n(H(X)-\epsilon)}\}$$

- ▶ Typical sequences: high likelihood to be observed

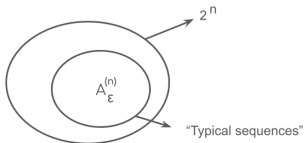


3 Basic properties of the typical set

| 28

For large n

- ▶ 1). $Prob(X^n \in \mathcal{A}_\epsilon^{(n)}) \stackrel{(LLN)}{\geq} 1 - \epsilon$ (easy - restricted)
- ▶ 2). $|\mathcal{A}_\epsilon^{(n)}| \approx 2^{n(H(X))} \ll 2^n$ (binary)
 - > from $\sum_{X^n \in \mathcal{A}_\epsilon^{(n)}} P(X^n) \approx 1$



- ▶ Compress: $X^n = (X_1, X_2, \dots, X_n) \xrightarrow{(\xi)} \xi(X^n) = (0, 1, 1, \dots, 0)$
 - > to vector of length $\ell = \mathcal{L}(\xi(X^n)) \ll n$
- ▶ Recall $\mathcal{A}_\epsilon^{(n)}$ contains no more than $2^{n(H(X)+\epsilon)}$ elements/vectors
- ▶ Take a sequence X^n , examine if $X^n \in \mathcal{A}_\epsilon^{(n)}$ or not
 - > a) if $X^n \in \mathcal{A}_\epsilon^{(n)}$, map X^n to index of $\lceil nH(X) + \epsilon \rceil$ bits
 - > b) if $X^n \notin \mathcal{A}_\epsilon^{(n)}$, map X^n to index of $\lceil \log |X^n| \rceil$ bits.
 - > c) add 1 bit to indicate if case a) or case b).
- ▶ This a compression based on typicality

- Define $\mathcal{L}(\xi(X^n))$: Length of the bit string X^n was mapped into

$$\begin{aligned} E[\mathcal{L}(\xi(X^n))] &= \sum_{X^n \in \mathcal{X}^n} \mathcal{L}(\xi(X^n))P(X^n) \\ &= \sum_{X^n \in \mathcal{A}_\epsilon^{(n)}} \mathcal{L}(\xi(X^n))P(X^n) + \sum_{X^n \notin \mathcal{A}_\epsilon^{(n)}} \mathcal{L}(\xi(X^n))P(X^n) \\ &\leq \sum_{X^n \in \mathcal{A}_\epsilon^{(n)}} (n(H(X) + \epsilon) + 1 + 1)P(X^n) \\ &\quad + \sum_{X^n \notin \mathcal{A}_\epsilon^{(n)}} (\log |X^n| + 1 + 1)P(X^n) \\ &\leq n(H(X) + \epsilon) + 2 + (\log |X^n| + 2)\epsilon \end{aligned}$$

- Normalize... $\frac{E[\mathcal{L}(\xi(X^n))]}{n} \leq H(X) + \epsilon + \frac{2}{n} + \epsilon(\log |X| + \frac{2}{n})$
- $\frac{E[\mathcal{L}(\xi(X^n))]}{n}$ is arbitrarily close to entropy (but scheme not scalable).

3 Compression of correlated sequences

| 31

- ▶ What if $\{X_i\}$ no longer iid ... identically distributed though

$$\text{e.g. } X_i \in \mathcal{X} = \begin{cases} r & 0.1, \\ s & 0.9 \end{cases}$$

- ▶ But X_i and X_{i+1} correlated.

$$H(X_i) = H(X_j) = H(X), \quad \forall i \neq j$$

- ▶ Definition: **Entropy rate**

$$\begin{aligned} > D_1: H(\mathcal{X}) = \lim_{n \rightarrow \infty} \frac{1}{n} H(X_1, X_2, \dots, X_n) = \lim_{n \rightarrow \infty} \frac{1}{n} H(X^n) \\ > D_2: H(\mathcal{X}) = \lim_{n \rightarrow \infty} H(X_n | X_1, X_2, \dots, X_{n-1}) \end{aligned}$$

