

# Information Theory Recap Session

Changyeol Lee (Yonsei University)

# Entropy

$X$ : a discrete random variable over  $\mathcal{X}$  with the PMF(probability mass function)  $p(\cdot)$ .

The **entropy** of  $X$ : a measure of the uncertainty of  $X$

$$H(X) = \sum_{x \in \mathcal{X}} p(x) \log \frac{1}{p(x)} = \mathbb{E}_{X \sim p} \left[ \log \frac{1}{p(X)} \right]$$

**Fact.**  $H(X) \geq 0$ .

## Conditional Entropy

$$H(Y|X) = \sum_{x \in \mathcal{X}} p(x) H(Y|X = x) = \sum_{x \in \mathcal{X}, y \in \mathcal{Y}} p(x, y) \log \frac{1}{p(y|x)} = \mathbb{E}_{(X,Y) \sim p} \left[ \log \frac{1}{p(Y|X)} \right]$$

## Chain Rule

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

# Kullback-Leibler Divergence (Relative Entropy)

**Kullback-Leibler divergence** between PMFs  $p$  and  $q$

$$D(p \parallel q) = \sum_{x \in \mathcal{X}} p(x) \log \frac{p(x)}{q(x)} = \mathbb{E}_{X \sim p} \left[ \log \frac{p(X)}{q(X)} \right] = \mathbb{E}_{X \sim p} \left[ \log \frac{1}{q(X)} \right] - \mathbb{E}_{X \sim p} \left[ \log \frac{1}{p(X)} \right]$$

\*  $D(p \parallel q) = \infty$  if  $\exists x \in \mathcal{X}$  s.t.  $p(x) > 0$  and  $q(x) = 0$ .

\*  $D(p \parallel q) \neq D(q \parallel p)$ , i.e., no symmetricity in general

\*  $D(p \parallel q) + D(q \parallel r) \not\cong D(p \parallel r)$  and  $D(p \parallel q) + D(q \parallel r) \not\leq D(p \parallel r)$  in general

Chain rule

$$D(p(x, y) \parallel q(x, y)) = D(p(x) \parallel q(x)) + D(p(y|x) \parallel q(y|x))$$

# Mutual Information

## Mutual information

- a measure of the amount of information that one RV contains about another RV

$$\begin{aligned} I(X; Y) &= \mathbb{E}_{(X,Y) \sim p} \left[ \log \frac{p(X, Y)}{p(X)p(Y)} \right] \\ &= D(p(x, y) \parallel p(x)p(y)) \\ &= H(X) - H(X|Y) = H(Y) - H(Y|X) \end{aligned}$$

the reduction in the uncertainty of  $X$  ( $Y$ )  
due to the knowledge of  $Y$  ( $X$ )

## Conditional Mutual Information

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

the mutual information of  $X_1$  and  $X_2$ , given  $X_3$ ;  
not the mutual information of  $X_1$  and  $X_2|X_3$ .

\*  $I(X; Y|Z) \not\leq I(X; Y)$  and  $I(X; Y|Z) \not\geq I(X; Y)$  in general

## Chain Rule

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

# Information Inequality

**Theorem.**  $D(p \parallel q) \geq 0$  with equality iff  $p = q$ .

**Corollary.**  $I(X; Y) \geq 0$  with equality iff  $X$  and  $Y$  are independent.

**Corollary.**  $H(X|Y) \leq H(X)$ , i.e., *conditioning only reduces entropy*.

**Corollary.**  $H(X) \leq \log|\mathcal{X}|$  with equality iff  $p$  is the uniform distribution.

# Data-processing Inequality

By chain rule,

$$I(X; Z) + I(X; Y|Z) = I(\mathbf{X}; \mathbf{Y}, \mathbf{Z}) = I(X; Y) + I(X; Z|Y)$$

**Theorem.** If  $X \rightarrow Y \rightarrow Z$ , then  $I(X; Y) \geq I(X; Z)$ .

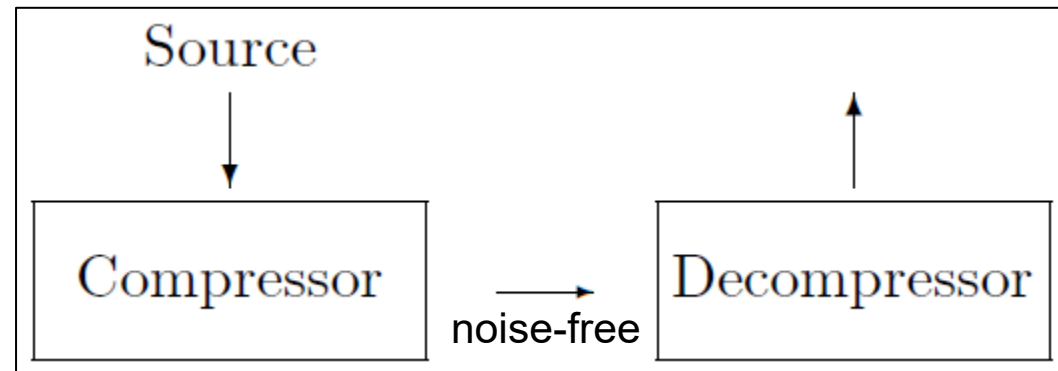
If  $X \rightarrow Y \rightarrow Z$ ,  
by definition,  
 $I(X; Z|Y) = 0$ .

**Theorem.** If  $X \rightarrow Y \rightarrow Z$ , then  $I(X; Y) \geq I(X; Y|Z)$ .

# Source Coding Theorem

We have  $n$  i.i.d. RVs.

What is the min #bits required to send the data only with negligible error?



# Upper/Lower bound on Compression

We can construct an instantaneous code given a length function  $\ell: \mathcal{X} \rightarrow \{0,1\}^*$  if

$$\sum_{x \in \mathcal{X}} 2^{-\ell(x)} \leq 1.$$

Shannon compression  $\ell(x) := \lceil -\log p(x) \rceil$  gives  $\mathbb{E}_{X \sim p}[\ell(X)] < H(X) + 1$ .

**Theorem.** Huffman compression is optimal, i.e.,  $\mathbb{E}[\ell_{\text{Huffman}}(X)] \leq \mathbb{E}[\ell_{\text{Uniquely Decodable Code}}(X)]$ .

---

Any uniquely decodable code with a length function  $\ell: \mathcal{X} \rightarrow \{0,1\}^*$  must satisfies

$$\sum_{x \in \mathcal{X}} 2^{-\ell(x)} \leq 1.$$

**Theorem.**  $\mathbb{E}[\ell_{\text{Uniquely Decodable Code}}(X)] \geq H(X)$  with equality iff  $p(x) = 2^{-\ell(x)}$  for all  $x \in \mathcal{X}$ .

---

$$H(X) \leq \mathbb{E}[\ell^*(X)] < H(X) + 1$$

# Upper/Lower bound on Compression

Consider a sequence of (possibly dependent) RVs  $X_1, X_2, \dots, X_n$  with joint distribution  $p$ .

Shannon compression gives  $\mathbb{E}_{(X_1, X_2, \dots, X_n) \sim p}[\ell(X_1, X_2, \dots, X_n)] < H(X_1, X_2, \dots, X_n) + 1$ .

If  $X_1, X_2, \dots, X_n$  are i.i.d.,  $H(X_1, X_2, \dots, X_n) = nH(X)$ .

Therefore, the expected length per symbol of an optimal compression is

$$H(X) \leq \frac{1}{n} \mathbb{E}[\ell^*(X_1, X_2, \dots, X_n)] < H(X) + \frac{1}{n}.$$

**$H(X)$  is the fundamental limit!**

Q. What is the fundamental limit if we allow small error in the compression scheme?



# AEP (Asymptotic Equipartition Property)

Consider a sequence of i.i.d. RVs  $X_1, X_2, \dots, X_n$ .

For any  $\epsilon > 0$ , for all sufficiently large  $n$ ,

$$\Pr[2^{-n(H(X)+\epsilon)} < p(X_1, X_2, \dots, X_n) < 2^{-n(H(X)-\epsilon)}] \geq 1 - \epsilon.$$

The **typical set**  $A_\epsilon^{(n)}$  w.r.t.  $p$  is the set of sequence  $\mathbf{x} = (x_1, x_2, \dots, x_n) \in \mathcal{X}^n$  such that

$$2^{-n(H(X)+\epsilon)} < p(\mathbf{x}) < 2^{-n(H(X)-\epsilon)}.$$

**AEP.** For sufficiently large  $n$ ,

$$\Pr[\mathbf{X} \in A_\epsilon^{(n)}] \geq 1 - \epsilon \quad \text{contains most of the probability}$$

$$(1 - \epsilon)2^{n(H(X)-\epsilon)} \leq |A_\epsilon^{(n)}| \leq 2^{n(H(X)+\epsilon)}. \quad \text{cardinality} \approx 2^{nH(X)}$$

**$A_\epsilon^{(n)}$  contains most of the probability and has a cardinality  $\approx 2^{nH(X)}$ .**

# Typical Set is a Smallest Set

Let  $B_\delta^{(n)} \subseteq \mathcal{X}^n$  be a smallest set with  $\Pr[\mathbf{X} \in B_\delta^{(n)}] \geq 1 - \delta$ .

**Lemma.**  $|B_\delta^{(n)}| \geq (1 - \epsilon - \delta)2^{n(H(X) - \epsilon)} \approx 2^{nH(X)}$ .

$A_\epsilon^{(n)}$  contains most of the probability and *essentially* has a smallest cardinality of  $\approx 2^{nH(X)}$ .

