

Part 1: Theory

L03: Basics of entropy (1/6)

[measures of information, intuition behind entropy]

Wolfgang Gatterbauer, Javed Aslam

cs7840 Foundations and Applications of Information Theory (fa24)

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

9/11/2024

Let's gain some intuition for "measures of information"

The following numeric examples with hats and 4 balls are based on Chapter 1.1 from [Moser'18] Information Theory (lecture notes, 6th ed). https://moser-isi.ethz.ch/cgi-bin/request_script.cgi?script=it

Let's gain some intuition: What is information?

What is information? Let's look at some sentences with "information":

1. "It will rain tomorrow."
2. "It will snow tomorrow."
3. "The name of the next president of the USA will be..."
 - a. ... Donald."
 - b. ... Donald Duck."
4. "Our university is called Northeastern University."



Let's gain some intuition: What is information?

What is information? Let's look at some sentences with "information":

1. "It will rain tomorrow."
2. "It will snow tomorrow."
3. "The name of the next president of the USA will be..."
 - a. ... Donald."
 - b. ... Donald Duck."
4. "Our university is called Northeastern University."

⇒ Information (in a sentence) is linked to surprise (which is the delta of knowledge before and after seeing the sentence).

Let's next try to quantify "information" 😊

Let's try to quantify "information"

EXAMPLE 1: A gambler throws a fair die with 4 sides {**A**, **B**, **C**, **D**}.



- "Side **C** comes up."

- The "pure" message U_1 that we care about in our abstraction is ...



Let's try to quantify "information"

EXAMPLE 1: A gambler throws a fair die with 4 sides {**A**, **B**, **C**, **D**}.



- "Side **C** comes up."
- message $U_1 = \text{"C"}$

EXAMPLE 2: A gambler throws a fair die with **6** sides {**A**, **B**, **C**, **D**, **E**, **F**}.

- "Side **C** comes up."
- message $U_2 = \text{"C"}$



what has changed ?

Let's try to quantify "information"

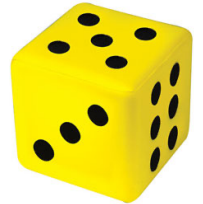
EXAMPLE 1: A gambler throws a fair die with 4 sides {**A**, **B**, **C**, **D**}.



- "Side **C** comes up."
- message $U_1 = \text{"C"}$
- There are **4** possible outcomes, each has a probability of $\frac{1}{4}$.

EXAMPLE 2: A gambler throws a fair die with 6 sides {**A**, **B**, **C**, **D**, **E**, **F**}.

- "Side **C** comes up."
- message $U_2 = \text{"C"}$
- There are **6** possible outcomes, each has a probability of $\frac{1}{6}$.



*⇒ 1) The number of possible outcomes should be linked to "information"
(we need more space to encode a message)*

Let's try to quantify "information"

EXAMPLE 1: A gambler throws a fair die with 4 sides {**A**, **B**, **C**, **D**}.



– "Side **C** comes up."

00 01 10 11

– message $U_1 = \text{"C"}$, or in above binary encoding $U_1 = \text{"10"}$

– There are **4** possible outcomes, each has a probability of $\frac{1}{4}$.

EXAMPLE 2: A gambler throws a fair die with 6 sides {**A**, **B**, **C**, **D**, **E**, **F**}.

– "Side **C** comes up."

000 001 010 011 100 101

– message $U_2 = \text{"C"}$, or in above binary encoding $U_2 = \text{"010"}$

– There are **6** possible outcomes, each has a probability of $\frac{1}{6}$.



⇒ 1) The number of possible outcomes should be linked to "information"
(we need more space to encode a message)

Let's try to quantify "information"

EXAMPLE 1: A gambler throws a fair die with 4 sides {**A**, **B**, **C**, **D**}.



- "Side **C** comes up."
- message $U_1 = \text{"C"}$
- There are 4 possible outcomes, each has a probability of $\frac{1}{4}$.

EXAMPLE 3: The gambler throws the 4-sided die **three times**.

- "The sequence of sides are: (**C**, **B**, **D**)"
- The message $U_3 = \text{"CBD"}$.



How many outcomes do we have now ?

Notice "**BCD**" is not the same as "**CBD**"

Let's try to quantify "information"

EXAMPLE 1: A gambler throws a fair die with 4 sides {**A**, **B**, **C**, **D**}.



- "Side **C** comes up."
- message $U_1 = \text{"C"}$
- There are **4** possible outcomes, each has a probability of $\frac{1}{4}$.

EXAMPLE 3: The gambler throws the 4-sided die **three times**.

- "The sequence of sides are: (**C**, **B**, **D**)"
- The message $U_3 = \text{"CBD"}$.
- Now we had **64** = $4 \cdot 4 \cdot 4 = 4^3$ possible outcomes.



16 times more!

How much more information did we learn in situation 3?



Let's try to quantify "information"

EXAMPLE 1: A gambler throws a fair die with 4 sides {**A**, **B**, **C**, **D**}.



- "Side **C** comes up."
- message $U_1 = \text{"C"}$
- There are 4 possible outcomes, each has a probability of $\frac{1}{4}$.

EXAMPLE 3: The gambler throws the 4-sided die **three times**.

- "The sequence of sides are: (**C**, **B**, **D**)"
- The message $U_3 = \text{"CBD"}$.
- Now we had $64 = 4 \cdot 4 \cdot 4 = 4^3$ possible outcomes.



We have 3 independent throws, the message U is 3 times as long, despite 4^3 possible total outcomes. Our information is 3 times as much.

⇒ 2) Information is additive in some sense

Hartley's measure of information [1928]



1 roll has 4 outcomes.



3 rolls have $64 = 4 \cdot 4 \cdot 4 = 4^3$ outcomes.

$$\log_4(4) = 1$$

$$\log_4(64) = 3$$

Hartley's insight: use the **logarithm of the number of possible outcomes r** to measure the amount of information in an outcome.

Hartley's measure
of information

$$H_0(U) = \log_b(n)$$

n = number of outcomes



Hartley's measure of information [1928]



1 roll has 4 outcomes.



3 rolls have $64 = 4 \cdot 4 \cdot 4 = 4^3$ outcomes.

$$\log_4(4) = 1$$

$$\log_4(64) = 3$$

Hartley's insight: use the **logarithm of the number of possible outcomes r** to measure the amount of information in an outcome.



Hartley's measure of information

$$H_0(U) = \log_b(n)$$

n = number of outcomes

The basis b of the logarithm is not really important.
(just unit of information, like 1 km = 1000 m)

$$\log_2(c) = 1.443 \cdot \log_e(c)$$

$$2^{1.443} = e \Leftrightarrow 1.443 = \log_2(e)$$

We will use: $\lg(c)$

$$e^z = (2^{1.443})^z = 2^{1.443 \cdot z}$$

Hartley's measure of information [1928]



1 roll has 4 outcomes.



3 rolls have $64 = 4 \cdot 4 \cdot 4 = 4^3$ outcomes.

$$\log_4(4) = 1$$

$$\log_4(64) = 3$$

Hartley's insight: use the **logarithm of the number of possible outcomes r** to measure the amount of information in an outcome.

Hartley's measure of information

$$H_0(U) = \log_b(n)$$

n = number of outcomes



For k independent trials,
the amount of information is:

$$\log_b(n^k) = ?$$

Hartley's measure of information [1928]



1 roll has 4 outcomes.



3 rolls have $64 = 4 \cdot 4 \cdot 4 = 4^3$ outcomes.

$$\log_4(4) = 1$$

$$\log_4(64) = 3$$

Hartley's insight: use the **logarithm of the number of possible outcomes r** to measure the amount of information in an outcome.

Hartley's measure of information

$$H_0(U) = \log_b(n)$$

n = number of outcomes



For k independent trials,
the amount of information is:

$$\log_b(n^k) = k \cdot \log_b(n)$$

the power of the **logarithm** 😊

Let's practice

EXAMPLE 4: A country has 1 million telephones. How long does the country's telephone numbers need to be?



Let's practice

EXAMPLE 4: A country has 1 million telephones. How long does the country's telephone numbers need to be?

$$\log_{10}(1,000,000) = 6$$

With 6 digits (like "123 456") we can represent 10^6 different telephones.

EXAMPLE 5: The current world population is 8,174,891,806 (as of Sat, September 7, 2024). How long must a binary telephone number be to connect to every person?

A tip: $2^{32} = 4,294, \dots, \dots$



Let's practice

EXAMPLE 4: A country has 1 million telephones. How long does the country's telephone numbers need to be?

$$\log_{10}(1,000,000) = 6$$

With 6 digits (like "123 456") we can represent 10^6 different telephones.

EXAMPLE 5: The current world population is 8,174,891,806 (as of Sat, September 7, 2024). How long must a binary telephone number be to connect to every person?

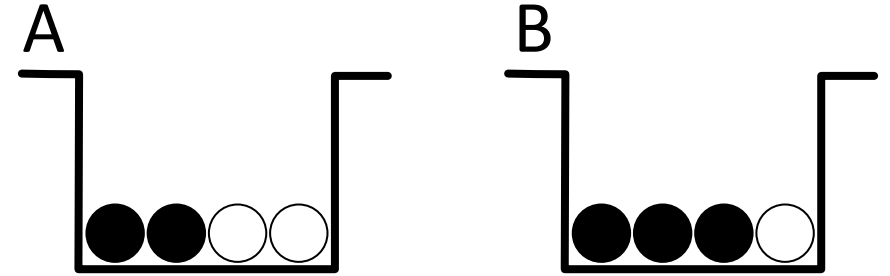
A tip: $2^{32} = 4,294, \dots, \dots$

$$\log_2(8,174,891,806) \approx 32.93$$

With 33 bits we can uniquely identify every person on the planet (today).

A problem with Hartley's information measure

EXAMPLE 6: we have two hats with indistinguishable black and white balls. There are 4 balls total in each hat.



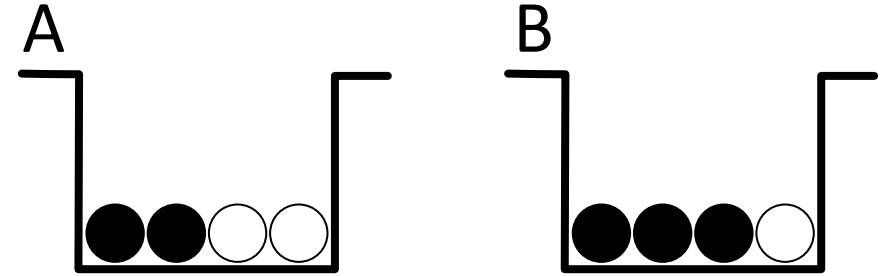
We randomly draw a ball from both hats. Let U_A, U_B be the color of the ball.

What does Hartley's information measure tell us
(maybe let's start with U_A)



A problem with Hartley's information measure

EXAMPLE 6: we have two hats with indistinguishable black and white balls. There are 4 balls total in each hat.



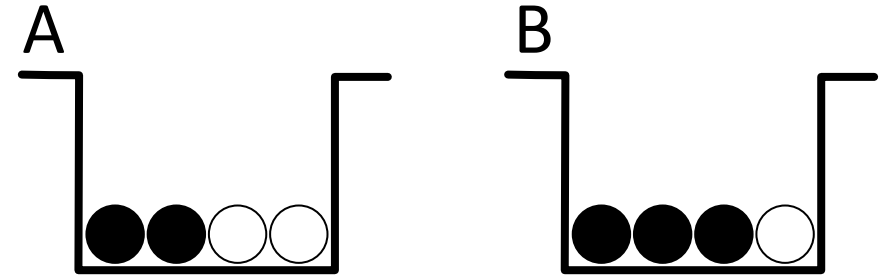
We randomly draw a ball from both hats. Let U_A, U_B be the color of the ball.

$$H_0(U_A) = \lg(2) = 1 \text{ bit} \quad (\text{we have 2 equally likely colors})$$

$$H_0(U_B) = ?$$

A problem with Hartley's information measure

EXAMPLE 6: we have two hats with indistinguishable black and white balls. There are 4 balls total in each hat.