- ▶ $H(\mathcal{X})$ sets the limit of compression of correlated sources
- ▶ In general $H(\mathcal{X}) \leq H(X)$ (equality when $\{X_i\}$ iid)
 - > Compressing correlated sources takes less space than iid data

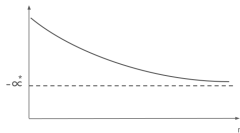
3 Compression of correlated sequences

| 32

- Note: D_1 same as D_2 (for stationary random process)

Proof:

- $H(X_{n+1}|X_1, X_2, \dots, X_n) =: \alpha_{n+1}$
- $H(\mathcal{X}) \stackrel{(D_2)}{=} \lim_{n \rightarrow \infty} \alpha_n$
- $\alpha_{n+1} = H(X_{n+1}|X_1, X_2, \dots, X_n) \stackrel{(cr)}{\leq} H(X_{n+1}|X_2, X_3, \dots, X_n) = H(X_n|X_1, X_2, \dots, X_{n-1}) = \alpha_n$ (stationary)
- $\alpha_{n+1} \leq \alpha_n, \Rightarrow \alpha_n \rightarrow \alpha^*$ when $n \rightarrow \infty$

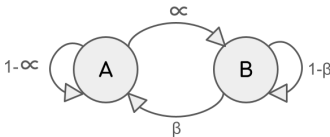


- $H(\mathcal{X}) \stackrel{(D_1)}{=} \frac{1}{n} H(X_1, X_2, \dots, X_n) \stackrel{(cr)}{=} \frac{1}{n} \sum_{i=1}^n H(X_i|X_1, X_2, \dots, X_{i-1}) = \frac{1}{n} \sum_{i=1}^n \alpha_i \xrightarrow{n \rightarrow \infty} \alpha^*$ (because sequence converges)
- $\Rightarrow D_1 = D_2$

3 Application of entropy rate: Markov chains

| 33

- Evaluate entropy rate for Binary Markov chains (BMC):



$$X_i = 0 (A) \text{ or } 1 (B), \quad P(X_i) \stackrel{(stnr)}{=} P(X_{i-1}), \quad P_0 = \frac{\beta}{\alpha + \beta}, \quad P_1 = \frac{\alpha}{\alpha + \beta} \quad (**)$$

- Entropy rate:

$$H(\mathcal{X}) \stackrel{(D_2)}{=} \lim_{n \rightarrow \infty} H(X_n | X_1, X_2, \dots, X_{n-1}) \stackrel{(mc)}{=} H(X_n | X_{n-1})$$

- $H(X_2 | X_1) = \dots = H(\alpha)P_0 + H(\beta)P_1$ (binary entropy)

- $H(\mathcal{X}) \stackrel{(**)}{=} \frac{H(\beta)\alpha + H(\alpha)\beta}{\alpha + \beta} \stackrel{(cre)}{<} H(X_2) = H\left(\frac{\alpha}{\alpha + \beta}\right) \stackrel{(sym)}{=} H\left(\frac{\beta}{\alpha + \beta}\right)$



- Equality for $\alpha = \beta = 0.5$, no correlation anymore.

> no gain, no difference between entropy rate and entropy.

- ① Introduction
- ② Definitions and Inequalities
- ③ Fundamental Limits of Compression
- ④ **Communication and Channel Capacity**
Channel Capacity for the BSC
Encoding and Decoding
- ⑤ Shannon's Channel Coding Theorem
- ⑥ Information Theory of Continuous Data

 $X \in \mathcal{X}$ finite $Y \in \mathcal{Y}$ finite

- ▶ Channel given by probability transition matrix:

$$\text{Prob}(Y_n = y | X_n = x), \forall x \in \mathcal{X}, y \in \mathcal{Y}$$

- ▶ Memoryless channel:

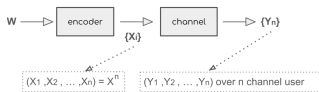
- > Y_n depends on X_n only, not on X_{n-1}, X_{n-2}, \dots
- > Assume stationarity: the index n does not matter

$$\text{Prob}(Y = y | X = x), \forall x \in \mathcal{X}, y \in \mathcal{Y}$$

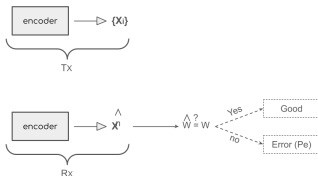
4 Communication and Channel Capacity

| 36

- Fundamental question: what is max reliable communication rate?



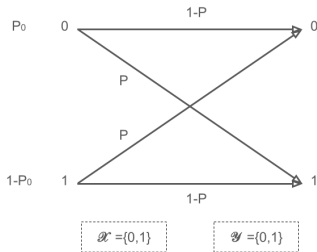
- $w \in \mathcal{W} = \{w^1, \dots, w^M\}$: message set
- $\{X_i\}$: stream of bits in which the message is encoded
- Rate = “number of information bits sent per channel use”
- $\log_2 M$: information bits per message
 - > as $M \uparrow$, messages too close, harder for RX to distinguish $\Rightarrow \hat{X}_i \neq X_i$
- Capacity C is maximal rate R such that $P_e \rightarrow 0$ as $n \rightarrow \infty$.



4 Channel Capacity for the BSC

| 37

- ▶ Let's explore the Binary Symmetric Channel (BSC)



- ▶ We will use (and prove later) that

$$C = \max_{P(X)} I(X; Y), \quad P(X) \text{ input distribution}$$

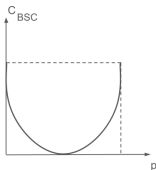
4 Channel Capacity for the BSC

| 38

- ▶ Start with $C = \max_{P(X)} I(X; Y)$
- ▶ $I(X; Y) = H(Y) - H(Y|X) \quad (*)$

$$\begin{aligned} H(Y|X) &= \sum_{x \in \mathcal{X}} H(Y|X=x)P(x) = H(Y|X=0)p_0 + (1-p_0)H(Y|X=1) \\ &= H_b(p)p_0 + (1-p_0)H_b(p) = H_b(p) \quad (**) \end{aligned}$$

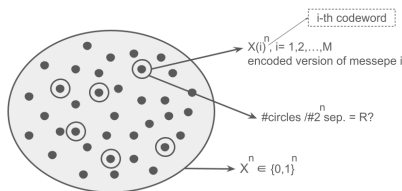
- ▶ $I(X; Y) \stackrel{(*),(**)}{=} H(Y) - H_b(p) \leq 1 - H_b(p)$
 - > $H(Y) \leq 1$ Maxed if $p_0 = 0.5$
- ▶ $C_{BSC} = \max_{P(X)} I(X; Y) = 1 - H_b(p)$
 - > achieved when $P(X=0) = P(X=1) = \frac{1}{2}$ which yields maximizing $P(Y=0) = P(Y=1) = \frac{1}{2}$



4 Encoding and Decoding

| 39

- ▶ Each message $i = \{1, \dots, M\}$, maps onto long vector $X^n(i) = \{0, 1\}^n$ from Binary Code \mathcal{C}



- > Recall: Entropy rate (affects goodness of code)

$$H = \lim_{n \rightarrow \infty} H(X_n | X_1, X_2, \dots, X_{n-1}) \stackrel{(if \text{ } mc)}{=} H(X_n | X_{n-1})$$

- ▶ Transmit vector $X^n(i)$, corrupted by channel, to yield Y^n
- ▶ Want to decode: $g(Y^n) = i \in \{1, 2, \dots, M\}$

4 Decoding Errors

| 40

► $\mathcal{C} = \text{codebook} = \{X^n(1), X^n(2), \dots, X^n(M)\}$

► Define probability error:

> Def. $\lambda_i = \text{Prob}(g(Y^n) \neq i \mid X^n = X^n(i))$

> Def. $\lambda^n = \max_{i=1,2,\dots,M} \lambda_i$ (worst case)

> Def. $P_e = \frac{1}{M} \sum_{i=1}^M \lambda_i$ (average case)

► We want

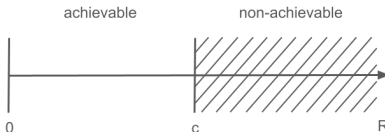
$$\lambda^n \xrightarrow{n \rightarrow \infty} 0$$

- ① Introduction
- ② Definitions and Inequalities
- ③ Fundamental Limits of Compression
- ④ Communication and Channel Capacity
- ⑤ Shannon's Channel Coding Theorem
- ⑥ Information Theory of Continuous Data
- ⑦ The Gaussian Channel

5 Shannon's Channel Coding Theorem

| 42

- ▶ A rate R is achievable iff there exists a series of (M, n) codes such that $\lambda^n \xrightarrow{n \rightarrow \infty} 0$



- ▶ **Theorem:** the maximum of all achievable rates R is given by

$$C = \max_{P(X)} I(X; Y)$$

5 Shannon's Channel Coding Theorem

| 43

Proof:

- ▶ Let us build a random code
- ▶ We generate M code words $X^n(1), X^n(2), \dots, X^n(M)$.
 - > According to $P(X^n) = \prod_{i=1}^n P(X_i)$, $X^n = (X_1, X_2, \dots, X_n)$
- ▶ Assume that message i is transmitted ($X^n(i)$ is transmitted)
- ▶ Decoder receives Y^n
- ▶ Decoder finds message i s.t. $(X^n(i), Y^n)$ is jointly typical
 - > Let's see what joint typicality (J.T.) means

- ▶ Def. Sequences X^n and Y^n are jointly typical if:
 - ▶ 1. X^n is typical³ with accuracy ϵ , i.e., if $|\frac{1}{n} \log P(X^n) - H(X)| < \epsilon$
 - > Recall $P(X^n) \approx 2^{-nH(X)}$
 - > Recall $|\mathcal{A}_{X,\epsilon}| \approx 2^{nH(X)}$
 - ▶ 2. Y^n is typical with accuracy ϵ , i.e., if $|\frac{1}{n} \log P(Y^n) - H(Y)| < \epsilon$
 - > Recall $P(Y^n) \approx 2^{-nH(Y)}$
 - > Recall $|\mathcal{A}_{Y,\epsilon}| \approx 2^{nH(Y)}$
 - ▶ 3. (X^n, Y^n) has to be such that $|\frac{1}{n} \log P(X^n, Y^n) - H(X, Y)| < \epsilon$
 - > $P(X^n, Y^n) \approx 2^{-nH(X,Y)}$
 - > $|\mathcal{A}_\epsilon| \approx 2^{nH(X,Y)}$

³Make histogram, consider as pdf, calculate entropy.

Key property of joint-typicality:

- ▶ Recall: channel represented by $P(Y|X)$ i.e., by $P(X, Y)$
- ▶ Let (X^n, Y^n) , drawn from $\prod_i^n P(X_i, Y_i)$ (i.e., channel related)
 - > Y^n has been produced after sending X^n through channel
 - > $\prod_i^n P(X_i, Y_i)$ since X_i iid and channel memoryless.