## Source coding theorem.

Consider a sequence of  $n$  i.i.d. RVs with entropy  $H$ .

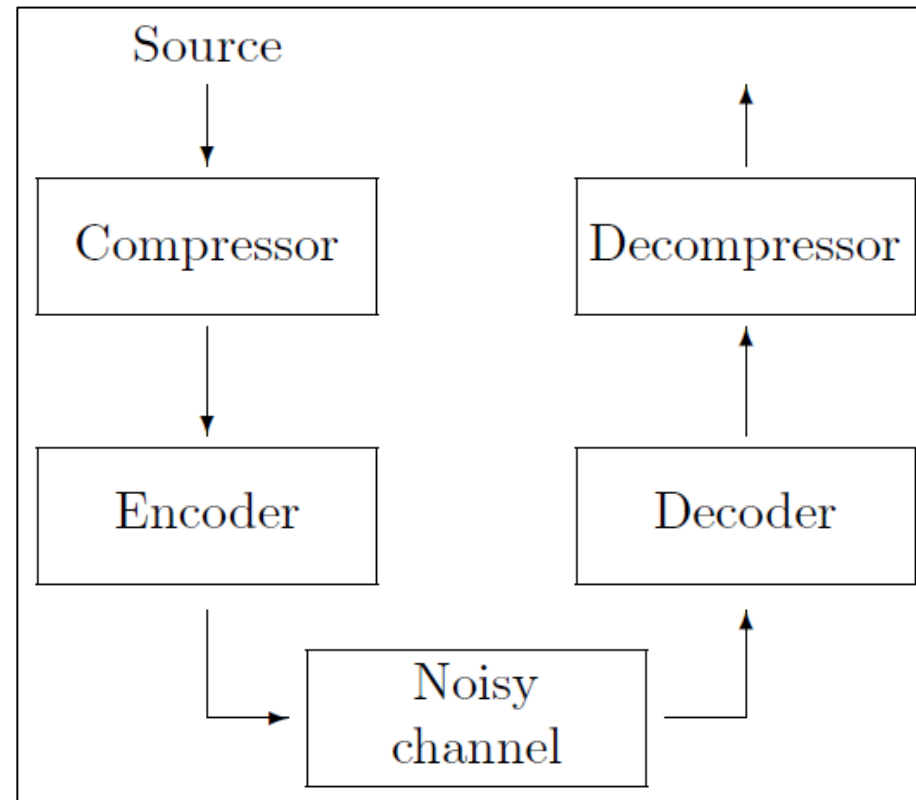
(1-1) Using *slightly more than*  $nH$  bits admits a lossless compression.

(1-2) With #bits *very close to*  $nH$ ,  $\Pr[\text{error}] \approx 0$  for sufficiently large  $n$

(2) With less than  $nH$  bits,  $\Pr[\text{error}] \approx 1$  for sufficiently large  $n$  (i.e., all information is lost).

# Channel Coding Theorem

Suppose we use  $n$  symbols to encode a date to cope with the channel noise. What is the max #data we can send only with negligible error?



# Motivation

We have a noisy channel  $p(y|x)$ .

Alice tosses a coin  $X$  and send  $X$  to Bob (using single bit).

- amount of information before being sent =  $H(X)$

Bob receives a bit  $Y$  through a noisy channel.

- amount of information that channel decreased =  $H(X|Y)$
- amount of information conveyed by the channel =  $H(X) - H(X|Y) = I(X, Y)$

Information channel capacity

- Assume for now that we wish to maximize the amount of information conveyed by the channel.
- We do this by choosing a *best* distribution of  $X$ .

# Channel and Channel Capacity

## Discrete memoryless channel $Q$

$\mathcal{X}$ : an input alphabet (a set of input symbols)

$\mathcal{Y}$ : an output alphabet (a set of output symbols)

$\{p(y|x)\}_{x \in \mathcal{X}}$ : a collection of (conditional) PMFs

Choose a distribution.  $\rightarrow$  Capable of sending  $I(X; Y)$  amount of information

Channel capacity (choose a *best* distribution)

$$C = \max_{\text{distribution over } \mathcal{X}} I(X; Y) \quad * \text{ maximum is well defined}$$

# Code, Decode, Rate, Error

$M, n \in \mathbb{N}$

$(M, n)$  **code** for a channel  $Q = (\mathcal{X}, p(y|x), \mathcal{Y})$  “encoding scheme”

- a set of indices(data)  $\{1, \dots, M\}$
- a set of codewords  $\{\mathbf{x}^{(1)}, \dots, \mathbf{x}^{(M)}\}$ , where  $\mathbf{x}^{(j)} \in \mathcal{X}^n$

A **decoder**  $g$  guess an index.

- An optimal decoder  $g^*$  chooses a posteriori most likely index.

( $\log M$ )-bit

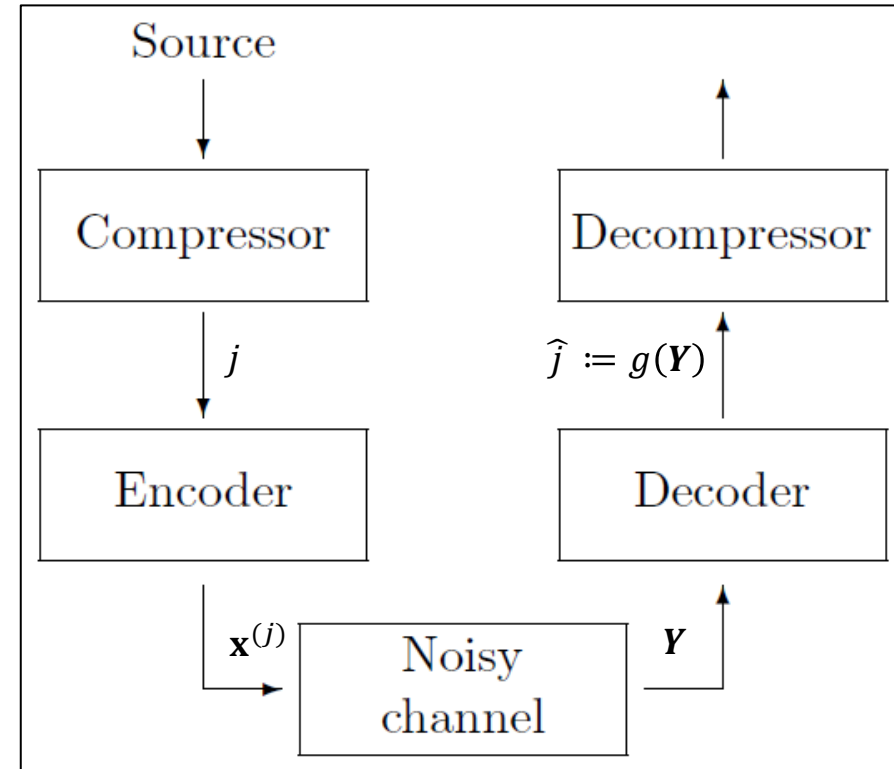
$(M, n)$  code sends an index with  $n$  transmissions.

**Rate**  $R$  of  $(M, n)$  code for  $Q$

$$R = \frac{\text{size of data}}{\text{\#transmission}} = \frac{\log M}{n}$$

**Maximal probability of error** (for a fixed channel  $Q$  and fixed  $(M, n)$  code for  $Q$ )

$$\lambda_{\max} := \max_{j \in \{1, \dots, M\}} \Pr_{Y \sim p(\cdot | \mathbf{x}^{(j)})} [g(Y) \neq j]$$



$$Y \sim p(\cdot | \mathbf{x}^{(j)}) = \prod_{i=1}^n p(y_i | x_i^{(j)})$$

# Achievable Rate

Let us fix  $n$ .

Increase rate = Send more information per transmission  
= Cause more error (possibly)

**Q.** What rate can we prove is *achievable*?

One way to show achievability = show existence of such code

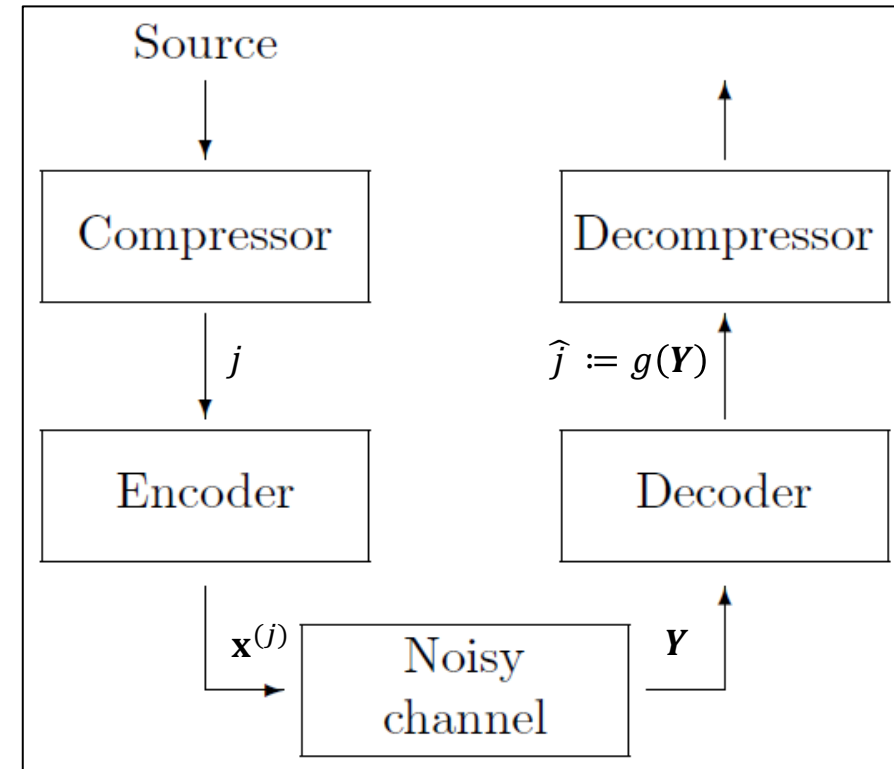
- Show the expected error of a random code is small.
- There must be a code with small error.