We randomly draw a ball from both hats. Let U_A, U_B be the color of the ball.

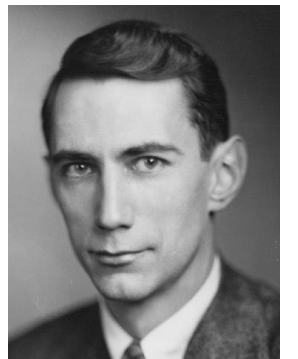
$$H_0(U_A) = \lg(2) = 1 \text{ bit}$$

$$H_0(U_B) = \lg(2) = 1 \text{ bit}$$

Problem: if $U = \text{black}$, then we get less information from U_B than from U_A (since we somehow expected that outcome)

⇒ 3) A proper measure of information should take into account the (possibly different) probabilities of the various outcomes.

This was the key insight of Claude Shannon [1948]

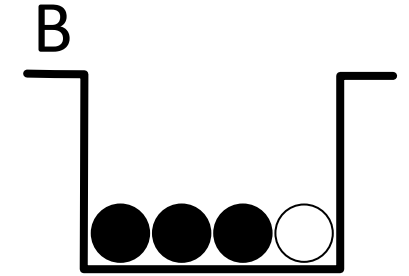


"Fixing" Hartley's information measure

Let's analyze the possible outcomes for U_B :

$U_B = \text{white}$:

What does Hartley tell us about the information we get after learning $U_B = \text{white}$?



"Fixing" Hartley's information measure

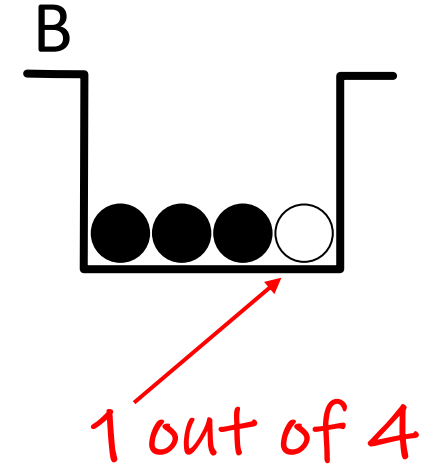
Let's analyze the possible outcomes for U_B :

$U_B = \text{white}$:

There is a $p = 1/4$ chance to draw a white ball.

That's the result of 1 out of $n = 4$ possible outcomes.

$$H_0(U_B) = ??? \quad ?$$



"Fixing" Hartley's information measure

Let's analyze the possible outcomes for U_B :

$U_B = \text{white}$:

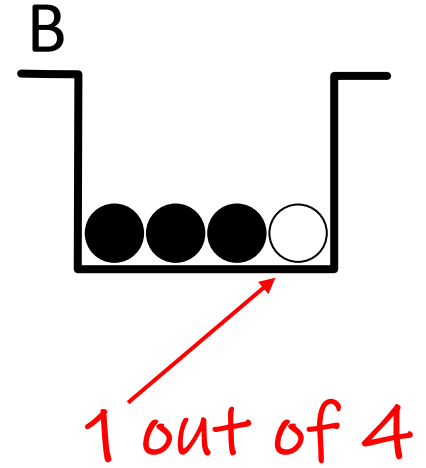
There is a $p = 1/4$ chance to draw a white ball.

That's the result of 1 out of $n = 4$ possible outcomes.

$$H_0(U_B) = \lg(4) = 2 \text{ bits}$$

$U_B = \text{black}$: $\lg\left(\frac{1}{p}\right)$

Hartley does not work directly.
What can we do?



"Fixing" Hartley's information measure

Let's analyze the possible outcomes for U_B :

$U_B = \text{white}$:

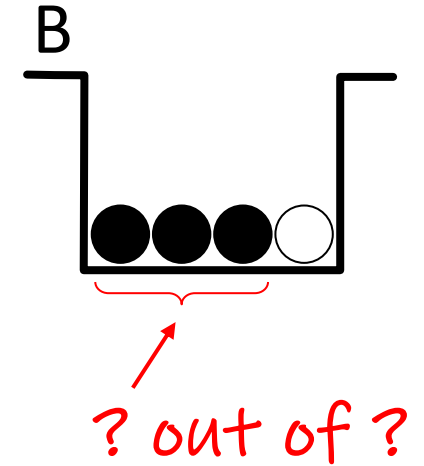
There is a $p = 1/4$ chance to draw a white ball.

That's the result of 1 out of $n = 4$ possible outcomes.

$$H_0(U_B) = \lg(4) = 2 \text{ bits}$$

$U_B = \text{black}$: $\lg\left(\frac{1}{p}\right)$

What is our chance p to draw a black ball? **?**



"Fixing" Hartley's information measure

Let's analyze the possible outcomes for U_B :

$U_B = \text{white}$:

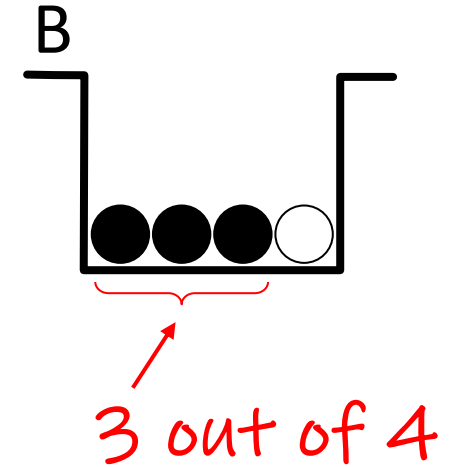
There is a $p = 1/4$ chance to draw a white ball.

That's the result of 1 out of $n = 4$ possible outcomes.

$$H_0(U_B) = \lg(4) = 2 \text{ bits}$$

$U_B = \text{black}$: $\lg\left(\frac{1}{p}\right)$

There is a $p = 3/4$ chance to draw a black ball.



What do we do with the $3/4$? ?

"Fixing" Hartley's information measure

Let's analyze the possible outcomes for U_B :

$U_B = \text{white}$:

There is a $p = 1/4$ chance to draw a white ball.

That's the result of 1 out of $n = 4$ possible outcomes.

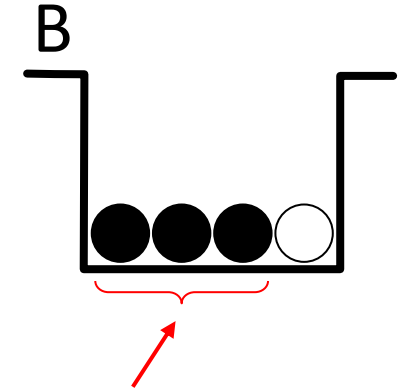
$$H_0(U_B) = \lg(4) = 2 \text{ bits}$$

$U_B = \text{black}$: $\lg\left(\frac{1}{p}\right)$

There is a $p = 3/4$ chance to draw a black ball.

That's the result of 1 out of $n = 4/3$ possible outcomes.

$$H_0(U_B) = \text{?}$$



3 out of 4
= 1 out of 4/3

For Hartley, we need to have 1 black ball (and have "1 out of r outcomes"). We get this by normalizing, i.e. dividing by 3...

"Fixing" Hartley's information measure

Let's analyze the possible outcomes for U_B :

$U_B = \text{white}$:

There is a $p = 1/4$ chance to draw a white ball.

That's the result of 1 out of $n = 4$ possible outcomes.

$$H_0(U_B) = \lg(4) = 2 \text{ bits}$$

$U_B = \text{black}$:

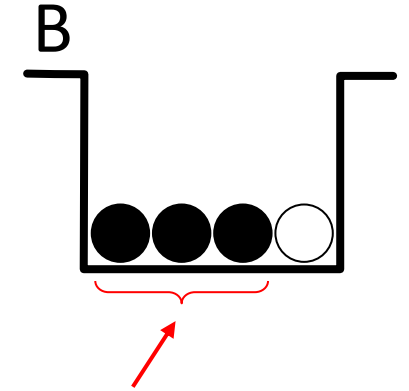
$$\lg\left(\frac{1}{p}\right)$$

#total balls /
#black balls

There is a $p = 3/4$ chance to draw a black ball.

That's the result of 1 out of $n = 4/3$ possible outcomes.

$$H_0(U_B) = \log_2\left(\frac{4}{3}\right) = 0.415 \text{ bits}$$



3 out of 4
= 1 out of 4/3

How do we combine these two possible outcomes to get one measure

?

"Fixing" Hartley's information measure

Let's analyze the possible outcomes for U_B :

$U_B = \text{white}$:

There is a $p = 1/4$ chance to draw a white ball.

That's the result of 1 out of $n = 4$ possible outcomes.

$$H_0(U_B) = \lg(4) = 2 \text{ bits}$$

$U_B = \text{black}$: $\lg\left(\frac{1}{p}\right)$

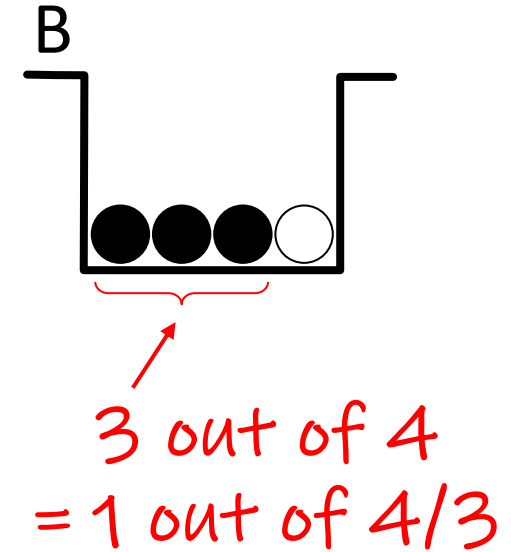
There is a $p = 3/4$ chance to draw a black ball.

That's the result of 1 out of $n = 4/3$ possible outcomes.

$$H_0(U_B) = \lg\left(\frac{4}{3}\right) = 0.415 \text{ bits}$$

Let's do "in expectation" 😊

$$\mathbb{E}[H_0(U_B)] = \frac{1}{4} \cdot \dots + \frac{3}{4} \cdot \dots$$



"Fixing" Hartley's information measure

Let's analyze the possible outcomes for U_B :

$U_B = \text{white}$:

There is a $p = 1/4$ chance to draw a white ball.

That's the result of 1 out of $n = 4$ possible outcomes.

$$H_0(U_B) = \lg(4) = 2 \text{ bits}$$

$U_B = \text{black}$: $\lg\left(\frac{1}{p}\right)$

There is a $p = 3/4$ chance to draw a black ball.

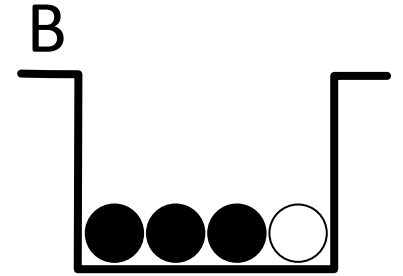
That's the result of 1 out of $n = 4/3$ possible outcomes.

$$H_0(U_B) = \lg\left(\frac{4}{3}\right) = 0.415 \text{ bits}$$

That's our expected amount of information we learn.

Let's do "in expectation":

$$\mathbb{E}[H_0(U_B)] = \frac{1}{4} \cdot 2 \text{ bits} + \frac{3}{4} \cdot 0.415 \text{ bits} = 0.811 \text{ bits}$$



"Fixing" Hartley's information measure

Let's analyze the possible outcomes:

$U_B = \text{white}$:

There is a $p = 1/4$ chance to draw a white ball.