Theorem

- ▶ 1. $\text{Prob}((X^n, Y^n) \in \mathcal{A}_\epsilon) \xrightarrow{n \rightarrow \infty} 1$
 - > From LLN (as before - convergence of sequence of logs of joint distr.)
- ▶ 2. Let $\tilde{X}^n, \tilde{Y}^n \sim P(X^n)P(Y^n)$
 - > i.e. \tilde{X}^n, \tilde{Y}^n are independent
 - > Observing \tilde{Y}^n . Going through all X^n . If \tilde{X}^n not the tx, then independent to \tilde{Y}^n . ($\tilde{X}^n, \tilde{Y}^n \sim P(X^n)P(Y^n)$)
 - > similarly: \tilde{Y}^n is not "channel associated" to transmitted \tilde{X}^n
 - > Then:

$$\text{Prob}((\tilde{X}^n, \tilde{Y}^n) \in \mathcal{A}_\epsilon) \approx 2^{-n(I(X, Y))} \quad (1)$$

► Sketch Proof of (1):

- > How many jointly typical pairs come from distribution $P(X^n)P(Y^n)$

$$P((\tilde{X}^n, \tilde{Y}^n) \in \mathcal{A}_\epsilon) = \sum_{(x^n, y^n) \in \mathcal{A}_\epsilon^n} P(x^n)P(y^n)$$

- > Number of summands is $|\mathcal{A}_\epsilon^n| \approx 2^{nH(X,Y)}$
> Pairs are jointly typical so $P(X^n) = 2^{-nH(X)}$, $P(Y^n) = 2^{-nH(Y)}$

$$P((\tilde{X}^n, \tilde{Y}^n) \in \mathcal{A}_\epsilon) \leq \frac{2^{nH(X,Y)}}{2^{nH(X)}2^{nH(Y)}} = 2^{-nI(X;Y)}$$

► Example:

$$\begin{aligned} |\mathcal{A}_{\epsilon,X}| &= |\mathcal{A}_{\epsilon,Y}| = 1000, \quad |\mathcal{A}_\epsilon| = 300 \\ \Rightarrow P((\tilde{X}^n, \tilde{Y}^n) \in \mathcal{A}_\epsilon) &\leq \frac{300}{10^6} = 3 \cdot 10^{-4} \end{aligned}$$

5 Joint Typicality and Channel Capacity

| 47

- ▶ ...Continue toward proving $C = \max_{P(X)} I(X; Y)$
- ▶ Consider random code \mathcal{C}
- ▶ Let us assume **message** $i = 1$ has been sent via $X^{(1)} \in \mathcal{C}$
- ▶ Let us assume that Y^n **was received**
- ▶ Error events : $E_i = \begin{cases} 1 & \text{If } Y^n \text{ is J.T. with } X^{(i)}, \\ 0 & \text{If not} \end{cases}$
- ▶ An error occurs if $E_1 = 0$ or $E_2 = 1$ or ... $E_M = 1$
- ▶ Thus

$$\begin{aligned} P_e &\stackrel{(\text{unionB})}{\leq} \text{Prob}(E_1 = 0) + \sum_{i=2}^M \text{Prob}(E_i = 1) \\ &\leq \epsilon + (2^{nR} - 1)2^{-n(I(X,Y)-3\epsilon)} \end{aligned}$$

- > Since $\text{Prob}(E_1 = 0) \approx \epsilon \rightarrow 0$ from Theorem (part 1)
- > Since $M = 2^{nR}$
- > Since $\text{Prob}(E_i = 1) \approx 2^{-nI(X,Y)}$ from Theorem (part 2)

- ▶ Have seen that

$$P_e \leq \epsilon + (2^{nR} - 1)2^{-n(I(X,Y)-3\epsilon)}$$

- ▶ Thus

$$P_e \leq 3\epsilon \rightarrow 0, \quad n \gg 1, \quad R < I(X, Y).$$

- ▶ Thus $R < I(X, Y)$

- ▶ Eventually ...

$$C = \max_{P(X)} I(X; Y)$$

- ① Introduction
- ② Definitions and Inequalities
- ③ Fundamental Limits of Compression
- ④ Communication and Channel Capacity
- ⑤ Shannon's Channel Coding Theorem
- ⑥ Information Theory of Continuous Data
Typical Sequences
- ⑦ The Gaussian Channel