One (naïve) way to construct a random code

- Fix a distribution  $p_x$  over  $\mathcal{X}$ . Sample i.i.d.  $n$  symbols from  $p_x$  for each codeword, independently.

**A.**  $R \lesssim I$  is achievable. i.e., sending  $2^{nI}$  is possible

If  $p_x$  is a *best* distribution,  $R \lesssim C$  is achievable.



$$\mathbf{Y} \sim p(\cdot | \mathbf{x}^{(j)}) = \prod_{i=1}^n p(y_i | x_i^{(j)})$$

# Theorem (part 1)

The following holds for any discrete memoryless channel  $Q = (\mathcal{X}, p(y|x), \mathcal{Y})$ .

For any  $\epsilon > 0$  and  $R < C := \max_{p_x} I(X; Y)$ , there exists a  $(M := \lceil 2^{nR} \rceil, n)$  code for  $Q$  such that

$\lambda_{\max} < \epsilon$  for all sufficiently large  $n$ .

## Showing the existence

- Fix any  $p_x$ . Generate a random  $(M', n)$  code according to  $p_x$  where  $M' = \lceil 2^{n(R+1/n)} \rceil$ :
  - For each  $j \in [M']$ , independently,  $\mathbf{X}^{(j)} = X_1^{(j)} X_2^{(j)} \dots X_n^{(j)}$  where  $X_i^{(j)} \sim p_x$  independently.
- Sample  $J \in [M']$  at random. Consider  $Y \sim p(\cdot | \mathbf{X}^{(J)})$  and a *jointly typical* decoder  $g$ .  
Sample a codeword  $\mathbf{X}^{(J)}$  at random
- **Claim.**  $\Pr_{\substack{(M,n) \text{ code,} \\ J, Y \sim p(\cdot | \mathbf{X}^{(J)})}} [g(\mathbf{Y}) \neq J]$  is small.  $\rightarrow \exists$  a  $(M, n)$  code with small  $\Pr_{J, Y \sim p(\cdot | \mathbf{X}^{(J)})} [g(\mathbf{Y}) \neq J]$ .
- Removing the worst half of codewords ensures  $\lambda_{\max} = \max_{j \in [M]} \Pr_{Y \sim p(\cdot | \mathbf{x}^{(j)})} [g(\mathbf{Y}) \neq j]$  is also small.  
Rate decreases by  $1/n$ .



# Intuitive Idea

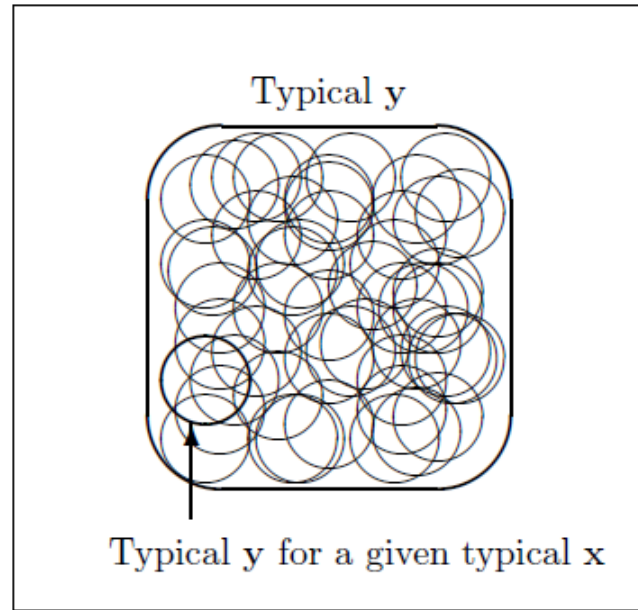
Consider sufficiently large  $n$ .

By AEP,

$$|\{\mathbf{y} \mid p(\mathbf{y}) \approx 2^{-nH(Y)}\}| \approx 2^{nH(Y)}$$

Similarly, given typical  $\mathbf{x}$ ,

$$|\{\mathbf{y} \mid p(\mathbf{y}|\mathbf{x}) \approx 2^{-nH(Y|X)}\}| \approx 2^{nH(Y|X)}$$



# Intuitive Idea

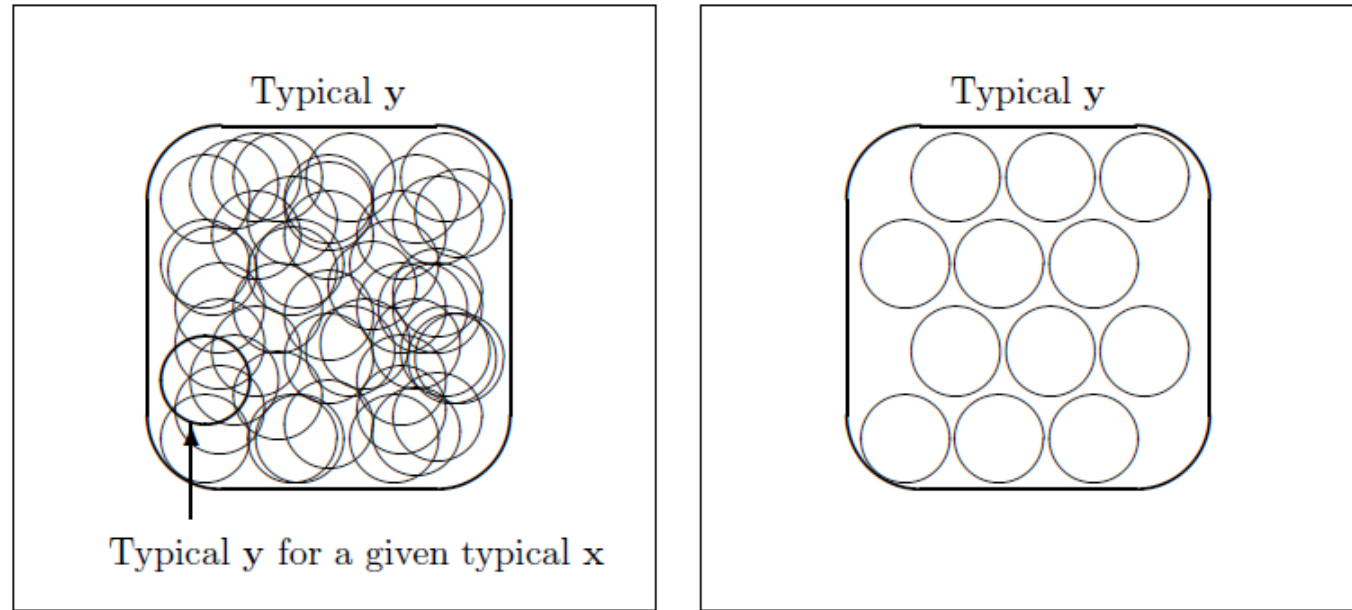
Consider sufficiently large  $n$ .

By AEP,

$$|\{\mathbf{y} \mid p(\mathbf{y}) \approx 2^{-nH(Y)}\}| \approx 2^{nH(Y)}$$

Similarly, given typical  $\mathbf{x}$ ,

$$|\{\mathbf{y} \mid p(\mathbf{y}|\mathbf{x}) \approx 2^{-nH(Y|X)}\}| \approx 2^{nH(Y|X)}$$



Ideally, we “pack”  $2^{nH(Y)} / 2^{nH(Y|X)} = 2^{nI(X;Y)}$  number of non-confusable typical  $\mathbf{y}$  for a given typical  $\mathbf{x}$ .

**Joint AEP** Independently sampling (about)  $2^{nI(X;Y)}$  #codewords suffices for non-confusability.

**Jointly-typical decoder**  $g$  (which is sub-optimal compared to  $g^*$ )

$g(\mathbf{y}) = j$  if  $(\mathbf{x}^{(j)}, \mathbf{y})$  is *jointly typical* and no other  $j'$  exists such that  $(\mathbf{x}^{(j')}, \mathbf{y})$  is *jointly typical*.

Otherwise, outputs arbitrary index.

# Joint AEP and Error Bound

Jointly typical set  $A_\epsilon^{(n)} = \{(\mathbf{x}, \mathbf{y}) \mid p(\mathbf{x}) \approx 2^{-nH(X)}, p(\mathbf{y}) \approx 2^{-nH(Y)}, p(\mathbf{x}, \mathbf{y}) \approx 2^{-nH(X,Y)}\}$

Joint AEP. For sufficiently large  $n$ ,

$$\Pr_{(\mathbf{X}, \mathbf{Y}) \sim p(\mathbf{x}, \mathbf{y})} \left[ (\mathbf{X}, \mathbf{Y}) \in A_\epsilon^{(n)} \right] \geq 1 - \epsilon \quad \text{contains most of the probability}$$

$$\left| A_\epsilon^{(n)} \right| \leq 2^{n(H(X,Y) + \epsilon)} \quad \text{cardinality } \lesssim 2^{nH(X)}$$

$$(1 - \epsilon)2^{-n(I(X;Y) + 3\epsilon)} \leq \Pr_{\substack{\mathbf{X} \sim p(\mathbf{x}) \\ \mathbf{Y} \sim p(\mathbf{y})}} \left[ (\mathbf{X}, \mathbf{Y}) \in A_\epsilon^{(n)} \right] \leq 2^{-n(I(X;Y) - 3\epsilon)} \quad \text{Independent } \mathbf{X}, \mathbf{Y} \text{ being jointly typical is exponentially small}$$

Error Bound. Fix any  $j$ .