That's the result of 1 out of $n = 4$ possible outcomes.

$$H_0(U_B) = \lg(4) = 2 \text{ bits}$$

$U_B = \text{black}$: $\lg\left(\frac{1}{p}\right)$

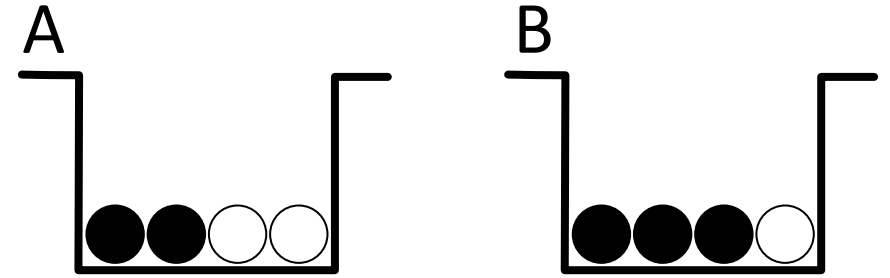
There is a $p = 3/4$ chance to draw a black ball.

That's the result of 1 out of $n = 4/3$ possible outcomes.

$$H_0(U_B) = \lg\left(\frac{4}{3}\right) = 0.415 \text{ bits}$$

Let's do "in expectation":

$$\mathbb{E}[H_0(U_B)] = \frac{1}{4} \cdot 2 \text{ bits} + \frac{3}{4} \cdot 0.415 \text{ bits} = 0.811 \text{ bits}$$



What would we get for hat A instead of hat B ?

"Fixing" Hartley's information measure

Let's analyze the possible outcomes:

$U_B = \text{white}$:

There is a $p = 1/4$ chance to draw a white ball.

That's the result of 1 out of $n = 4$ possible outcomes.

$$H_0(U_B) = \lg(4) = 2 \text{ bits}$$

$U_B = \text{black}$: $\lg\left(\frac{1}{p}\right)$

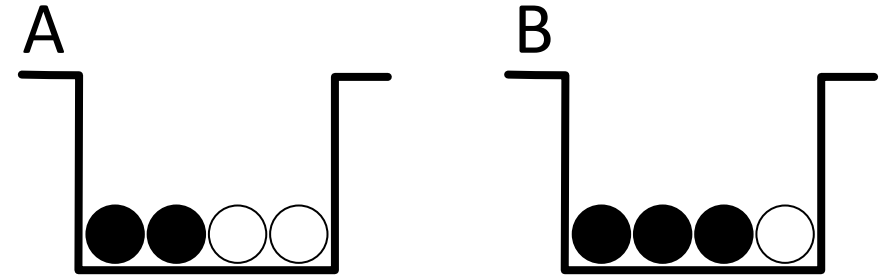
There is a $p = 3/4$ chance to draw a black ball.

That's the result of 1 out of $n = 4/3$ possible outcomes.

$$H_0(U_B) = \lg\left(\frac{4}{3}\right) = 0.415 \text{ bits}$$

Let's do "in expectation":

$$\mathbb{E}[H_0(U_B)] = \frac{1}{4} \cdot 2 \text{ bits} + \frac{3}{4} \cdot 0.415 \text{ bits} = 0.811 \text{ bits}$$



Notice that 1 bit was the min unit of information for the Hartley measure. Expectation allowed us to go lower!

1 bit for hat A

0.811 bits hat B

"Fixing" Hartley's information measure

Let's analyze the possible outcomes:

$U_B = \text{white}$:

There is a $p = 1/4$ chance to draw a white ball.

That's the result of 1 out of $n = 4$ possible outcomes.

$$H_0(U_B) = \lg(4) = 2 \text{ bits}$$

$U_B = \text{black}$: $\lg\left(\frac{1}{p}\right)$

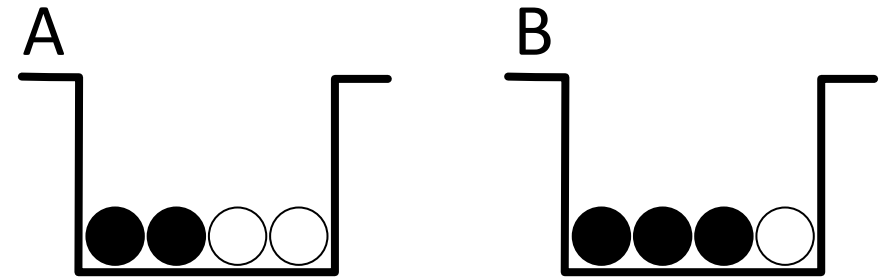
There is a $p = 3/4$ chance to draw a black ball.

That's the result of 1 out of $n = 4/3$ possible outcomes.

$$H_0(U_B) = \lg\left(\frac{4}{3}\right) = 0.415 \text{ bits}$$

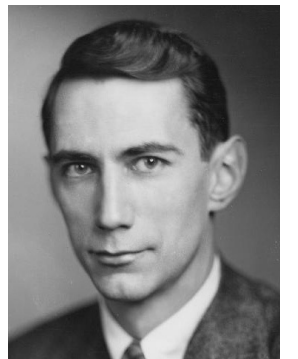
Let's do "in expectation":

$$\mathbb{E}[H_0(U_B)] = \frac{1}{4} \cdot \lg(4) + \frac{3}{4} \cdot \lg\left(\frac{4}{3}\right)$$



*This is Claude Shannon's
measure of information*

1 bit for hat A
= 0.811 bits hat B



Shannon's entropy

Shannon's measure of information as expected Hartley information (averaged over all possible outcomes)

$$H(\mathbf{p}) = \sum_{i=1}^r p_i \cdot \lg\left(\frac{1}{p_i}\right) = - \sum_{i=1}^r p_i \cdot \lg(p_i) = \mathbb{E} \left[\lg\left(\frac{1}{p_i}\right) \right]$$

$H_0(U)$

p_i = probability of the i -th possible outcome

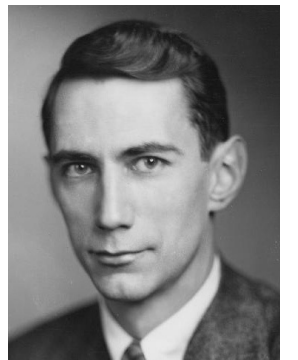
Uncertainty: Normalized number of outcomes, for option i to be "1 out of ... outcomes"

1948:

A Mathematical Theory of Communication

By C. E. SHANNON

$$H = -K \sum_{i=1}^n p_i \log p_i$$



Shannon's entropy

Shannon's measure of information as expected Hartley information (averaged over all possible outcomes)

$$H(\mathbf{p}) = \sum_{i=1}^r p_i \cdot \lg\left(\frac{1}{p_i}\right) = - \sum_{i=1}^r p_i \cdot \lg(p_i) = \mathbb{E} \left[\lg\left(\frac{1}{p_i}\right) \right]$$

$H_0(U)$

p_i = probability of the i -th possible outcome

Uncertainty: Normalized number of outcomes, for option i to be "1 out of ... outcomes"

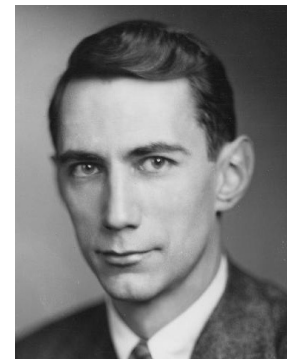
1928: **Transmission of Information**
By R. V. L. HARTLEY

$$H = Kn,$$
$$H = n \log s$$



1948: **A Mathematical Theory of Communication**
By C. E. SHANNON

$$H = -K \sum_{i=1}^n p_i \log p_i$$



Ralph Hartley. Transmission of information, The Bell System Technical Journal, 1928. <https://doi.org/10.1002/j.1538-7305.1928.tb01236.x>

Claude Shannon. A Mathematical Theory of Communication, The Bell System Technical Journal, 1948. <https://doi.org/10.1002/j.1538-7305.1948.tb01338.x>

Gatterbauer, Aslam. Foundations and Applications of Information Theory: <https://northeastern-datalab.github.io/cs7840/>

Shannon's entropy

Shannon's measure of information as expected Hartley information (averaged over all possible outcomes)

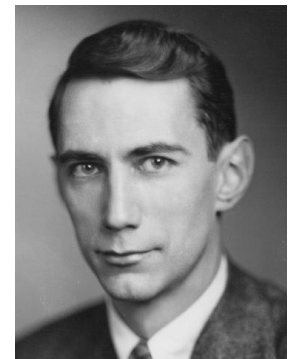
$$H(\mathbf{p}) = \sum_{i=1}^r p_i \cdot \lg\left(\frac{1}{p_i}\right) = - \sum_{i=1}^r p_i \cdot \lg(p_i) = \mathbb{E} \left[\lg\left(\frac{1}{p_i}\right) \right]$$

$H_0(U)$

p_i = probability of the i -th possible outcome

Uncertainty: Normalized number of outcomes, for option i to be "1 out of ... outcomes"

- 1) The **number of possible outcomes** should be linked to "information" H_0
- 2) Information is **additive** in some sense H_0
- 3) A proper measure of information should take into account the **different probabilities of the outcomes.** H



Ralph Hartley. Transmission of information, The Bell System Technical Journal, 1928. <https://doi.org/10.1002/j.1538-7305.1928.tb01236.x>

Claude Shannon. A Mathematical Theory of Communication, The Bell System Technical Journal, 1948. <https://doi.org/10.1002/j.1538-7305.1948.tb01338.x>

Gatterbauer, Aslam. Foundations and Applications of Information Theory: <https://northeastern-datalab.github.io/cs7840/>

Part 1: Theory

L05: Basics of entropy (2/6)

[conditional entropy, binary entropy, maximum entropy]

Wolfgang Gatterbauer, Javed Aslam

cs7840 Foundations and Applications of Information Theory (fa24)

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

9/18/2024

Pre-class conversations

- Last class recapitulation
- To be posted: Online Python notebook (feedback **very** welcome, also possibly useful for your own scribes)
- Feedback on newly posted scribes on Piazza over weekend
- Any feedback on organization on course website (Canvas, Piazza)?

- Today:
 - Keep pen & paper ready for hands-on calculus, logarithm
 - also see Schneider's "Information Theory Primer, With an Appendix on Logarithms"
 - Intuition behind entropy (and variants)

Properties of information (entropy) by example

Shannon entropy for unbiased outcomes

EXAMPLE 1: What is the entropy in a roll of an unbiased 8-sided die?



$$H(\mathbf{p}) = \sum_{i=1}^r p_i \cdot \lg\left(\frac{1}{p_i}\right)$$

Shannon entropy for unbiased outcomes

EXAMPLE 1: What is the entropy in a roll of an unbiased 8-sided die?



$$H(\mathbf{p}) = \boxed{\sum_{i=1}^r p_i \cdot \lg\left(\frac{1}{p_i}\right)} = \underbrace{\left(\sum_{i=1}^r p_i\right)}_1 \cdot \lg\left(\frac{1}{p_i}\right) = \lg\left(\frac{1}{p_i}\right) \text{ ?}$$

Shannon entropy for unbiased outcomes = Hartley measure

EXAMPLE 1: What is the entropy in a roll of an unbiased 8-sided die?



$$H(\mathbf{p}) = \sum_{i=1}^r p_i \cdot \lg\left(\frac{1}{p_i}\right) = \underbrace{\left(\sum_{i=1}^r p_i\right)}_1 \cdot \lg\left(\frac{1}{p_i}\right) = \lg\left(\frac{1}{p_i}\right) = H_0\left(\frac{1}{p_i}\right)$$

number of outcomes

Entropy is exactly the Hartley information measure for unbiased outcomes 😊

Shannon entropy for unbiased outcomes = Hartley measure

EXAMPLE 1: What is the entropy in a roll of an unbiased 8-sided die?



$$H(\mathbf{p}) = \sum_{i=1}^r p_i \cdot \lg\left(\frac{1}{p_i}\right) = \underbrace{\left(\sum_{i=1}^r p_i\right)}_1 \cdot \lg\left(\frac{1}{p_i}\right) = \lg\left(\frac{1}{p_i}\right) = H_0\left(\frac{1}{p_i}\right) = \lg(8) = 3$$

number of outcomes

Entropy is exactly the Hartley information measure for unbiased outcomes 😊

Characterization of the Hartley information measure

Shannon entropy for uniform sampling from n choices.

$$H_0(r) = H_0\left(\frac{1}{p_i}\right) = \lg(n)$$

two independent uniformly distributed Rvs,
with alphabet size m and n

The Hartley function only depends on the number of elements in a set, and hence can be viewed as a function on natural numbers. Rényi showed that the Hartley function in base 2 is the only function mapping natural numbers to real numbers that satisfies

1. $H_0(mn) = H_0(m) + H_0(n)$ (additivity) $\lg(m \cdot n) = \lg(m) + \lg(n)$
2. $H_0(m) \leq H_0(m + 1)$ (monotonicity)
3. $H_0(2) = 1$ (normalization)

Condition 1 says that the uncertainty of the Cartesian product of two finite sets A and B is the sum of uncertainties of A and B . Condition 2 says that a larger set has larger uncertainty.