6 Information Theory of Continuous Data

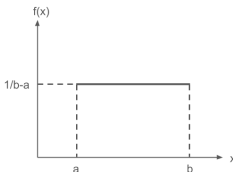
| 50

- ▶ X is continuous R.V with distribution $f(x)$ (PDF)
- ▶ $F(x) = \text{Prob}(X \leq x)$ CDF $\rightarrow f(x) = \frac{dF(x)}{dx}$
- ▶ $\int_{-\infty}^{+\infty} f(x) dx = 1$

Definition: Differential entropy $X \sim f(X)$

$$h(x) \triangleq -E[\log f(x)] = \int_{-\infty}^{+\infty} -f(x) \log f(x) dx = \int_{\text{Supp}(X)} -f(x) \log f(x) dx$$

- ▶ $X \sim U[a, b]$
 $a, b \in \mathbb{R}, a \leq b \Rightarrow h(x) = -E[\log f(x)] = -E[\log \frac{1}{b-a}] = \log(b-a)$



Example: Gaussian pdf (maximizes Entropy)

► $X \sim N(0, \sigma^2)$, $f(X) = \frac{1}{\sqrt{2\pi\sigma^2}} e^{-\frac{X^2}{2\sigma^2}}$

$$\begin{aligned} h(X) &= - \int_{-\infty}^{+\infty} f(X) \log f(X) = - \int_{-\infty}^{+\infty} \frac{1}{\sqrt{2\pi\sigma^2}} e^{-\frac{X^2}{2\sigma^2}} \log_2 \frac{1}{\sqrt{2\pi\sigma^2}} e^{-\frac{X^2}{2\sigma^2}} dx \\ &= - \log \frac{1}{\sqrt{2\pi\sigma^2}} - \int_{-\infty}^{+\infty} -\frac{X^2 \log_2 e}{2\sigma^2} f(x) dx \\ &= - \log \frac{1}{\sqrt{2\pi\sigma^2}} + \frac{\log_2 e}{2\sigma} \int_{-\infty}^{+\infty} X^2 f(X) dx \end{aligned}$$

► Since $f(X)$ and X^2 are symmetric, then second term goes to zero

$$\Rightarrow h(X) = \log \sqrt{2\pi e \sigma^2}$$

> The higher the σ^2 , the more disorder we have, the higher $h(X)$

- ▶ 1) Central limit theorem: X_1, \dots, X_n i.i.d arbitrary

- > $\bar{X}_n = \frac{1}{n} \sum_{i=1}^n X_n$ empirical average
 $\Rightarrow \sqrt{n}\bar{X}_n \rightarrow \bar{X} \sim N(0, \sigma^2)$

- ▶ 2) Uncorrelated vs. independent

- > $E[XY] = 0 \Leftarrow f(X,Y)=f(X)f(Y)$
 \nRightarrow

- ▶ 3) Sum of independent Gaussians

- > X, Y jointly Gaussian $\Rightarrow Z = X \pm Y \sim N(\mu_X \pm \mu_Y, \sigma_X^2 + \sigma_Y^2)$

- ▶ 4) Shift and scale

$$Z \sim N(0, 1) \Rightarrow X = aZ + b \sim N(b, a^2) \quad a, b \text{ constants}$$

- ▶ 5) Multivariate Gaussian $\underline{X} = \{X_1, X_2, \dots, X_n\} \sim N(\underline{\mu}, \Sigma)$

- > with $\underline{\mu} = E[\underline{X}], \Sigma = E[(\underline{X} - \underline{\mu})(\underline{X} - \underline{\mu})^T]$

$$f(\underline{X}) = \frac{1}{\sqrt[n]{2\pi|\Sigma|}} e^{-\underline{X}^T \Sigma^{-1} \underline{X}}$$