- $\Pr[(\mathbf{X}^{(j)}, \mathbf{Y}) \text{ is not jointly typical}] < \epsilon$
  - For fixed  $j' \neq j$ ,  $\Pr[(\mathbf{X}^{(j')}, \mathbf{Y}) \text{ is jointly typical}] \lesssim 2^{-nI(X;Y)}$
- $\rightarrow \Pr[\exists j' \neq j : (\mathbf{X}^{(j')}, \mathbf{Y}) \text{ is jointly typical}] \lesssim (M' - 1)2^{-nI(X;Y)} \lesssim \epsilon$
- Union bound Holds if  $R \lesssim I(X;Y)$
- $$\Pr_{\substack{(M,n) \text{ code,} \\ J, \mathbf{Y} \sim p(\cdot | \mathbf{X}^{(J)})}} [g(\mathbf{Y}) \neq J] \lesssim 2\epsilon$$

## Theorem (part 2)

The following holds for any discrete memoryless channel  $Q = (\mathcal{X}, p(y|x), \mathcal{Y})$ .

Any  $(\lceil 2^{nR} \rceil, n)$  code with  $R > C := \max_{p_x} I(X; Y)$  has  $\lambda_{\text{avg}}$  bounded away from 0 for all  $n$ .

Proof sketch)

- For fixed encoder and decoder, we have  $J \rightarrow \mathbf{X}^{(J)} \rightarrow \mathbf{Y} \rightarrow g(\mathbf{Y})$ .

$$\begin{aligned} H(J) &= H(J|g(\mathbf{Y})) + I(J; g(\mathbf{Y})) \\ &\leq H(J|g(\mathbf{Y})) + I(\mathbf{X}^{(J)}; \mathbf{Y}) \quad \text{Data processing} \\ &\leq 1 + \Pr[J \neq g(\mathbf{Y})] \log(M - 1) + I(\mathbf{X}^{(J)}; \mathbf{Y}) \quad \text{Fano's inequality} \\ &\leq 1 + \Pr[J \neq g(\mathbf{Y})] nR + nC \quad \begin{array}{l} \text{Doing } n \text{ transmissions;} \\ I(X; Y) \text{ per transmission.} \end{array} \end{aligned}$$

Assuming  $J$  is sampled from a uniform distribution,  $H(J) = \log \lceil 2^{nR} \rceil \approx nR$  and thus

$$\Pr[J \neq g(\mathbf{Y})] \geq 1 - \frac{C}{R} - \frac{1}{nR}.$$

## Stronger Theorem (part 2)

The following holds for any discrete memoryless channel  $Q = (\mathcal{X}, p(y|x), \mathcal{Y})$ .

Any  $(\lceil 2^{nR} \rceil, n)$  code with  $R > C := \max_{p_x} I(X; Y)$  has  $\lambda_{\text{avg}} \approx 1$  for all sufficiently large  $n$ .

# Rate Distortion Theory

We have  $n$  i.i.d. RVs.

**Source coding theorem** says

with a distortion function  $d(\mathbf{x}, \hat{\mathbf{x}}) := \mathbb{I}_{\mathbf{x} \neq \hat{\mathbf{x}}}$ ,

- if  $R < H$ , then  $\mathbb{E}[d(\mathbf{X}, \hat{\mathbf{X}})] \leq 1 - \epsilon$  is not possible (when  $n$  is large);
- if  $R > H$ , then  $\mathbb{E}[d(\mathbf{X}, \hat{\mathbf{X}})] \leq \epsilon$  is possible (when  $n$  is large).

**Rate distortion theory** says

with a separable distortion function  $d(\cdot, \cdot)$ , i.e.,  $d(\mathbf{x}, \hat{\mathbf{x}}) = \frac{1}{n} \sum_{i=1}^n d(x_i, \hat{x}_i)$ ,

- if  $R < ?$ , then  $\mathbb{E}[d(\mathbf{X}, \hat{\mathbf{X}})] \leq D$  is not possible;
- if  $R > ?$ , then  $\mathbb{E}[d(\mathbf{X}, \hat{\mathbf{X}})] \leq D$  is possible.

# Rate and Distortion

Consider a (per symbol) distortion function  $d(x, \hat{x})$ , e.g.,  $d(x, \hat{x}) = \mathbb{I}_{x \neq \hat{x}}$  and  $d(x, \hat{x}) < \infty$ .

A compression scheme  $(f, g)_{n,R}$  compresses  $n$  symbols with  $R$  bits per symbol. i.e.,  $nR (\in \mathbb{N})$  bits in total.

- Codebook (for decoding):  $\{j, \hat{\mathbf{x}}^{(j)}\}_{j=1}^{2^{nR}}$

**Distortion** of  $(f, g)_{n,R} = \mathbb{E}_{\mathbf{X}} \left[ d \left( \mathbf{X}, g(f(\mathbf{X})) \right) \right]$ .

$(R, D)$  is **achievable** if for any  $\delta > 0$ ,

there exist a scheme  $(f, g)_{n,R}$  that compresses  $n$  i.i.d. symbols with distortion  $\leq D + \delta$ .

**Rate distortion function**  $R(D) = \min_{(R,D) \text{ achievable}} R$

Assume, given a distribution  $p(x)$  over  $\mathcal{X}$ , we wish to find a “test channel”  $(\mathcal{X}, p(\hat{x}|x), \mathcal{X})$  with distortion at most  $D$  such that it conveys a minimum amount of information.

$$R^{(I)}(D) := \min_{\substack{p(\hat{x}|x) : \\ \mathbb{E}_{\substack{X \sim p(x) \\ \hat{X} \sim p(\cdot|X)}} [d(X, \hat{X})] \leq D}} I(X; \hat{X})$$

# Rate Distortion Theory

**Theorem.**  $R(D) = R^{(I)}(D)$ .

proof sketch)

- Any  $(f, g)_{n,R}$  with distortion  $\leq D$  must satisfy  $R \geq R^{(I)}(D)$ .
  - $R^{(I)}(D)$  is nonincreasing and convex in  $D$ .
  - $d$  is separable.
  - Data inequalities.
- $(R^{(I)}(D), D)$  is achievable.
  - Ideally, construct a codebook  $\{j, \hat{\mathbf{x}}^{(j)}\}_{j=1}^{2^{nR^{(I)}(D)}}$  that “covers” all typical  $\mathbf{x}$  within distortion  $D$ .
  - **(Joint) Distortion AEP**

**Independently sampling  $2^{nR}$  number of typical  $\hat{\mathbf{x}}$  suffices for covering all typical  $\mathbf{x}$ .**



# Achievability

One way to show achievability = show existence of such scheme

- Construct a random scheme  $(f, g)_{n,R}$  where  $R > R^{(I)}(D)$ . Show the expected distortion  $\leq D$ .

One way to construct a random scheme

- Fix a distribution  $p(\hat{x}|x)$  with distortion =  $D \rightarrow p(\hat{x})$  is fixed.  
Fix a test channel.
- For each  $j \in [2^{nR}]$ , independently,  $\mathbf{X}^{(j)} = X_1^{(j)} X_2^{(j)} \dots X_n^{(j)}$  where  $X_i^{(j)} \sim p(\hat{x})$  independently.
  - Let  $g(j) = \hat{\mathbf{x}}^{(j)}$  for each  $j \in [2^{nR}]$ .
- Consider a (jointly) distortion typical  $f$ .
  - $f(\mathbf{x}) = \text{any } j \text{ such that } (\mathbf{x}, \hat{\mathbf{x}}^{(j)}) \text{ is distortion typical.}$
  - If there is no  $j$  such that  $(\mathbf{x}, \hat{\mathbf{x}}^{(j)})$  is distortion typical, outputs arbitrary  $j$ .

**Claim.** Random scheme has distortion  $\leq D + \delta$ .  $\rightarrow$  There is a scheme with distortion  $\leq D + \delta$ .

# (Joint) Distortion AEP

## (Jointly) Distortion typical set

$$A_{d(\cdot, \cdot), \epsilon}^{(n)} = \{ (\mathbf{x}, \mathbf{y}) \mid p(\mathbf{x}) \approx 2^{-nH(X)}, p(\mathbf{y}) \approx 2^{-nH(Y)}, p(\mathbf{x}, \mathbf{y}) \approx 2^{-nH(X,Y)}, d(\mathbf{x}, \mathbf{y}) \approx \mathbb{E}[d(\mathbf{X}, \mathbf{Y})] \}$$

Since  $d(\mathbf{x}, \mathbf{y}) = \frac{1}{n} \sum_{i=1}^n d(x_i, y_i)$  is separable, the law of large number can be applied.

**Joint (Distortion) AEP.** For sufficiently large  $n$ ,

$$\Pr_{(\mathbf{X}, \mathbf{Y}) \sim p(\mathbf{x}, \mathbf{y})} \left[ (\mathbf{X}, \mathbf{Y}) \in A_{d, \epsilon}^{(n)} \right] \geq 1 - \epsilon \quad \text{contains most of the probability}$$

## Distortion Bound.

- $X$  is distortion typical with some  $\hat{X} \rightarrow$  Contribution to distortion  $\approx D$
- $\Pr[X \text{ is not distortion typical with any } \hat{X}] \cdot \max d(\mathbf{x}, \hat{\mathbf{x}}) \rightarrow$  Contribution to distortion  $\approx 0$

This probability goes to zero *exponentially* fast if  $R > I(X; \hat{X})$ .