Learning partial information

EXAMPLE 2: We roll an unbiased 8-sided die with sides $\{1, 2, \dots, 8\}$.



We then get a message with the information that the outcome of a roll is even.

How much information did we learn? **?**

Learning partial information

EXAMPLE 2: We roll an unbiased 8-sided die with sides $\{1, 2, \dots, 8\}$.



We then get a message with the information that the outcome of a roll is even. How much information did we learn?

- **Before** the message: ?
- **After** the message: ?

Learning partial information

EXAMPLE 2: We roll an unbiased 8-sided die with sides $\{1, 2, \dots, 8\}$.



We then get a message with the information that the outcome of a roll is even. How much information did we learn?

- Before the message: There are 8 choices: $\{1, 2, 3, 4, 5, 6, 7, 8\}$
- After the message: There are 4 choices: $\{2, 4, 6, 8\}$

How much information did we have before?
How much information did we have after



Learning partial information

EXAMPLE 2: We roll an unbiased 8-sided die with sides $\{1, 2, \dots, 8\}$.



We then get a message with the information that the outcome of a roll is even. How much information did we learn?

- Before the message: There are 8 choices: $\{1, 2, 3, 4, 5, 6, 7, 8\}$ $H_0(8) = 3$ bits
- After the message: There are 4 choices: $\{2, 4, 6, 8\}$ $H_0(4) = 2$ bits

Let's think about encodings *(binary encoding with atypical 1-indexing)*

Before: $\{1, 2, 3, 4, 5, 6, 7, 8\}$
 ↓ ↓ ↓ ↓ ↓ ↓ ↓ ↓
 000 001 010 011 100 101 110 111

After: $\{1, 2, 3, 4, 5, 6, 7, 8\}$
 ↓ ↓ ↓ ↓ ↓ ↓ ↓ ↓
 000 001 010 011 100 101 110 111

Do you notice something



Learning partial information

EXAMPLE 2: We roll an unbiased 8-sided die with sides $\{1, 2, \dots, 8\}$.



We then get a message with the information that the outcome of a roll is even. How much information did we learn?

- Before the message: There are 8 choices: $\{1, 2, 3, 4, 5, 6, 7, 8\}$ $H_0(8) = 3$ bits
- After the message: There are 4 choices: $\{2, 4, 6, 8\}$ $H_0(4) = 2$ bits

Let's think about encodings

Before: $\{1, 2, 3, 4, 5, 6, 7, 8\}$
 ↓ ↓ ↓ ↓ ↓ ↓ ↓ ↓
 000 001 010 011 100 101 110 111

After: $\{1, 2, 3, 4, 5, 6, 7, 8\}$
 ↓ ↓ ↓ ↓ ↓ ↓ ↓ ↓
 000 00**1** 010 01**1** 100 10**1** 110 11**1**

We have learned 1 bit! **??1**

"Grouping rule": Dividing the outcomes into two (last bit), randomly choose one group (e.g. 1), and then randomly pick an element from that group (e.g. 10), does not change the entropy

Learning partial information

EXAMPLE 2: We roll an unbiased 8-sided die with sides $\{1, 2, \dots, 8\}$.



We then get a message with the information that the outcome of a roll is even. How much information did we learn?

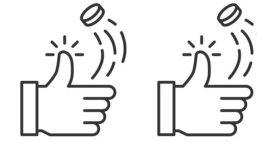
- Before the message: There are 8 choices: $\{1, 2, 3, 4, 5, 6, 7, 8\}$ $H_0(8) = 3$ bits
- After the message: There are 4 choices: $\{2, 4, 6, 8\}$ $H_0(4) = 2$ bits

Recall: **information is additive**:

1 flip of a 2-sided coin has 2 outcomes.



2 flips have $2^2 = 4$ outcomes.



3 flips have $2^3 = 8$ outcomes.



$$\lg(2) = 1$$

$$\lg(4) = 2$$

$$\lg(8) = 3$$

+1 bit

+1 bit

Learning partial information

EXAMPLE 2: We roll an unbiased 8-sided die with sides $\{1, 2, \dots, 8\}$.



We then get a message with the information that the outcome of a roll is even. How much information did we learn?

- Before the message: There are 8 choices: $\{1, 2, 3, 4, 5, 6, 7, 8\}$ $H_0(8) = 3$ bits
- After the message: There are 4 choices: $\{2, 4, 6, 8\}$ $H_0(4) = 2$ bits

The power of the logarithm: transform multiplication into addition

Uncertainty before – Uncertainty after

$$\underbrace{\lg(8) - \lg(4)}$$

$$\lg\left(\frac{8}{4}\right) = \lg(2) = 1 \text{ bit}$$

$$H(U) = \boxed{\lg\left(\frac{n}{m}\right)}$$

Information content in a message U that reduces the number of unbiased outcomes from n to m

Learning partial information

EXAMPLE 3: We roll an unbiased 8-sided die with sides $\{1, 2, \dots, 8\}$.



We then get 4 messages, one after the other: $U_1 =$ "The outcome of the roll is not **1**", $U_2 =$ "... not **3**", $U_3 =$ "... not **5**", $U_4 =$ "... not **7**".

How much information do we learn from each individual message?



$$H(U_1) = ?$$

$$H(U_2|U_1) = ? \quad \text{These are called "conditional entropies!"}$$

$$H(U_3|U_{1,2}) = ?$$

$$H(U_4|U_{1-3}) = ?$$

Learning partial information

EXAMPLE 3: We roll an unbiased 8-sided die with sides $\{1, 2, \dots, 8\}$.



We then get 4 messages, one after the other: $U_1 =$ "The outcome of the roll is not **1**", $U_2 =$ "... not **3**", $U_3 =$ "... not **5**", $U_4 =$ "... not **7**".

How much information do we learn from each individual message?

$$H(U_1) = \lg\left(\frac{8}{7}\right) = 0.193 \text{ bits}$$

$$H(U_2|U_1) = ?$$

$$H(U_3|U_{1,2}) = ?$$

$$H(U_4|U_{1-3}) = ?$$

Learning partial information

EXAMPLE 3: We roll an unbiased 8-sided die with sides $\{1, 2, \dots, 8\}$.



We then get 4 messages, one after the other: $U_1 =$ "The outcome of the roll is not **1**", $U_2 =$ "... not **3**", $U_3 =$ "... not **5**", $U_4 =$ "... not **7**".

How much information do we learn from each individual message?

... and all of them together?



$$H(U_1) = \lg\left(\frac{8}{7}\right) = 0.193 \text{ bits}$$

$$H(U_2|U_1) = \lg\left(\frac{7}{6}\right) = 0.222 \text{ bits}$$

$$H(U_3|U_{1,2}) = \lg\left(\frac{6}{5}\right) = 0.263 \text{ bits}$$

$$H(U_4|U_{1-3}) = \lg\left(\frac{5}{4}\right) = 0.322 \text{ bits}$$

Learning partial information

EXAMPLE 3: We roll an unbiased 8-sided die with sides $\{1, 2, \dots, 8\}$.



We then get 4 messages, one after the other: $U_1 =$ "The outcome of the roll is not **1**", $U_2 =$ "... not **3**", $U_3 =$ "... not **5**", $U_4 =$ "... not **7**".

How much information do we learn from each individual message?
... and all of them together?

$$H(U_1) = \lg\left(\frac{8}{7}\right) = 0.193 \text{ bits}$$

$$H(U_2|U_1) = \lg\left(\frac{7}{6}\right) = 0.222 \text{ bits}$$

$$H(U_3|U_{1,2}) = \lg\left(\frac{6}{5}\right) = 0.263 \text{ bits}$$

$$H(U_4|U_{1-3}) = \lg\left(\frac{5}{4}\right) = 0.322 \text{ bits}$$

$$H(\{U_1, U_2, U_3, U_4\}) = \mathbf{1 \text{ bit}}$$

*How come that the **SUM** of these numbers turns out to be sooooo nice?*



Learning partial information

EXAMPLE 3: We roll an unbiased 8-sided die with sides $\{1, 2, \dots, 8\}$.



We then get 4 messages, one after the other: $U_1 =$ "The outcome of the roll is not **1**", $U_2 =$ "... not **3**", $U_3 =$ "... not **5**", $U_4 =$ "... not **7**".

How much information do we learn from each individual message?

... and all of them together?

$$H(U_1) = \lg\left(\frac{8}{7}\right) = 0.193 \text{ bits}$$

$$H(U_2|U_1) = \lg\left(\frac{7}{6}\right) = 0.222 \text{ bits}$$

$$H(U_3|U_{1,2}) = \lg\left(\frac{6}{5}\right) = 0.263 \text{ bits}$$

$$H(U_4|U_{1-3}) = \lg\left(\frac{5}{4}\right) = 0.322 \text{ bits}$$

$$H(\{U_1, U_2, U_3, U_4\}) = 1 \text{ bit}$$

$$H(\{U_1, U_2, U_3, U_4\})$$

$$= H(U_1) + H(U_2|U_1) + H(U_3|U_{1,2}) + H(U_4|U_{1-3})$$

This is called the "chain rule"

Learning partial information

EXAMPLE 3: We roll an unbiased 8-sided die with sides $\{1, 2, \dots, 8\}$.



We then get 4 messages, one after the other: $U_1 =$ "The outcome of the roll is not **1**", $U_2 =$ "... not **3**", $U_3 =$ "... not **5**", $U_4 =$ "... not **7**".

How much information do we learn from each individual message?

... and all of them together?

$$H(U_1) = \lg\left(\frac{8}{7}\right) = 0.193 \text{ bits}$$

$$H(U_2|U_1) = \lg\left(\frac{7}{6}\right) = 0.222 \text{ bits}$$

$$H(U_3|U_{1,2}) = \lg\left(\frac{6}{5}\right) = 0.263 \text{ bits}$$

$$H(U_4|U_{1-3}) = \lg\left(\frac{5}{4}\right) = 0.322 \text{ bits}$$

$$H(\{U_1, U_2, U_3, U_4\}) = 1 \text{ bit}$$

$$H(\{U_1, U_2, U_3, U_4\})$$

$$= H(U_1) + H(U_2|U_1) + H(U_3|U_{1,2}) + H(U_4|U_{1-3})$$

$$= \lg\left(\frac{8}{7}\right) + \lg\left(\frac{7}{6}\right) + \lg\left(\frac{6}{5}\right) + \lg\left(\frac{5}{4}\right)$$

= ?

Learning partial information

EXAMPLE 3: We roll an unbiased 8-sided die with sides $\{1, 2, \dots, 8\}$.



We then get 4 messages, one after the other: $U_1 =$ "The outcome of the roll is not **1**", $U_2 =$ "... not **3**", $U_3 =$ "... not **5**", $U_4 =$ "... not **7**".

How much information do we learn from each individual message?