Def: Joint and Conditional Entropy

- ▶ Def: $h(X, Y) \triangleq -E[\log f(X, Y)] = -\int_{-\infty}^{+\infty} f(X, Y) \log f(X, Y) dx dy$
- ▶ Def: $h(X|Y) = -\int_{-\infty}^{+\infty} f(X, Y) \log f(X|Y) dx dy$

Def: KL-divergence and Mutual Information

- ▶ Def: $D(f(x)||g(x)) = \int_{-\infty}^{+\infty} f(X) \log \frac{f(X)}{g(X)} dx$
- ▶ Def: $I(X, Y) = D(f(X, Y)||f(X)f(Y))$

Also similar to before

$$I(X, Y) = h(X) - h(X|Y) = h(Y) - h(Y|X)$$

$$h(Y|X) \leq h(Y), \quad h(X|Y) \leq h(X)$$

► **Chain Rule**

$$\begin{aligned}h(X_1, \dots, X_n) &= \sum_{i=1}^n h(X_i | X_1, \dots, X_{i-1}) \\&= \sum_{i=1}^n h(X_i | X_{i+1}, \dots, X_n) \leq \sum_{i=1}^n h(X_i)\end{aligned}$$

- > Equality for X_i independent
- > Remember $h(X, Y) = h(X) + h(Y|X)$.

► **Scale and Shift** $X \sim f(X)$

- > $Y = X + C$ (C fixed) $\Rightarrow h(Y) = h(X)$
- > $Y = aX$ (a fixed) $\Rightarrow h(Y) = h(X) + \log_2 a$ (just plug in def.)

6 Maximization of Differential Entropy

| 55

- ▶ Recall that $h(X + C) = h(X) \rightarrow$ we focus w.l.o.g. on $E[X] = 0$
- ▶ Focus on $\text{var}(x) = \sigma^2$ (fixed) otherwise ill-posed problem
- ▶ $\max h(x) = ?$
- ▶ Proof as before but now let $g(x) = \frac{1}{\sqrt{2\pi\sigma^2}} e^{-\frac{x^2}{2\sigma^2}}$ (Gaussian)
- ▶ As before use that $D(f(X)||g(X)) > 0$ to get

$$h(X) \leq \frac{1}{2} \log(2\pi e\sigma^2)$$

- > with equality for $f(x) = g(x)$ (gaussian)
- > Whereas recall: For discrete case $H(X) \leq \log_2 |\mathcal{X}|$ (uniform)

6 ...Typical Sequences (similar to discrete, BUT....)

| 56

- ▶ $X^n = \{X_1, \dots, X_n\}$ iid
- ▶ $\frac{1}{n} \sum_{i=1}^n \log_2 f(X_i) \xrightarrow{n \rightarrow +\infty} E[\log_2 f(X_i)] = -h(X)$
- ▶ $\Delta_\epsilon^{(n)} = \{X^n \in \mathbb{R}^n \mid | -\frac{1}{n} \sum_{i=1}^n \log(f(X_i)) - h(X) | < \epsilon\}$

- ▶ Prop. 1

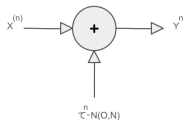
$$P(X^n \in \Delta_\epsilon^{(n)}) \stackrel{(LLN)}{\geq} 1 - \epsilon, \text{ as } n \rightarrow \infty$$

- ▶ Prop. 2

$$\text{Vol}(\Delta_\epsilon^{(n)}) \approx 2^{n(h(X))}$$

- > Where volume of set $S \subset \mathbb{R}^n$ is $\text{Vol}(S) = \int_S dx^n$

- 1 Introduction
- 2 Definitions and Inequalities
- 3 Fundamental Limits of Compression
- 4 Communication and Channel Capacity
- 5 Shannon's Channel Coding Theorem
- 6 Information Theory of Continuous Data
- 7 The Gaussian Channel
Capacity of Gaussian channel