Since  $d(x, \hat{x})$  is bounded.

\* Stronger sense of typicality upper-bounds the distortion w.h.p., i.e.,  $\Pr[d(\mathbf{X}, \hat{\mathbf{X}}) > D + \delta] \approx 0$

# Characterization of $R(D)$

We showed  $R(D) = \min_{p(\hat{x}|x) : \substack{\mathbb{E}_{X \sim p(x)} [d(X, \hat{X})] \leq D \\ \hat{X} \sim p(\hat{x}|X)}} I(X; \hat{X})$ .

Solve the minimization problem of a convex function over the convex set of some distributions.

→ We obtain an optimal  $p(\hat{x}|x)$ .

**Blahut-Arimoto algorithm** computes two alternating minimization iteratively.

- Also converges to an optimal  $p(\hat{x}|x)$ .

## Channel Coding Theorem (part 3)

When given a discrete memoryless channel with a bounded separable distortion function, the distortion  $D$  is achievable iff  $C > R(D)$ .

# Differential Entropy

Continuous Random variable

# Differential Entropy

Consider a continuous random variable  $X$  with density  $f$ .

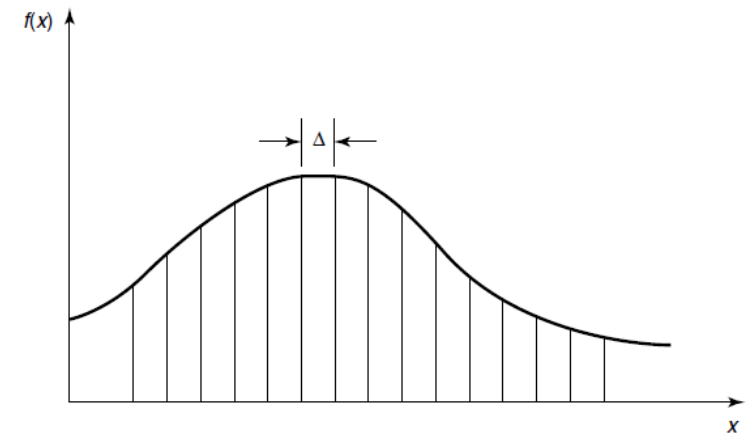
Divide range of  $X$  into bins of length  $\Delta$ .

Let  $x_i$  be a value such that  $f(x_i) \cdot \Delta = \int_{i\Delta}^{(i+1)\Delta} f(x) dx$ .

Consider a discrete random variable  $X^\Delta$  where  $X^\Delta = x_i$  with probability  $f(x_i)\Delta$ .

$$H(X^\Delta) = - \sum_{i=-\infty}^{\infty} f(x_i)\Delta \log(f(x_i)\Delta) = - \sum_{i=-\infty}^{\infty} f(x_i)\Delta \log f(x_i) - \log \Delta$$

As  $\Delta \rightarrow 0$ , ,  $H(X^\Delta) + \log \Delta \rightarrow - \int_{-\infty}^{\infty} f(x) \log f(x) dx$ . if  $f(x) \log f(x)$  is Riemann integrable



## Differential Entropy

$$h(X) = h(f) = - \int_{-\infty}^{\infty} f(x) \log f(x) dx$$

# Properties of Differential Entropy

- $h \neq$  amount of information (or entropy of quantized continuous RV)
  - $h(f) < 0$  is possible. Consider a density function  $f$  that corresponds to  $U[0,1/4]$ .

$$h(f) = - \int_0^{1/4} 4 \log 4 \, dx = -2$$

- Translation does not change  $h$ .  $h(X) = h(X + c)$ .
- Scaling does change  $h$ .  $h(aX) = h(X) + \log|a|$

- Differential entropy  $h(X)$  of gaussian RV  $X \sim N(0, \sigma^2) = \frac{1}{2} \log(2\pi\sigma^2 e)$

$$X \sim f(x) = \frac{1}{\sqrt{2\pi\sigma^2}} e^{-\frac{x^2}{2\sigma^2}}$$

$$\begin{aligned} h(f) &= - \int_{-\infty}^{\infty} f(x) \log f(x) \, dx = - \int_{-\infty}^{\infty} f(x) \left( \log \frac{1}{\sqrt{2\pi\sigma^2}} - \frac{x^2}{2\sigma^2} \log e \right) dx \\ &= \frac{1}{2} \log(2\pi\sigma^2) + \frac{\mathbb{E}_{X \sim N(0, \sigma^2)}[X^2]}{2\sigma^2} \log e = \frac{1}{2} \log(2\pi\sigma^2) + \frac{1}{2} \log e = \frac{1}{2} \log(2\pi\sigma^2 e) \end{aligned}$$

# KL Divergence and Mutual Information

**Kullback-Leibler divergence** between PDFs  $f$  and  $g$

$$D(f \parallel g) = \int \int f(x_1) \log \frac{f(x_1)}{g(x_2)} dx_1 dx_2$$

**Mutual information** between  $X$  and  $Y$  with joint density  $f(x, y)$

$$\begin{aligned} I(X; Y) &= \int \int f(x, y) \log \frac{f(x, y)}{f(x)f(y)} dx dy \\ &= D(f(x, y) \parallel f(x)f(y)) \\ &= h(X) - h(X|Y) = h(Y) - h(Y|X) \\ &\approx I(X^\Delta; Y^\Delta) \end{aligned}$$

**Theorem.**  $D(f \parallel g) \geq 0$  with  $=$  iff  $f = g$ .

**Corollary.**  $I(X; Y) \geq 0$  with  $=$  iff  $X$  and  $Y$  are independent.

**Corollary.**  $h(X|Y) \leq h(X)$  with  $=$  iff  $X$  and  $Y$  are independent.

**Corollary.** If  $X$  be a RV with support  $[-a, a]$ ,  $h(X) \leq h(U[-a, a])$  with equality iff  $X \sim U[-a, a]$ .

**Corollary.** If  $X$  be a RV with a variance  $\sigma^2$ ,  $h(X) \leq h(N(0, \sigma^2))$  with equality iff  $X \sim N(0, \sigma^2)$ .

# AEP

Consider a sequence of i.i.d. RVs  $X_1, X_2, \dots, X_n$ .

For any  $\epsilon > 0$ , for all sufficiently large  $n$ ,

$$\Pr[2^{-n(h(X)+\epsilon)} < f(X_1, X_2, \dots, X_n) < 2^{-n(h(X)-\epsilon)}] \geq 1 - \epsilon.$$

The **typical set**  $A_\epsilon^{(n)}$  w.r.t.  $f$  is the set of sequence  $\mathbf{x} = (x_1, x_2, \dots, x_n) \in S^n$  such that

$$2^{-n(h(X)+\epsilon)} < f(\mathbf{x}) < 2^{-n(h(X)-\epsilon)}. \quad S: \text{support of } f$$

**AEP.** For sufficiently large  $n$ ,

$$\Pr[\mathbf{X} \in A_\epsilon^{(n)}] \geq 1 - \epsilon \quad \text{and} \quad (1 - \epsilon)2^{n(h(X)-\epsilon)} \leq \text{Vol}(A_\epsilon^{(n)}) \leq 2^{n(h(X)+\epsilon)}.$$

contains most of the probability volume  $\approx 2^{nh(X)}$

Let  $B_\delta^{(n)} \subseteq S^n$  be a smallest set with  $\Pr[\mathbf{X} \in B_\delta^{(n)}] \geq 1 - \delta$ .

**Lemma.**  $\text{Vol}(B_\delta^{(n)}) \geq (1 - \epsilon - \delta)2^{n(h(X)-\epsilon)} \approx 2^{nh(X)}$ .

**$A_\epsilon^{(n)}$  contains most of the probability and essentially has a smallest volume of  $\approx 2^{nh(X)}$ .**



# Gaussian Channel

Previously, we considered discrete-time discrete-space input channel.

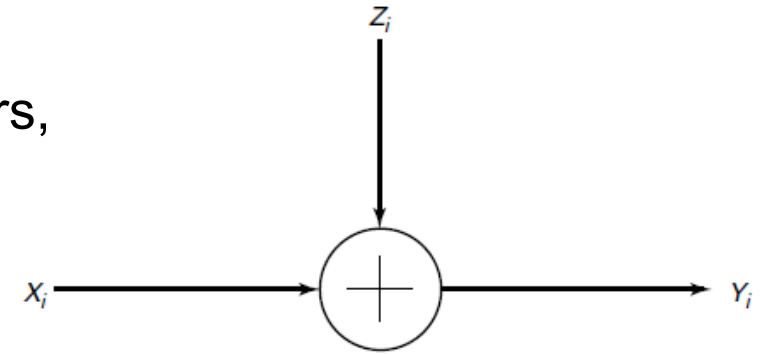
Gaussian channel is a discrete-time continuous-space input channel.

We also consider a continuous-time continuous-space input channel with bandlimit.

# (Discrete-time) Gaussian Channel

In a Gaussian channel, the input space is continuous, e.g., real numbers, and Gaussian noise is added to the input.

Note that it is a discrete-time channel.