... and all of them together?

$$H(U_1) = \lg\left(\frac{8}{7}\right) = 0.193 \text{ bits}$$

$$H(U_2|U_1) = \lg\left(\frac{7}{6}\right) = 0.222 \text{ bits}$$

$$H(U_3|U_{1,2}) = \lg\left(\frac{6}{5}\right) = 0.263 \text{ bits}$$

$$H(U_4|U_{1-3}) = \lg\left(\frac{5}{4}\right) = 0.322 \text{ bits}$$

$$H(\{U_1, U_2, U_3, U_4\}) = 1 \text{ bit}$$

$$H(\{U_1, U_2, U_3, U_4\})$$

$$= H(U_1) + H(U_2|U_1) + H(U_3|U_{1,2}) + H(U_4|U_{1-3})$$

$$= \lg\left(\frac{8}{7}\right) + \lg\left(\frac{7}{6}\right) + \lg\left(\frac{6}{5}\right) + \lg\left(\frac{5}{4}\right)$$

$$= \lg\left(\frac{8}{\cancel{7}} \cdot \frac{\cancel{7}}{\cancel{6}} \cdot \frac{\cancel{6}}{\cancel{5}} \cdot \frac{\cancel{5}}{4}\right) = \lg\left(\frac{8}{4}\right) = 1 \text{ bit}$$

Again, the logarithm 😊

Maximum Entropy distributions

Binary Entropy Function

X is a Bernoulli RV with $p(x) = \begin{cases} p, & \text{if } x = 1 \\ 1 - p, & \text{if } x = 0 \end{cases}$

Biased coin flip:



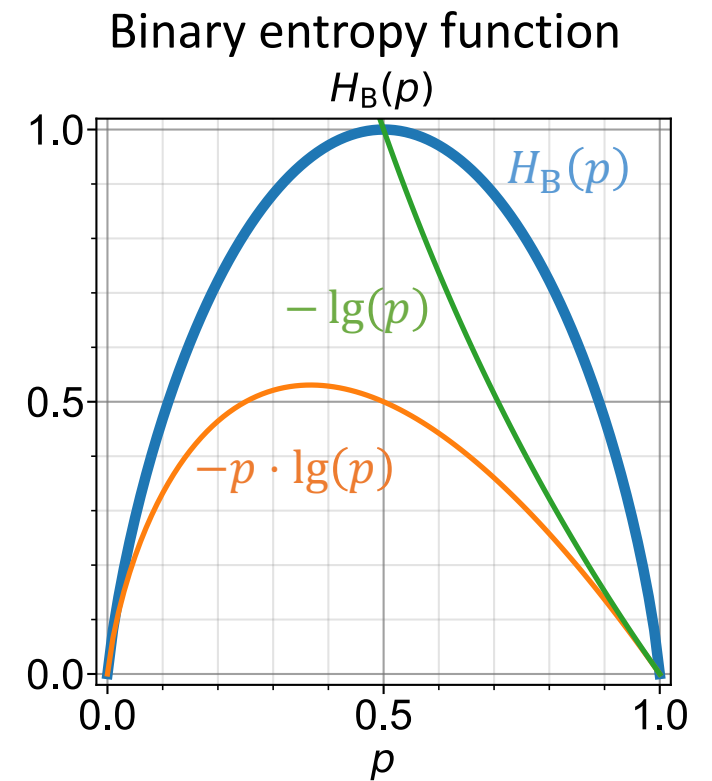
$$H_B(p) = ?$$

Binary Entropy Function


X is a Bernoulli RV with $p(x) = \begin{cases} p, & \text{if } x = 1 \\ 1 - p, & \text{if } x = 0 \end{cases}$



$$H_B(p) = -p \cdot \lg(p) - (1 - p) \cdot \lg(1 - p)$$

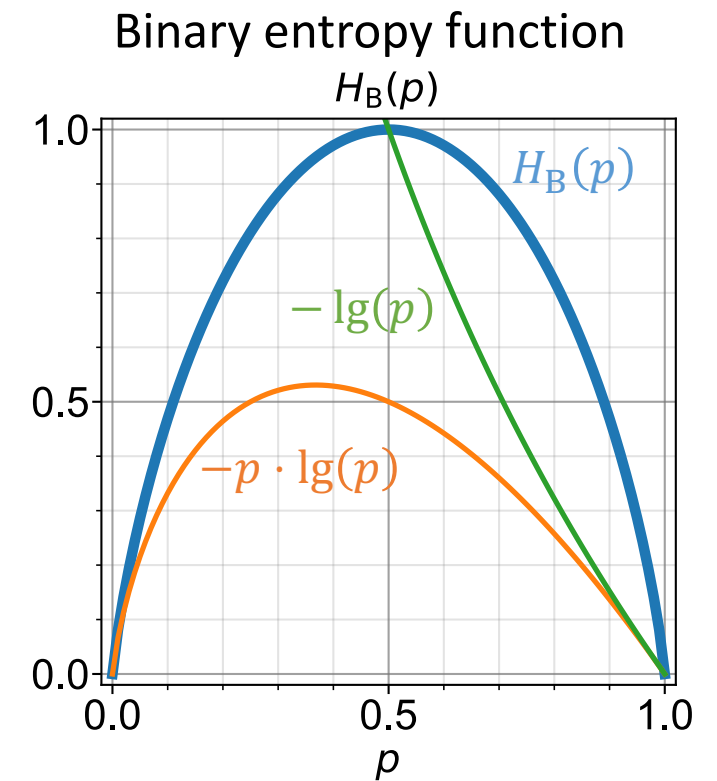


Binary Entropy Function

X is a Bernoulli RV with $p(x) = \begin{cases} p, & \text{if } x = 1 \\ 1 - p, & \text{if } x = 0 \end{cases}$ 

$$H_B(p) = -p \cdot \lg(p) - (1 - p) \cdot \lg(1 - p)$$

How to choose p in order to maximize entropy ?



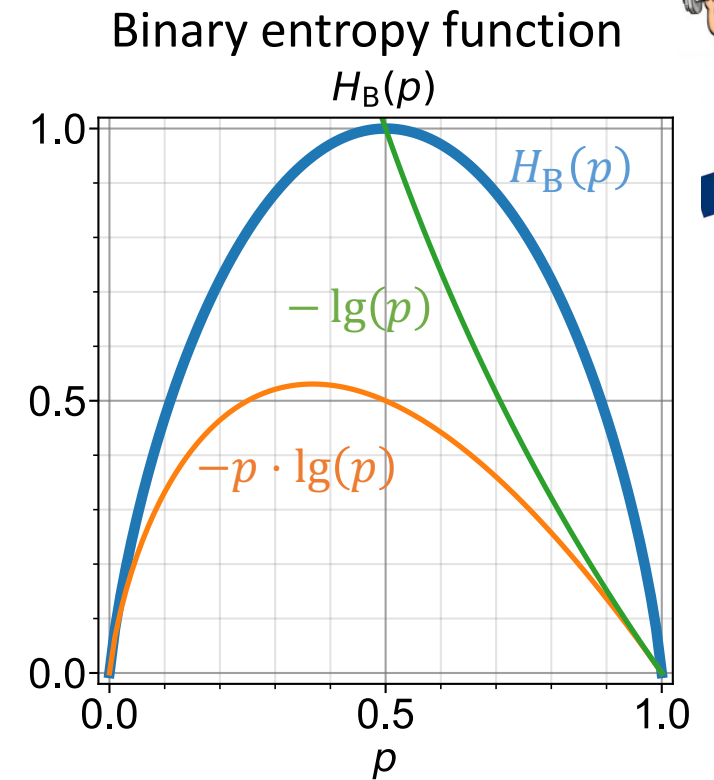
Binary Entropy Function

X is a Bernoulli RV with $p(x) = \begin{cases} p, & \text{if } x = 1 \\ 1 - p, & \text{if } x = 0 \end{cases}$

$$H_B(p) = -p \cdot \lg(p) - (1 - p) \cdot \lg(1 - p)$$

How to choose p in order to maximize entropy?

$$\frac{dH}{dp} = ?$$



Understanding "change of basis"

$$\lg(x) = \log_2(x) = \frac{\ln(x)}{\ln(2)}$$



Calculus
cheat
sheet

$$\ln(x)' =$$

$$\lg(x)' = \left(\frac{\ln(x)}{\ln(2)} \right)' =$$

$$(x \cdot \lg(x))' =$$

$$\lg(1 - x)' =$$

$$((1 - x) \cdot \lg(1 - x))' =$$



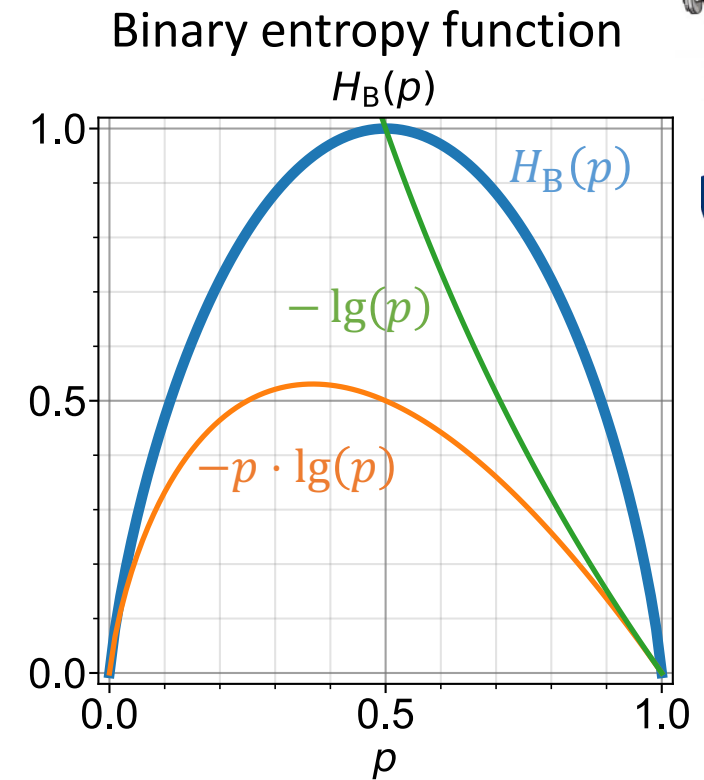
Binary Entropy Function

X is a Bernoulli RV with $p(x) = \begin{cases} p, & \text{if } x = 1 \\ 1 - p, & \text{if } x = 0 \end{cases}$

$$H_B(p) = -p \cdot \lg(p) - (1 - p) \cdot \lg(1 - p)$$

How to choose p in order to maximize entropy?

$$\frac{dH}{dp} = ?$$



Understanding "change of basis"

$$\lg(x) = \log_2(x) = \frac{\ln(x)}{\ln(2)} \quad \text{definition}$$

$$2^{\log_2(x)} = x \quad \text{apply } \ln(\dots) \text{ on both sides}$$

$$\ln(2^{\log_2(x)}) = \ln(x) \quad \ln(a^b) = b \cdot \ln(a)$$

$$\log_2(x) \cdot \ln(2) = \ln(x)$$

Calculus
cheat
sheet

$$\ln(x)' =$$

$$\lg(x)' = \left(\frac{\ln(x)}{\ln(2)}\right)' =$$

$$(x \cdot \lg(x))' =$$

$$\lg(1 - x)' =$$

$$((1 - x) \cdot \lg(1 - x))' =$$



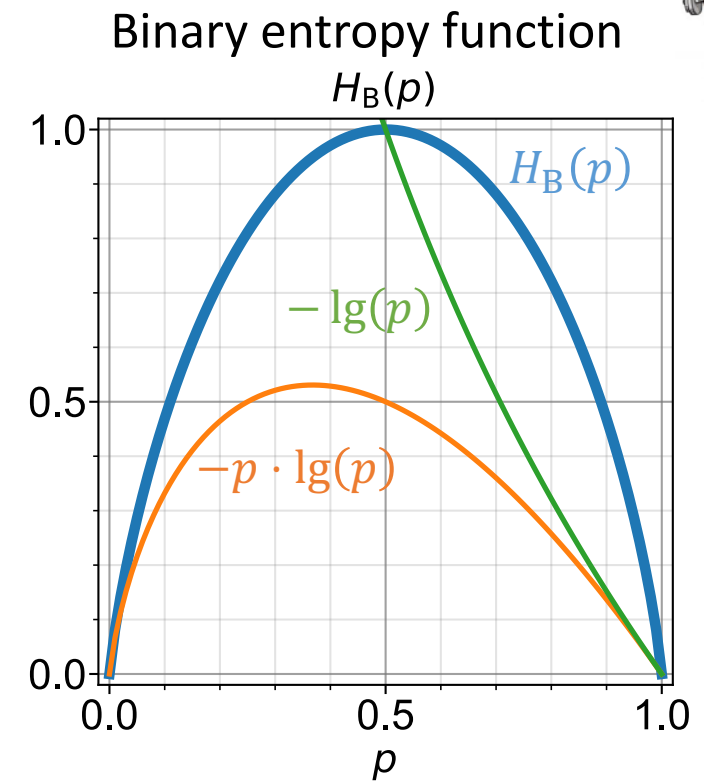
Binary Entropy Function

X is a Bernoulli RV with $p(x) = \begin{cases} p, & \text{if } x = 1 \\ 1 - p, & \text{if } x = 0 \end{cases}$

$$H_B(p) = -p \cdot \lg(p) - (1 - p) \cdot \lg(1 - p)$$

How to choose p in order to maximize entropy?