- ▶ Why Gaussian?
 - > Mathematical Tractability and CLT!
 - > Useful for the worst case!
 - > AWGN Z_i a good assumption
- ▶ $w \in \mathcal{W} = \{1, 2, 3, \dots, M\}$ message, which is encoded to $X^n(w) = (X_1, X_2, \dots, X_n)$
- ▶ AWGN Channel

$$Y_i = X_i + Z_i \quad \forall i = \{1, \dots, n\}, \text{Rate } R = \frac{\log_2 M}{n}$$

Q: what is R_{\max} s.t. $p(\hat{w} = w) \xrightarrow{n \rightarrow \infty} 1$

Looking for

$$C = \max_{P(X)} I(X; Y)$$

- ▶ Under power constraint $\begin{cases} E[X_i^2] \leq P & \forall i \text{ (instantaneous)} \\ \frac{1}{n} \sum_{i=1}^n X_i^2 \leq P & \text{average} \end{cases}$
- ▶ Let us compute:

$$\begin{aligned} I(X; Y) &= h(Y) - h(Y|X) \\ &= h(Y) - h(X + Z|X) = h(Y) - h(Z) \\ &= h(Y) - \frac{1}{2} \log(2\pi eN) \leq \frac{1}{2} \log(2\pi e\sigma_y^2) - \frac{1}{2} \log(2\pi eN) \end{aligned}$$

- ▶ $h(y)$ is always upper bounded by Gaussian entropy
- ▶ If $X \sim N(0, P)$ then Y Gaussian with $\sigma_y^2 = \sigma_x^2 + \sigma_z^2 \leq P + N$

$$I(X; Y) \leq \frac{1}{2} \log(2\pi e(P + N)) - \frac{1}{2} \log(2\pi eN)$$

- ▶ Thus $C = \frac{1}{2} \log(1 + \frac{P}{N})$ bits/sec/Hz or bits/channel use

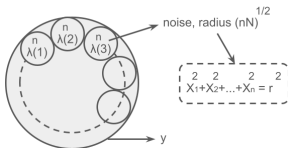
7 Capacity of Gaussian channel – Intuition via Sphere Packing | 60

- ▶ $\mathcal{W} = \{1, 2, \dots, M\} \xrightarrow{\text{encoder}} \mathcal{C} = \{X^n(1), X^n(2), \dots, X^n(M)\}$ codebook
- ▶ $Y^n = X^n(w) + Z^n$ and

$$\sum_{i=1}^n X_i^n(w)^2 = nP \quad (\text{power constraint } \forall w)$$

- ▶ When n large

$$\frac{1}{n} \sum_{i=1}^n Z_i^2 \approx N \quad \frac{1}{n} \sum_{i=1}^n Y_i^2 \approx N + P$$



- ▶ Small Hypersphere of radius $\sum_{i=1}^n Z_i^2 \approx nN$
- ▶ Big Hypersphere of radius $\sum_{i=1}^n Y_i^2 \approx n(N + P)$
- ▶ Decoder works if small bubbles do not touch (and fails otherwise)

7 Capacity of gaussian channel – Sphere packing argument

| 61

- ▶ volume of n -dimensional hyper-sphere is $vol(n) = \alpha_n r^n$
 - > volume of hyper-sphere of radius $\sqrt{n(N+P)}$ is $\alpha_n(n(N+P))^{\frac{n}{2}}$
 - > volume of "noisy bubble" is $\alpha_n(nN)^{\frac{n}{2}}$
- ▶ M (number of bubbles/messages) can not be bigger than

$$M \leq \frac{\alpha_n(n(N+P))^{\frac{n}{2}}}{\alpha_n(nN)^{\frac{n}{2}}} = \left(1 + \frac{P}{N}\right)^{\frac{n}{2}}$$

- ▶ This corresponds to $\log\left(\left(1 + \frac{P}{N}\right)^{\frac{n}{2}}\right)$ bits, over n dimensions
- ▶ Thus

$$R \leq \frac{1}{2} \log\left(1 + \frac{P}{N}\right) \quad \text{per real dimension}$$