- If no constraint on the input, it has infinite capacity.
  - Even if the noise variance  $> 0$ , it can transmit infinitely many numbers almost perfectly.
- Common constraint: upper bound on the average power of the input  $(x_1, x_2, \dots, x_n)$

$$\frac{1}{n} \sum_{i=1}^n x_i^2 \leq P$$

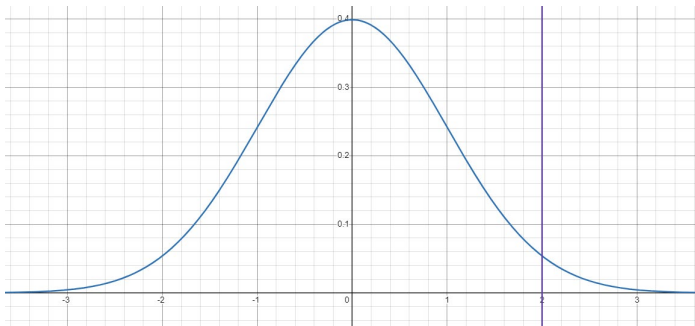
# Example

There is a Gaussian channel  $Q$  with input-space as  $\mathbb{R}$ , a Gaussian noise  $N(0,1)$  and power constraint  $P = 4$ .

Alice want to send the result of a random coin with  $x \in \mathbb{R}$  (i.e., length 1 codeword) through  $Q$ . Alice encode the result of the coin to a codeword  $X$  where  $X = 2$  if HEAD,  $X = -2$ , otherwise. Bob receives a noisy number  $Y(= X + Z)$ . Bob decodes  $Y$  to HEAD if  $Y > 0$ , TAIL, otherwise.

## Error-probability

$$\begin{aligned}\Pr[\text{error}] &= \Pr[Y \leq 0, X = 2] + \Pr[Y > 0, X = -2] \\ &= \Pr[Z \leq -2, X = 2] + \Pr[Z > 2, X = -2] \\ &= \Pr[Z > 2] \approx 0.022\end{aligned}$$



# Channel Capacity

The information capacity of the Gaussian channel  $N(0, \sigma^2)$  with power constraint  $P$

$$C = \max_{f(x): \mathbb{E}_{X \sim f}[X^2] \leq P} I(X; Y)$$

We have  $I(X; Y) = h(Y) - h(Y | X) = h(Y) - h(Z | X) = h(Y) - h(Z)$ .

Since  $Z \sim N(0, \sigma^2)$ , we have  $h(Z) = \frac{1}{2} \log(2\pi\sigma^2 e)$

Moreover, the variance of  $Y$  is  $\mathbb{E}[Y^2] = \mathbb{E}[X^2 + 2XZ + Z^2] = P + \sigma^2$ .

$$h(Y) \leq h(N(0, P + \sigma^2)) = \frac{1}{2} \log(2\pi(P + \sigma^2)e)$$

Therefore,

$$C = \frac{1}{2} \log\left(1 + \frac{P}{\sigma^2}\right)$$

# Code, Decode, Rate, Error

$(M, n)$  **code** for a channel  $N(0, \sigma^2)$  with power constraint  $P$

- a set of indices(data)  $\{1, \dots, M\}$
- a set of codewords  $\{\mathbf{x}^{(1)}, \dots, \mathbf{x}^{(M)}\}$ 
  - $\mathbf{x}^{(j)} = x_1^{(j)} x_1^{(j)} \dots x_n^{(j)}$  such that  $x_1^{(j)2} + x_2^{(j)2} + \dots + x_n^{(j)2} \leq nP$

A **decoder**  $g$  guess an index among  $[M]$ .

**Rate**  $R$  of  $(M, n)$  code for  $Q$

$$R = \frac{\text{size of data}}{\text{\#transmission}} = \frac{\log M}{n}$$

**Maximal probability of error** (for a fixed channel  $Q$  and fixed  $(M, n)$  code for  $Q$ )

$$\lambda_{\max} := \max_{j \in \{1, \dots, M\}} \Pr_{\mathbf{Z} \sim N^n(0, \sigma^2)} [g(\mathbf{Z} + \mathbf{x}^{(j)}) \neq j]$$

# Theorem

The following holds for any Gaussian channel  $Q$  with  $N(0, \sigma^2)$  with power constraint  $P$ .

1) For any  $\epsilon > 0$  and  $R < C$ , there exists a  $(M := \lceil 2^{nR} \rceil, n)$  code for  $Q$  such that

$$\lambda_{\max} < \epsilon \text{ for all sufficiently large } n.$$

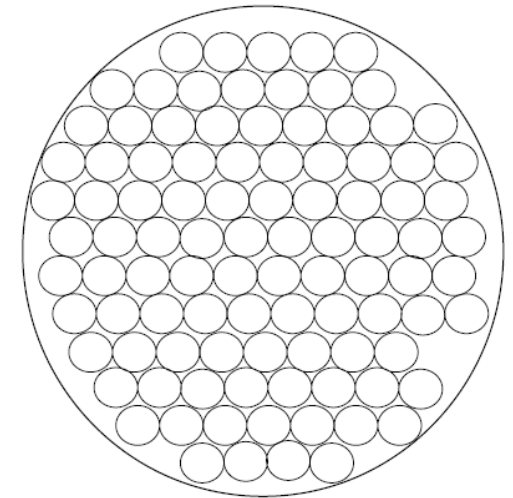
2) Any  $(\lceil 2^{nR} \rceil, n)$  code with  $R > C$  has  $\lambda_{\text{avg}}$  bounded away from 0 for all  $n$ .

For part 1, it uses Joint AEP.

Ideally, we “pack” non-confusable typical  $\mathbf{y}$  for a given typical  $\mathbf{x}$ .

Joint AEP says independent sampling suffices to achieve this.

For part 2, it utilizes data processing and Fano’s inequality.



# Continuous-time Gaussian Channel and Bandlimit

Previously, we considered discrete-time channels. ( $n$  usage of given channel)

Now, consider a **continuous-time Gaussian channel**.

- an input *signal*  $x(t)$
- *additive white Gaussian noise*  $Z(t)$
- power constraint  $P$  (defined in a continuous manner)

Consider a **bandlimited** (continuous-time) Gaussian channel

- Channel cuts out all frequencies greater than  $W$  (e.g., by applying a bandpass filter)

## Nyquist-Shannon's Theorem

Sampling a signal that is bandlimited to  $W$  at a sampling rate  $\frac{1}{2W}$  is sufficient for the reconstruction.

**Discretize input signal. → Send through discrete-time channel for multiple times.**

# Capacity of Bandlimited Gaussian Channel

Consider a continuous-time Gaussian channel  $Q$  with

- bandwidth  $W$  Hz, power  $P$ , and *power spectral density* of noise  $N_0/2$  W/Hz.

By Nyquist-Shannon's theorem, it is equivalent to  $2W$  usage (per sec) of a discrete-time Gaussian channel  $Q'$  with power constraint  $P/2W$ , and noise  $N(0, N_0/2)$  of

Note the capacity of  $Q'$  is  $C' = \frac{1}{2} \log \left( 1 + \frac{P}{WN_0} \right)$ .

Then the capacity of  $Q$  is

$$C = 2W \cdot C' = W \log \left( 1 + \frac{P}{WN_0} \right)$$



# Statistics

Type gives a stronger sense of AEP.

# Type

$\text{type}(\mathbf{x})$ : **type** of a sequence  $\mathbf{x}$ , i.e., the frequency of each symbol in  $\mathcal{X}$

$\text{typeclass}(t)$ : **type class of type  $t$** , i.e., a set of sequence (of length  $n$ ) whose type is  $t$

Observation.

Exponential #sequence of length  $n$  ( $= |\mathcal{X}|^n$ ).

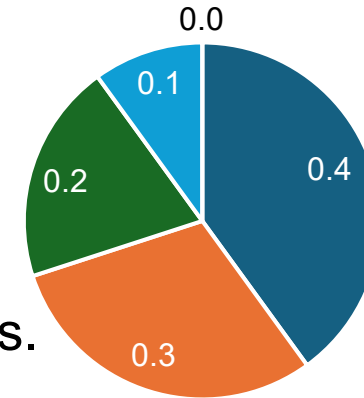
Polynomial #types of a sequence of length  $n$  ( $\leq (n + 1)^{|\mathcal{X}|}$ ).

→ ∃ a type whose type class contains *exponentially* many sequences.

In fact, every type class contains

$$\frac{1}{(n + 1)^{|\mathcal{X}|}} 2^{nH(t)} \leq |\text{typeclass}(t)| \leq 2^{nH(t)}$$

type of *acbdaabcba*



type is a distribution over  $\mathcal{X}$

■ a ■ b ■ c ■ d ■ e~z

type class of the above type

*aaaabbbccd*

*aaaabbbcdc*

...

*dccbbbaaaa*

# “AEP” for Type and Universal Compression

Consider a sequence of i.i.d. RVs  $\mathbf{X} = X_1 X_2 \cdots X_n$  where  $X_i \sim p$ .

Note that  $\text{type}(\mathbf{X})$  is a random distribution over  $\mathcal{X}$ .

Let typical sequence be a sequence  $\mathbf{x} \in \mathcal{X}^n$  such that  $D(\text{type}(\mathbf{x}) \parallel p) \approx 0$ .

“**AEP**” (for type). If  $n$  is sufficiently large,  $\Pr[D(\text{type}(\mathbf{X}) \parallel p) \approx 0] \approx 1$ .

**For almost every sequence, the sample frequencies are close to the true probability.**

## Corollary (Universal Codes).

Even if  $p$  is unknown, we can compress an i.i.d. source with (very close to)  $H(p)$ -bit per symbol.