$$\frac{dH}{dp} = ?$$



Calculus cheat sheet

$$\ln(x)' = \frac{1}{x}$$

$$\lg(x)' = \left(\frac{\ln(x)}{\ln(2)}\right)' = \frac{1}{x \cdot \ln(2)}$$

$$(x \cdot \lg(x))' = \cancel{x} \frac{1}{\cancel{x} \ln(2)} + \lg(x)$$

$$\lg(1 - x)' = -\frac{1}{(1-x) \cdot \ln(2)}$$

$$((1 - x) \cdot \lg(1 - x))' = -\frac{1}{\ln(2)} - \lg(1 - x)$$

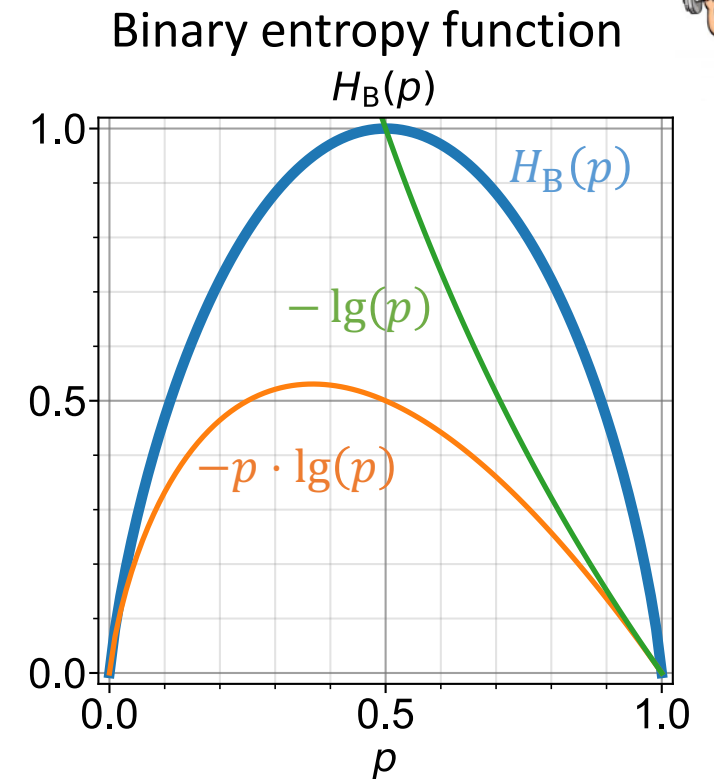
Binary Entropy Function

X is a Bernoulli RV with $p(x) = \begin{cases} p, & \text{if } x = 1 \\ 1 - p, & \text{if } x = 0 \end{cases}$

$$H_B(p) = \underbrace{-p \cdot \lg(p)}_{\text{red bracket}} - \underbrace{(1 - p) \cdot \lg(1 - p)}_{\text{red bracket}}$$

How to choose p in order to maximize entropy?

$$\frac{dH}{dp} = \underbrace{-\cancel{\frac{1}{\ln(2)}} - \lg(p)}_{\text{red bracket}} + \underbrace{\cancel{\frac{1}{\ln(2)}} + \lg(1 - p)}_{\text{red bracket}}$$



Calculus
cheat
sheet

$$\ln(x)' = \frac{1}{x}$$

$$\lg(x)' = \left(\frac{\ln(x)}{\ln(2)}\right)' = \frac{1}{x \cdot \ln(2)}$$

$$(x \cdot \lg(x))' = \cancel{x} \frac{1}{\cancel{x} \ln(2)} + \lg(x)$$

$$\lg(1 - x)' = -\frac{1}{(1-x) \cdot \ln(2)}$$

$$((1 - x) \cdot \lg(1 - x))' = -\frac{1}{\ln(2)} - \lg(1 - x)$$

Binary Entropy Function

X is a Bernoulli RV with $p(x) = \begin{cases} p, & \text{if } x = 1 \\ 1 - p, & \text{if } x = 0 \end{cases}$

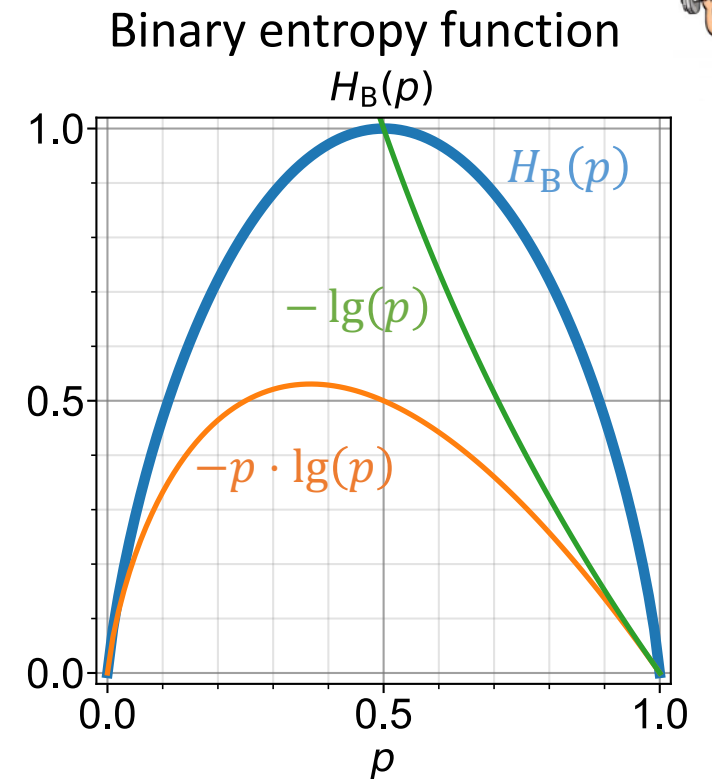
$$H_B(p) = \underbrace{-p \cdot \lg(p)}_{\text{red bracket}} - \underbrace{(1 - p) \cdot \lg(1 - p)}_{\text{red bracket}}$$

How to choose p in order to maximize entropy?

$$\frac{dH}{dp} = \cancel{\frac{1}{\ln(2)}} - \lg(p) + \cancel{\frac{1}{\ln(2)}} + \lg(1 - p) = 0$$

$$\Leftrightarrow \lg\left(\frac{1-p}{p}\right) = 0 \Leftrightarrow \left(\frac{1-p}{p}\right) = 1 \Leftrightarrow \boxed{p = \frac{1}{2}}$$

$$\frac{d^2H}{dp^2} = \text{?}$$



Calculus
cheat
sheet

$$\ln(x)' = \frac{1}{x}$$

$$\lg(x)' = \left(\frac{\ln(x)}{\ln(2)}\right)' = \frac{1}{x \cdot \ln(2)}$$

$$(x \cdot \lg(x))' = \cancel{x} \frac{1}{\cancel{x} \ln(2)} + \lg(x)$$

$$\lg(1 - x)' = -\frac{1}{(1-x) \cdot \ln(2)}$$

$$((1 - x) \cdot \lg(1 - x))' = -\frac{1}{\ln(2)} - \lg(1 - x)$$

Binary Entropy Function

X is a Bernoulli RV with $p(x) = \begin{cases} p, & \text{if } x = 1 \\ 1 - p, & \text{if } x = 0 \end{cases}$

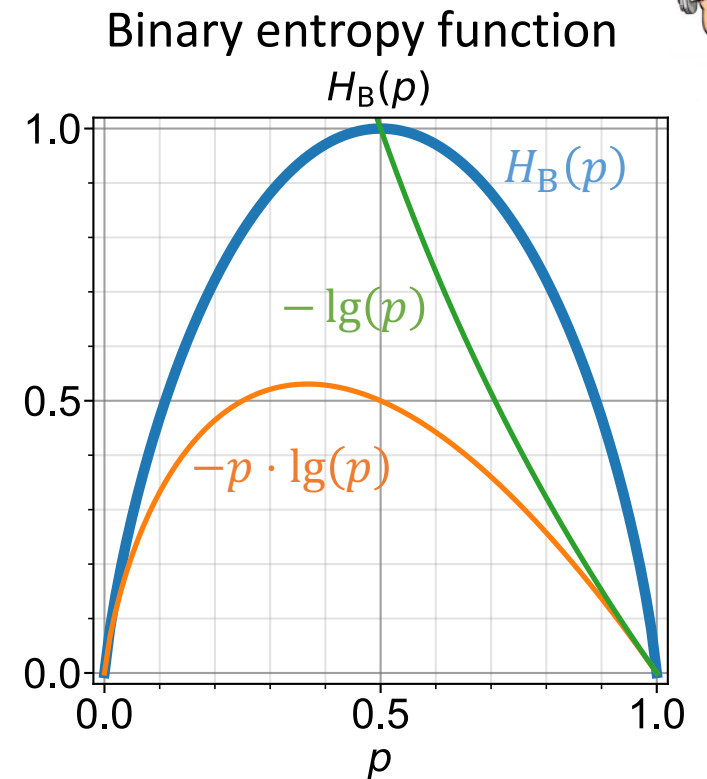
$$H_B(p) = \underbrace{-p \cdot \lg(p)}_{\text{red bracket}} - \underbrace{(1 - p) \cdot \lg(1 - p)}_{\text{red bracket}}$$

How to choose p in order to maximize entropy?

$$\frac{dH}{dp} = \cancel{\frac{1}{\ln(2)}} - \lg(p) + \cancel{\frac{1}{\ln(2)}} + \lg(1 - p) = 0$$

$$\Leftrightarrow \lg\left(\frac{1-p}{p}\right) = 0 \Leftrightarrow \left(\frac{1-p}{p}\right) = 1 \Leftrightarrow \boxed{p = \frac{1}{2}}$$

$$\frac{d^2H}{dp^2} = -\frac{1}{p \cdot \ln(2)} - \frac{1}{(1-p) \cdot \ln(2)} < 0 \text{ concave}$$



Calculus
cheat
sheet

$$\ln(x)' = \frac{1}{x}$$

$$\lg(x)' = \left(\frac{\ln(x)}{\ln(2)}\right)' = \frac{1}{x \cdot \ln(2)}$$

$$(x \cdot \lg(x))' = \cancel{x} \frac{1}{\ln(2)} + \lg(x)$$

$$\lg(1 - x)' = -\frac{1}{(1-x) \cdot \ln(2)}$$

$$((1 - x) \cdot \lg(1 - x))' = -\frac{1}{\ln(2)} - \lg(1 - x)$$

Asking ChatGPT for help

Can you calculate step-by-step the derivative for p of:

$$H(p) = -p \cdot \lg(p) - (1-p) \cdot \lg[(1-p)]$$

Repeat but assume log is to the basis 2 instead of e

Sure! Let's differentiate the function

$$H(p) = -p \cdot \log_2(p) - (1-p) \cdot \log_2(1-p),$$

where \log_2 denotes the logarithm base 2. The process is similar to before, so let's go through it step-by-step.

Step 1: Rewrite the Function

The function is already clear:

$$H(p) = -p \log_2(p) - (1-p) \log_2(1-p).$$

Step 2: Differentiate $H(p)$

We will differentiate each term separately.

Term 1: Differentiate $-p \log_2(p)$

Using the product rule:

$$\frac{d}{dp}[-p \log_2(p)] = -\left(\log_2(p) + p \cdot \frac{1}{p \ln(2)}\right) = -\log_2(p) - \frac{1}{\ln(2)}.$$

Term 2: Differentiate $-(1-p) \log_2(1-p)$

Again, using the product rule:

$$\frac{d}{dp}[-(1-p) \log_2(1-p)] = -\left(-\log_2(1-p) + (1-p) \cdot \frac{-1}{(1-p) \ln(2)}\right) = \log_2(1-p) - \frac{1}{\ln(2)}.$$

Step 3: Combine the Derivatives

Now we can combine the derivatives of both terms:

$$H'(p) = -\log_2(p) - \frac{1}{\ln(2)} + \log_2(1-p) - \frac{1}{\ln(2)}.$$

Step 4: Simplify

Combining the constant terms and rearranging gives:

$$H'(p) = \log_2(1-p) - \log_2(p) - \frac{2}{\ln(2)}.$$

This can also be expressed using properties of logarithms:

$$H'(p) = \log_2\left(\frac{1-p}{p}\right) - \frac{2}{\ln(2)}.$$

Final Result

Thus, the derivative of $H(p)$ with respect to p is:

$$H'(p) = \log_2\left(\frac{1-p}{p}\right) - \frac{2}{\ln(2)}.$$



ChatGPT made the *same* mistake as me!

Asking ChatGPT for help

Can you calculate step-by-step the derivative for p of:

$$H(p) = -p \cdot \lg(p) - (1-p) \cdot \lg[(1-p)]$$

Repeat but assume log is to the basis 2 instead of e

Sure! Let's differentiate the function

$$H(p) = -p \cdot \log_2(p) - (1-p) \cdot \log_2(1-p),$$

where \log_2 denotes the logarithm base 2. The process is similar to before, so let's go through it step-by-step.

Step 1: Rewrite the Function

The function is already clear:

$$H(p) = -p \log_2(p) - (1-p) \log_2(1-p).$$

Step 2: Differentiate $H(p)$

We will differentiate each term separately.

Term 1: Differentiate $-p \log_2(p)$

Using the product rule:

$$\frac{d}{dp}[-p \log_2(p)] = -\left(\log_2(p) + p \cdot \frac{1}{p \ln(2)}\right) = -\log_2(p) - \frac{1}{\ln(2)}.$$

Term 2: Differentiate $-(1-p) \log_2(1-p)$

Again, using the product rule:

$$\frac{d}{dp}[-(1-p) \log_2(1-p)] = -\left(-\log_2(1-p) + (1-p) \cdot \frac{-1}{(1-p) \ln(2)}\right) = \log_2(1-p) - \frac{1}{\ln(2)}.$$

Step 3: Combine the Derivatives

Now we can combine the derivatives of both terms:

$$H'(p) = -\log_2(p) - \frac{1}{\ln(2)} + \log_2(1-p) - \frac{1}{\ln(2)}.$$

Step 4: Simplify

Combining the constant terms and rearranging gives:

$$H'(p) = \log_2(1-p) - \log_2(p) - \frac{2}{\ln(2)}.$$

This can also be expressed using properties of logarithms:

$$H'(p) = \log_2\left(\frac{1-p}{p}\right) - \frac{2}{\ln(2)}.$$

Final Result

Thus, the derivative of $H(p)$ with respect to p is:

$$H'(p) = \log_2\left(\frac{1-p}{p}\right) - \frac{2}{\ln(2)}.$$

You had some error in earlier step 2 term 2. Please check the minuses

Let's go through the differentiation of the term $-(1-p) \log_2(1-p)$ again carefully to ensure we get the signs correct.

Step 1: Differentiate $-(1-p) \log_2(1-p)$

Using the product rule, we have:

$$\frac{d}{dp}[-(1-p) \log_2(1-p)] = -\left(\frac{d}{dp}(1-p) \cdot \log_2(1-p) + (1-p) \cdot \frac{d}{dp}[\log_2(1-p)]\right).$$

1. The derivative of $1-p$ is -1 .
2. The derivative of $\log_2(1-p)$ is:

$$\frac{d}{dp}[\log_2(1-p)] = -\frac{1}{(1-p) \ln(2)}.$$

Step 2: Apply the Product Rule

Putting it all together:

$$\frac{d}{dp}[-(1-p) \log_2(1-p)] = -\left(-\log_2(1-p) + (1-p) \cdot \left(-\frac{1}{(1-p) \ln(2)}\right)\right).$$

Step 3: Simplify the Expression

This simplifies to:

$$\log_2(1-p) + \frac{1}{\ln(2)}.$$

Final Result for the Derivative of $H(p)$

Now, let's combine everything correctly for the overall derivative $H'(p)$:

1. The derivative of the first term $-p \log_2(p)$:

$$-\log_2(p) - \frac{1}{\ln(2)}.$$

2. The corrected derivative of the second term $-(1-p) \log_2(1-p)$:

$$\log_2(1-p) + \frac{1}{\ln(2)}.$$

Combining these:

$$H'(p) = \left(-\log_2(p) - \frac{1}{\ln(2)}\right) + \left(\log_2(1-p) + \frac{1}{\ln(2)}\right).$$

The $-\frac{1}{\ln(2)}$ and $+\frac{1}{\ln(2)}$ cancel out, leading to:

$$H'(p) = \log_2(1-p) - \log_2(p).$$

Using properties of logarithms:

$$H'(p) = \log_2\left(\frac{1-p}{p}\right).$$

Final Result

Thus, the corrected derivative of $H(p)$ is:

$$H'(p) = \log_2\left(\frac{1-p}{p}\right).$$

Thank you for your patience, and I apologize for the confusion earlier!

Maximum Entropy distribution for nonbinary discrete RV

Entropy of a discrete distribution with n outcomes $H(\mathbf{p}) = -\sum_{i=1}^n p_i \cdot \lg(p_i)$

Which choice of $\mathbf{p} = (p_1, \dots, p_n)$ maximizes the entropy?



Maximum Entropy distribution for nonbinary discrete RV

Entropy of a discrete distribution with n outcomes $H(\mathbf{p}) = - \sum_{i=1}^n p_i \cdot \lg(p_i)$

Which choice of $\mathbf{p} = (p_1, \dots, p_n)$ maximizes the entropy?

Can be solved with constrained optimization:

$$\max_{\mathbf{p}} [H(\mathbf{p})] \text{ s.t. } \sum_{i=1}^n p_i = 1$$

Form the **Lagrangian**:

?

Maximum Entropy distribution for nonbinary discrete RV

Entropy of a discrete distribution with n outcomes $H(\mathbf{p}) = -\sum_{i=1}^n p_i \cdot \lg(p_i)$

Which choice of $\mathbf{p} = (p_1, \dots, p_n)$ maximizes the entropy?

Can be solved with constrained optimization:

$$\max_{\mathbf{p}} [H(\mathbf{p})] \quad \text{s.t.} \quad \sum_{i=1}^n p_i = 1$$

Form the Lagrangian:

$$J(\mathbf{p}, \lambda) = -\sum_{i=1}^n p_i \cdot \lg(p_i) + \lambda \left(\sum_{i=1}^n p_i - 1 \right)$$

$$\frac{\partial J}{\partial p_i} = ?$$

Calculus
exercise

$$\ln(x)' = \frac{1}{x}$$

$$\lg(x)' = \left(\frac{\ln(x)}{\ln(2)} \right)' = \frac{1}{x \cdot \ln(2)}$$

$$(x \cdot \lg(x))' = \frac{1}{\ln(2)} + \lg(x)$$

Maximum Entropy distribution for nonbinary discrete RV

Entropy of a discrete distribution with n outcomes $H(\mathbf{p}) = -\sum_{i=1}^n p_i \cdot \lg(p_i)$

Which choice of $\mathbf{p} = (p_1, \dots, p_n)$ maximizes the entropy?

Can be solved with constrained optimization: $\max_{\mathbf{p}} [H(\mathbf{p})]$ s.t. $\sum_{i=1}^n p_i = 1$

Form the Lagrangian:

$$J(\mathbf{p}, \lambda) = -\sum_{i=1}^n p_i \cdot \lg(p_i) + \lambda \left(\sum_{i=1}^n p_i - 1 \right)$$

$\frac{\partial J}{\partial \lambda} = \sum p_i - 1 = 0$

$$\frac{\partial J}{\partial p_i} = -\frac{1}{\ln(2)} - \lg(p_i) + \lambda = 0$$

$$\Leftrightarrow \lg(p_i) = \lambda - \frac{1}{\ln(2)} \Leftrightarrow p_i = 2^{\lambda - \frac{1}{\ln(2)}}$$

What next?



Calculus
exercise

$$\ln(x)' = \frac{1}{x}$$

$$\lg(x)' = \left(\frac{\ln(x)}{\ln(2)} \right)' = \frac{1}{x \cdot \ln(2)}$$

$$(x \cdot \lg(x))' = \frac{1}{\ln(2)} + \lg(x)$$

Maximum Entropy distribution for nonbinary discrete RV

Entropy of a discrete distribution with n outcomes $H(\mathbf{p}) = -\sum_{i=1}^n p_i \cdot \lg(p_i)$

Which choice of $\mathbf{p} = (p_1, \dots, p_n)$ maximizes the entropy?

Can be solved with constrained optimization:

$$\max_{\mathbf{p}} [H(\mathbf{p})] \quad \text{s.t.} \quad \sum_{i=1}^n p_i = 1$$

Form the Lagrangian:

$$J(\mathbf{p}, \lambda) = -\sum_{i=1}^n p_i \cdot \lg(p_i) + \lambda \left(\sum_{i=1}^n p_i - 1 \right)$$

$$\frac{\partial J}{\partial p_i} = -\frac{1}{\ln(2)} - \lg(p_i) + \lambda = 0$$

$$\Leftrightarrow \lg(p_i) = \lambda - \frac{1}{\ln(2)} \Leftrightarrow p_i = 2^{\lambda - \frac{1}{\ln(2)}} =: C$$

we are done 😊, all p_i are identical!

$$\sum_{i=1}^n p_i = 1 \Leftrightarrow \sum_{i=1}^n C = 1 \Leftrightarrow C = \frac{1}{n}$$

Calculus
exercise

$$\ln(x)' = \frac{1}{x}$$

$$\lg(x)' = \left(\frac{\ln(x)}{\ln(2)} \right)' = \frac{1}{x \cdot \ln(2)}$$

$$(x \cdot \lg(x))' = \frac{1}{\ln(2)} + \lg(x)$$

Part 1: Theory

L06: Basics of entropy (3/6)

[joint entropy, conditional entropy, mutual information, cross entropy]

Wolfgang Gatterbauer, Javed Aslam

cs7840 Foundations and Applications of Information Theory (fa24)

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

9/23/2024

Pre-class conversations

- Last class recapitulation
- Today:
 - Intuition behind entropy with examples continued
 - Together with the general principles of entropy
 - Then we are changing back to compression

Properties of information (entropy) by example (continued)

Learning partial information



EXAMPLE 4: We again roll the unbiased 8-sided die with sides $\mathcal{X}=\{1,2, \dots,8\}$.