# Sanov's Theorem

Consider a sequence of i.i.d. dices  $\mathbf{X} = X_1 X_2 \cdots X_n$  where  $X_i \sim p$ .

**Q.** What is the probability being  $\sum_i X_i \geq 4n$ ?

**A1.** Central limit theorem, i.e., the distribution of the sample mean  $\rightarrow$  a normal distribution.

- Poor approximation...

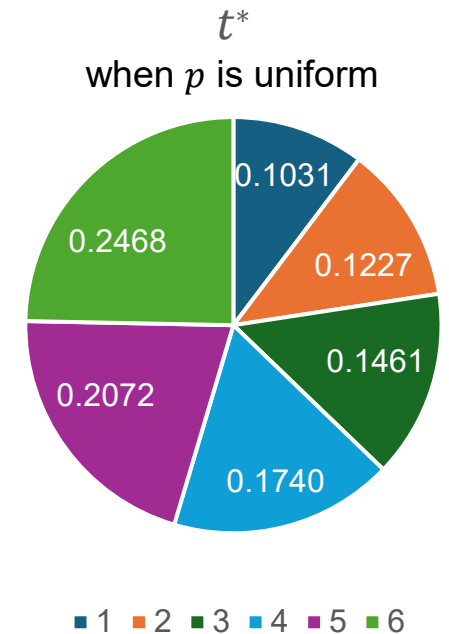
**A2.** Let  $\mathcal{T}$  be the set of types of sequence whose sum is at least  $4n$ , i.e.,

$$\mathcal{T} = \{\text{type}(\mathbf{x}) \mid \text{sum}(\mathbf{x}) \geq 4n\}.$$

**Sanov's theorem** says

$$\Pr[\text{type}(\mathbf{X}) \in \mathcal{T}] \leq |\mathcal{T}| \cdot 2^{-nD(t^* \| p)} \leq (n+1)^{|\mathcal{X}|} \cdot 2^{-nD(t^* \| p)},$$

where  $t^* \in \operatorname{argmin}_{t \in \mathcal{T}} D(t \| p)$ .



**The probability measure of  $\mathcal{T}$  is essentially determined by  $t^*$ .**

= probability of large deviation

# Conditional Limit Theorem

Consider a sequence of i.i.d. dices  $\mathbf{X} = X_1 X_2 \cdots X_n$  where  $X_i \sim p$ .

**Q.** Suppose  $\sum_i X_i \geq 4n$ . What can we say about the marginal probability distribution?

**A.** Let  $\mathcal{T}$  be the set of types of sequence whose sum is at least  $4n$ , i.e.,

$$\mathcal{T} = \{\text{type}(\mathbf{x}) \mid \text{sum}(\mathbf{x}) \geq 4n\},$$

and let  $t^* \in \underset{t \in \mathcal{T}}{\text{argmin}} D(t \parallel p)$  be a “closest” distribution in  $\mathcal{T}$  to  $p$ .

**Conditional limit theorem** says, for any  $x \in \mathcal{X}$ , if  $\mathcal{T}$  is a closed convex set.

$$\Pr[X_1 = x \mid \text{type}(\mathbf{X}) \in \mathcal{T}] \rightarrow t^*(x).$$

**The probability measure of  $\mathcal{T}$  is not only determined by  $t^*$   
but also concentrated near  $t^*$  i.e., the *conditional type* is close to  $t^*$ .**

# Another Proof of Joint AEP

product of marginal distributions of  $p$

Consider a sequence of pairs of i.i.d. RVs  $(\mathbf{X}, \mathbf{Y}) = (X_1 Y_1)(X_2 Y_2) \cdots (X_n Y_n)$  where  $(X_i Y_i) \sim p_x p_y$

Recall a jointly typical set  $A_\epsilon^{(n)} = \{(\mathbf{x}, \mathbf{y}) \mid p_x(\mathbf{x}) \approx 2^{-nH(X)}, p_y(\mathbf{y}) \approx 2^{-nH(Y)}, p(\mathbf{x}, \mathbf{y}) \approx 2^{-nH(X,Y)}\}$ .

Let  $\mathcal{T}$  be the set of types of typical sequence, i.e.,  $\mathcal{T} = \{\text{type}(\mathbf{x}, \mathbf{y}) \mid (\mathbf{x}, \mathbf{y}) \in A_\epsilon^{(n)}\}$ .

By applying Sanov's theorem,

$$\Pr_{(\mathbf{X}, \mathbf{Y}) \sim \prod p_x p_y} [\text{type}(\mathbf{X}, \mathbf{Y}) \in \mathcal{T}] \leq (n+1)^{|\mathcal{X}|} \cdot 2^{-nD(t^* \parallel p_x p_y)},$$

where  $t^* \in \underset{t \in \mathcal{T}}{\text{argmin}} D(t \parallel p_x p_y)$ .

In fact,  $t^* = p$  and thus  $\Pr_{(\mathbf{X}, \mathbf{Y}) \sim \prod p_x p_y} [\text{type}(\mathbf{X}, \mathbf{Y}) \in \mathcal{T}] \leq 2^{-n(I(X;Y)+\epsilon)}$ . when  $n$  is sufficiently large.

Also, by conditional limit theorem, its conditional type is likely to be close to  $p$ . when  $n$  is sufficiently large.

# Two Hypothesis Testing

Trade-off between  $\alpha := \Pr[\text{choosed } P_2 \text{ but } P_1 \text{ is true}]$  and  $\beta := \Pr[\text{choosed } P_1 \text{ but } P_2 \text{ is true}]$ .

**Neyman-Pearson lemma** says the optimum test is a (log-)likelihood ratio test.

$$\frac{P_1(\mathbf{X})}{P_2(\mathbf{X})} \geq \tau \Leftrightarrow D(\text{type}(\mathbf{X}) \parallel P_2) - D(\text{type}(\mathbf{X}) \parallel P_1) \geq \frac{1}{n} \log \tau$$

Let  $\mathcal{T}$  be the set of types that satisfies the above, i.e.,  $t \in \mathcal{T}$  accepts  $P_1$  and  $t \notin \mathcal{T}$  accepts  $P_2$ .

Suppose  $P_1$  was the true distribution, i.e.,  $\mathbf{X} \sim \prod P_1$ . Sanov's theorem gives

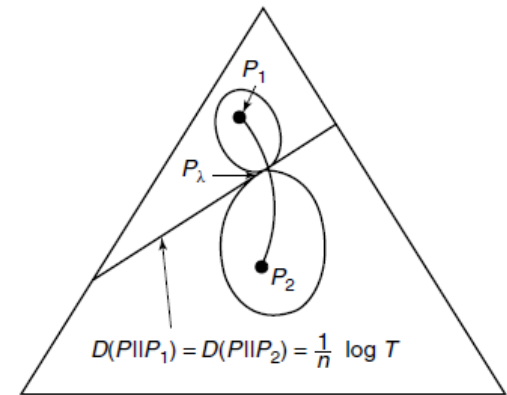
$$\alpha = \Pr[\text{type}(\mathbf{X}) \notin \mathcal{T}] = 2^{-n(D(t_1^* \parallel P_1) - \epsilon_1)}.$$

Similarly, if  $P_2$  was the true distribution, we have

$$\beta = \Pr[\text{type}(\mathbf{X}) \in \mathcal{T}] = 2^{-n(D(t_2^* \parallel P_2) - \epsilon_2)}.$$

We can also show  $t_1^* = \underset{t \notin \mathcal{T}}{\operatorname{argmin}} D(t \parallel P_1) = \underset{t \in \mathcal{T}}{\operatorname{argmin}} D(t \parallel P_2) = t_2^*$ .

**Chernoff-Stein lemma** says if  $\alpha < \epsilon$  is small,  $\beta \approx 2^{-n(D(P_1 \parallel P_2) - \epsilon)}$  is best possible. by AEP for KL-divergence.



# Kolmogorov Complexity

Source coding says  $\approx H(X)$  bits are required to describe  $X$ .

What is the shortest length of a program that describes (or outputs)  $X$ ?

“Approximately equal to its entropy”



# Kolmogorov Complexity

Consider any universal Turing machine  $\mathcal{U}$ .

$K_{\mathcal{U}}(\mathbf{x}) := \min_{p: \mathcal{U}(p)=\mathbf{x}} \ell(p)$  is the length of a shortest program that prints  $\mathbf{x}$  and halts (w.r.t.  $\mathcal{U}$ ).

Consider another Turing machine  $\mathcal{A}$  and let  $p_{\mathcal{A}}$  be a program for  $\mathcal{A}$  that prints  $\mathbf{x}$  and halts.

Consider a program  $s_{\mathcal{A}}$  for  $\mathcal{U}$  that simulates  $\mathcal{A}$  on  $\mathcal{U}$ .

Consider an input  $s_{\mathcal{A}}p_{\mathcal{A}}$  to  $\mathcal{U}$ ; it prints  $\mathbf{x}$  and halts. Therefore,

$$K_{\mathcal{U}}(\mathbf{x}) \leq K_{\mathcal{A}}(\mathbf{x}) + c_{\mathcal{A}}$$

where  $c_{\mathcal{A}} := \ell(s_{\mathcal{A}})$  is a constant.

$K(\mathbf{x})$  differs by a constant for any two universal Turing machines.

# Kolmogorov Complexity

Consider a Kolmogorov complexity when the length of  $x$  is additionally given.

$K(\mathbf{x} \mid \ell(\mathbf{x}))$  is the length of a shortest program such that when given  $n := \ell(\mathbf{x})$ , prints  $\mathbf{x}$  and halts.

Consider a program  $p$ : “print the first  $n$ -bit  $x_1x_2 \cdots x_n$ ”; its length  $\ell(p)$  is  $n + c$ .

$$K(\mathbf{x} \mid n) \leq n + c$$

However, we cannot say  $K(\mathbf{x}) \leq n + c$ , since if  $n$  is unknown,  $p$  does not know when to stop.

Consider a program  $p'$ : “read the first  $2\lceil \log n \rceil + 2$  bits and decide  $n$ ; print the next  $n$ -bit.”

if  $n = 5 = 101_{(2)}$ , we can describe  $n$  as 11001101 with **01** meaning ‘,’

We can upper bound  $K(\mathbf{x}) \leq K(\mathbf{x} \mid n) + 2 \log n + c'$ .

# Kolmogorov Complexity (Information Theory)

Consider a sequence of i.i.d. (binary) RVs  $\mathbf{X} = X_1 X_2 \cdots X_n$ .

Source coding theorem says  $\frac{1}{n} \mathbb{E}[K(\mathbf{X} | n)] \geq H(X)$ . A shortest program is a compression of  $\mathbf{X}$ .

Consider any type  $t$  and its type class  $\text{typeclass}(t)$ . We index each  $\mathbf{x} \in \text{typeclass}(t)$ .

Consider a program  $p$ : “print  $i$ -th string  $\mathbf{x}$  of  $\text{typeclass}(t)$ ” Recall  $|\text{typeclass}(t)| \leq 2^{nH(t)}$

- To describe a type,  $|\mathcal{X}| \log n$  bits suffice. To describe an index,  $nH(t)$  bits suffices.

Since the sample frequencies are close to the true probability, we have

$$\frac{1}{n} \mathbb{E}[K(\mathbf{X} | n)] \leq H(X) + \frac{|\mathcal{X}| \log n}{n} + \frac{c}{n}.$$

Since  $K(\mathbf{x}) \leq K(\mathbf{x} | n) + 2 \log n + c'$ ,

$$H(X) \leq \frac{1}{n} \mathbb{E}[K(\mathbf{X})] \leq H(X) + \frac{1}{n} (\log |\mathcal{X}| + c')$$

\* also holds without expectation

# Kolmogorov Complexity (Theory of Computation)

There is no program that decides  $K(\mathbf{x}) = k$  (for input  $\mathbf{x}, k$ ).

- Suppose t.c. there is such program  $p$ .
- Fix large  $k$  such that it satisfies  $k > \ell(p) + \log k + c$ . (by running  $p$ )
- Consider a program  $q$ : “iterates until it finds a string  $\mathbf{y}$  where  $K(\mathbf{y}) > k$ ; print  $\mathbf{y}$ ”  
(e.g., in lexicographical order)
- $\ell(q) = \ell(p) + \log k + c$
- However,  $k < K(\mathbf{y}) \leq \ell(q) = \ell(p) + \log k + c$ , which is a contradiction. “Berry paradox”

Shares the essential spirit with the noncomputability of the Halting problem (Chaitin's number) and Gödel's incompleteness theorem.

# Side Notes

# Sufficient Statistics

$\{f_\theta(x)\}$ : a family of PMFs indexed by  $\theta$

$X$ : a sample from a distribution in  $\{f_\theta(x)\}$ .

$T(X)$ : any statistics (such as sample mean or sample variance).

Then  $\theta \rightarrow X \rightarrow T(X)$ , and by the data-processing inequality, for any distribution on  $\theta$ ,

$$I(\theta; X) \geq I(\theta; T(X)).$$

If it holds with equality, i.e.,  $\theta \rightarrow T(X) \rightarrow X$ , no information is lost.

Once we know  $T(X)$ ,  
the remaining randomness in  $X$   
does not depend on  $\theta$ .

We say  $T(X)$  is a **sufficient statistics** for  $\theta$ .

Example)

$\mathbf{X} = X_1, X_2, \dots, X_{10}$  be an i.i.d. sequence of coin w.p.  $\theta$  (chosen randomly).

Let  $T(\mathbf{X}) = X_1 + \dots + X_{10}$  be the #1's.

$T(\mathbf{X})$  is a sufficient statistics for  $\theta$ .

$$I(\theta; \mathbf{X}) = H(\theta) - H(\theta|\mathbf{X}) = H(\theta) - H(\theta|T(\mathbf{X})) = I(\theta; T(\mathbf{X}))$$

# Sufficient Statistics

$\{f_\theta(x)\}$ : a family of PMFs indexed by  $\theta$

$X$ : a sample from a distribution in  $\{f_\theta(x)\}$ .

$T(X)$ : any statistics (such as sample mean or sample variance).

Then  $\theta \rightarrow X \rightarrow T(X)$ , and by the data-processing inequality, for any distribution on  $\theta$ ,

$$I(\theta; X) \geq I(\theta; T(X)).$$

If it holds with equality, i.e.,  $\theta \rightarrow T(X) \rightarrow X$ , no information is lost.

Once we know  $T(X)$ ,  
the remaining randomness in  $X$   
does not depend on  $\theta$ .

We say  $T(X)$  is a **sufficient statistics** for  $\theta$ .

If  $T$  is a function of every other sufficient statistic  $U$ , i.e.,  $\theta \rightarrow X \rightarrow U(X) \rightarrow T(X)$ ,

$$I(\theta; X) \left( \geq I(\theta; U(X)) \right) \geq I(\theta; T(X)).$$

If it holds with equality, i.e.,  $\theta \rightarrow T(X) \rightarrow U(X) \rightarrow X$ ,

we say  $T(X)$  is a **minimal sufficient statistics** for  $\theta$ .

# Fano's Inequality

**Theorem.** For any estimator  $\hat{X}$  s.t.  $X \rightarrow Y \rightarrow \hat{X}$ , we have

$$H(X|Y) \leq H(X|\hat{X}) \leq H(1_{\hat{X} \neq X}) + \Pr[\hat{X} \neq X] \log|\mathcal{X}| \leq 1 + \Pr[\hat{X} \neq X] \log|\mathcal{X}|.$$

## Probability of Error and Entropy

**Lemma.** If  $X$  and  $X'$  are i.i.d.,  $\Pr[X = X'] \geq 2^{-H(X)}$  with equality iff  $X$  has a uniform distribution.

**Corollary.** If  $X \sim p$  and  $X' \sim q$  are independent and  $\mathcal{X} = \mathcal{X}'$ ,

$$\Pr[X = X'] \geq 2^{-H(p) - D(p||q)}$$

$$\Pr[X = X'] \geq 2^{-H(q) - D(q||p)}$$



# Stochastic Process and Entropy Rate

*Stochastic process*  $\{X_i\}$ : an indexed sequence of RVs with arbitrary dependence

**Stationary** stochastic process: joint distribution of any subset is invariant w.r.t. shifts in index

$$\Pr[(X_1, X_2, \dots, X_n) = (x_1, x_2, \dots, x_n)] = \Pr[(X_{1+\ell}, X_{2+\ell}, \dots, X_{n+\ell}) = (x_1, x_2, \dots, x_n)]$$

## Entropy rate

Definition 1 (entropy per symbol).

$$H(\mathcal{X}) = \lim_{n \rightarrow \infty} \frac{1}{n} H(X_1, X_2, \dots, X_n) \quad \text{when the limit exists}$$

Definition 2 (conditional entropy of the last).

$$H'(\mathcal{X}) = \lim_{n \rightarrow \infty} H(X_n | X_{n-1}, X_{n-2}, \dots, X_1) \quad \text{when the limit exists}$$

**Theorem.** For a stationary stochastic process,  $H(\mathcal{X}) = H'(\mathcal{X})$ .

# General AEP

## AEP

For any i.i.d. process, in probability,

$$-\frac{1}{n} \log p(X_1, \dots, X_n) \rightarrow H(X)$$

## General AEP (chapter 16)

For any stationary *ergodic* process, with probability 1,

$$-\frac{1}{n} \log p(X_1, \dots, X_n) \rightarrow H(\mathcal{X})$$

# Entropy Rate of Stationary Markov Chain

With initial dist. as stationary dist.  $\mu$ , Markov chain is a stationary process.

$$H(\mathcal{X}) = H'(\mathcal{X}) = \lim_{n \rightarrow \infty} H(X_n | X_{n-1}, X_{n-2}, \dots, X_1) = \lim_{n \rightarrow \infty} H(X_n | X_{n-1}) = H(X_2 | X_1)$$

Markovity stationarity

# Entropy Rate of Functions of Markov Chain

Let  $\{X_i\}$  be a stationary Markov chain. Consider  $\{Y_i\}$  where  $Y_i = \phi(X_i)$ .  $\{Y_i\}$  does not necessarily form a Markov chain.

How to know  $H(Y_n | Y_{n-1}, Y_{n-2}, \dots, Y_1) \approx H(\mathcal{Y})$  for any  $n$ ?

$$H(Y_n | Y_{n-1}, Y_{n-2}, \dots, Y_1, X_1) \leq H(\mathcal{Y}) \leq H(Y_n | Y_{n-1}, Y_{n-2}, \dots, Y_1)$$

$$\lim_{n \rightarrow \infty} H(Y_n | Y_{n-1}, Y_{n-2}, \dots, Y_1, X_1) = H(\mathcal{Y}) = \lim_{n \rightarrow \infty} H(Y_n | Y_{n-1}, Y_{n-2}, \dots, Y_1)$$

Relates to a *hidden Markov model* (HMM)

# Thank You