We get two messages: U_1 that the outcome of a roll is even, U_2 that the outcome of the same roll is ≤ 4 . **How much information did we learn after each message?**

$$H(U_1) = ?$$

$$H(U_2) = ?$$

$$H(U_2|U_1) = ?$$



Learning partial information



EXAMPLE 4: We again roll the unbiased 8-sided die with sides $\mathcal{X}=\{1,2, \dots,8\}$.

We get two messages: U_1 that the outcome of a roll is even, U_2 that the outcome of the same roll is ≤ 4 . How much information did we learn after each message?

$$H(U_1) = \lg\left(\frac{8}{4}\right) = 1 \text{ bit}$$

$$H(U_2) = \lg\left(\frac{8}{4}\right) = 1 \text{ bit}$$

$$H(U_2|U_1) = \lg\left(\frac{4}{2}\right) = 1 \text{ bit}$$

$$H(U_2|U_1) = H(U_2) = 1$$

?

Learning partial information



EXAMPLE 4: We again roll the unbiased 8-sided die with sides $\mathcal{X}=\{1,2, \dots,8\}$.

We get two messages: U_1 that the outcome of a roll is even, U_2 that the outcome of the same roll is ≤ 4 . How much information did we learn after each message?

$$H(U_1) = \lg\left(\frac{8}{4}\right) = 1 \text{ bit}$$

$$H(U_2) = \lg\left(\frac{8}{4}\right) = 1 \text{ bit}$$

$$H(U_2|U_1) = \lg\left(\frac{4}{2}\right) = 1 \text{ bit}$$

messages are independent

$$H(U_2|U_1) = H(U_2) = 1$$

1	2
3	4
5	6
7	8

How do the messages reduce the possible outcomes?



Learning partial information



EXAMPLE 4: We again roll the unbiased 8-sided die with sides $\mathcal{X}=\{1,2, \dots,8\}$.

We get two messages: U_1 that the outcome of a roll is even, U_2 that the outcome of the same roll is ≤ 4 . How much information did we learn after each message?

$$H(U_1) = \lg\left(\frac{8}{4}\right) = 1 \text{ bit}$$

$$H(U_2) = \lg\left(\frac{8}{4}\right) = 1 \text{ bit}$$

$$H(U_2|U_1) = \lg\left(\frac{4}{2}\right) = 1 \text{ bit}$$

messages are independent

$$H(U_2|U_1) = H(U_2) = 1$$

U_1
→

1	2
3	4
5	6
7	8

Learning partial information



EXAMPLE 4: We again roll the unbiased 8-sided die with sides $\mathcal{X}=\{1,2, \dots,8\}$.

We get two messages: U_1 that the outcome of a roll is even, U_2 that the outcome of the same roll is ≤ 4 . How much information did we learn after each message?

$$H(U_1) = \lg\left(\frac{8}{4}\right) = 1 \text{ bit}$$

$$H(U_2) = \lg\left(\frac{8}{4}\right) = 1 \text{ bit}$$

$$H(U_2|U_1) = \lg\left(\frac{4}{2}\right) = 1 \text{ bit}$$

messages are independent

$$H(U_2|U_1) = H(U_2) = 1$$

the events are independent

$$p(U_2|U_1) = \underbrace{p(U_2)} = \frac{1}{2}$$

probability of the event $X \leq 4$

	U_1
	→
1	2
3	4
5	6
7	8
	U_2
	↑

Learning partial information



EXAMPLE 4: We again roll the unbiased 8-sided die with sides $\mathcal{X}=\{1,2, \dots,8\}$.

We get two messages: U_1 that the outcome of a roll is even, U_2 that the outcome of the same roll is ≤ 4 . How much information did we learn after each message?

$$H(U_1) = \lg\left(\frac{8}{4}\right) = 1 \text{ bit}$$

$$H(U_2) = \lg\left(\frac{8}{4}\right) = 1 \text{ bit}$$

$$H(U_2|U_1) = \lg\left(\frac{4}{2}\right) = 1 \text{ bit}$$

messages are independent

$$H(U_2|U_1) = H(U_2) = 1$$

the events are independent

$$p(U_2|U_1) = \underbrace{p(U_2)} = \frac{1}{2}$$

probability of the event $X \leq 4$

	U_1 →		
	1	2	
	3	4	
	5	6	
	7	8	
			U_2 ↑

$$H(\{U_1, U_2\}) = H(U_1) + H(U_2|U_1)$$

$$= H(U_1) + H(U_2)$$

U_1 and U_2 are independent

{ 1, 2, 3, 4, 5, 6, 7, 8 }
000 **001** 010 **011** 100 101 110 111

We learned 2 bits independently

0?1

Properties of information (entropy) abstracted

Entropy

Given a discrete RV X with probability mass function (PMF) $p(x) = \mathbb{P}[X = x]$, for $x \in \mathcal{X}$. Entropy is defined as:

$$H(X) = \mathbb{E} \left[\lg \left(\frac{1}{p(X)} \right) \right] = \sum_x p(x) \cdot \lg \left(\frac{1}{p(x)} \right)$$

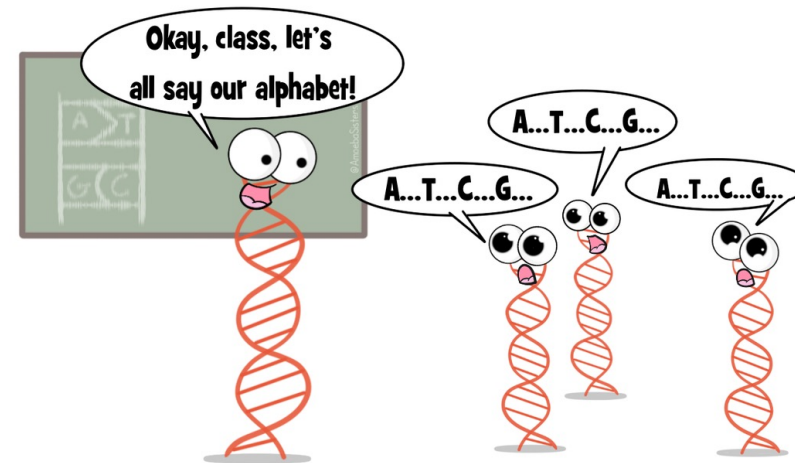
Alternative notation: $p(X) = p_X(x)$. Also: $\mathbb{E}_p[\dots]$ or $\mathbb{E}_X[\dots]$ or $\mathbb{E}_{X \sim p}[\dots]$ for the expected value operator w.r. to the distribution p

Entropy is **label-invariant**, meaning that it depends only on the probability distribution and not on the actual values that the random variable X can take.

$$\mathcal{X} = \{1, 2, 3, 4\}$$



$$\mathcal{X} = \{\mathbf{A}, \mathbf{T}, \mathbf{G}, \mathbf{C}\}$$



Joint Entropy

treat (X, Y) just like a single vector-valued RV $Z = \langle X, Y \rangle$

Given two RVs X and Y with PMF $p(X, Y)$, their **joint entropy** is:

$$H(X, Y) = \mathbb{E} \left[\lg \left(\frac{1}{p(X, Y)} \right) \right] = \sum_x \sum_y p(x, y) \cdot \lg \left(\frac{1}{p(x, y)} \right)$$

Other notation: $p(X, Y) = p_{X, Y}(x, y)$.
Also: $\mathbb{E}_{X, Y \sim p}[\dots]$ or $\mathbb{E}_{X, Y \sim p}[\dots]$ or $\mathbb{E}_p[\dots]$

If X and Y are independent:

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

How can we prove that? 

Joint Entropy

Given two RVs X and Y with PMF $p(X, Y)$, their **joint entropy** is:

$$H(X, Y) = \mathbb{E} \left[\lg \left(\frac{1}{p(X, Y)} \right) \right] = \sum_x \sum_y p(x, y) \cdot \lg \left(\frac{1}{p(x, y)} \right)$$

Other notation: $p(X, Y) = p_{X,Y}(x, y)$.
Also: $\mathbb{E}_{X,Y \sim p}[\dots]$ or $\mathbb{E}_{X,Y \sim p}[\dots]$ or $\mathbb{E}_p[\dots]$

If X and Y are independent:

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

$$\begin{aligned} H(X, Y) &= \mathbb{E} \left[\lg \left(\frac{1}{p(X, Y)} \right) \right] \\ &= \mathbb{E} \left[\lg \left(\frac{1}{p(X) \cdot p(Y)} \right) \right] \\ &= \mathbb{E} \left[\lg \left(\frac{1}{p(X)} \right) + \lg \left(\frac{1}{p(Y)} \right) \right] \\ &= \mathbb{E} \left[\lg \left(\frac{1}{p(X)} \right) \right] + \mathbb{E} \left[\lg \left(\frac{1}{p(Y)} \right) \right] \\ &= H(X) + H(Y) \end{aligned}$$


Conditional Entropy, Chain rule of Entropy

Given two RVs X and Y with PMF $p(X, Y)$, their joint entropy is:

$$H(X, Y) = \mathbb{E} \left[\lg \left(\frac{1}{p(X, Y)} \right) \right] = \sum_x \sum_y p(x, y) \cdot \lg \left(\frac{1}{p(x, y)} \right)$$

If X and Y are **not independent**:

$$H(X, Y) = \cancel{H(X) + H(Y)}$$

What do we need to do? 

If X and Y are **not independent**, observing X might contain already some information about Y , so simply adding the information from each would **overcount**.

Conditional Entropy, Chain rule of Entropy

Given two RVs X and Y with PMF $p(X, Y)$, their joint entropy is:

$$H(X, Y) = \mathbb{E} \left[\lg \left(\frac{1}{p(X, Y)} \right) \right] = \sum_x \sum_y p(x, y) \cdot \lg \left(\frac{1}{p(x, y)} \right)$$

If X and Y are not independent:

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

Conditional entropy $H(Y|X)$: the expected amount of information needed to describe the outcome of RV Y given that the value of another RV X is known

$$H(Y|X) = \sum_x p(x) \cdot H(Y|X = x)$$

$$\begin{aligned} H(X, Y) &= \sum_x \sum_y p(x, y) \cdot \lg \left(\frac{1}{p(x, y)} \right) \\ &= \sum_x \sum_y p(x) \cdot p(y|x) \cdot \lg \left(\frac{1}{p(x) \cdot p(y|x)} \right) \\ &= \sum_x \sum_y p(x) \cdot p(y|x) \cdot \lg \left(\frac{1}{p(x)} \right) + \sum_x \sum_y p(x) \cdot p(y|x) \cdot \lg \left(\frac{1}{p(y|x)} \right) \\ &= \underbrace{\sum_x p(x) \cdot \lg \left(\frac{1}{p(x)} \right)}_{H(X)} \cdot \underbrace{\sum_y p(y|x)}_1 + \sum_x p(x) \cdot \underbrace{\sum_y p(y|x) \cdot \lg \left(\frac{1}{p(y|x)} \right)}_{H(Y|X = x)} \end{aligned}$$

DEFINITION of **conditional entropy** $H(Y|X) = \sum_{x,y} p(x, y) \cdot \lg \left(\frac{1}{p(y|x)} \right)$

Chain rule for entropy

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

... obvious generalization to
a **chain of** (not necessarily
independent) **observations**

If not independent:

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

conditional entropy

Learning partial information

EXAMPLE 5: We again roll the unbiased 8-sided die with sides $\mathcal{X}=\{1,2, \dots,8\}$.



We then get a message U : "The outcome of the roll is even, and by the way, the next president of the US will be ...". Assuming two equally likely outcomes for the election, **how much information did we learn?**

?

Learning partial information

EXAMPLE 5: We again roll the unbiased 8-sided die with sides $\mathcal{X}=\{1,2, \dots,8\}$.



We then get a message U : "The outcome of the roll is even, and by the way, the next president of the US will be ...". Assuming two equally likely outcomes for the election, how much information did we learn?

- We still learn $3-2=1$ bit about the roll of the die X .
 - We also learn 1 bit about the election outcome.
- } We learned 2 bits (U contains 2 bits)

Wasn't information supposed to be additive?



Learning partial information

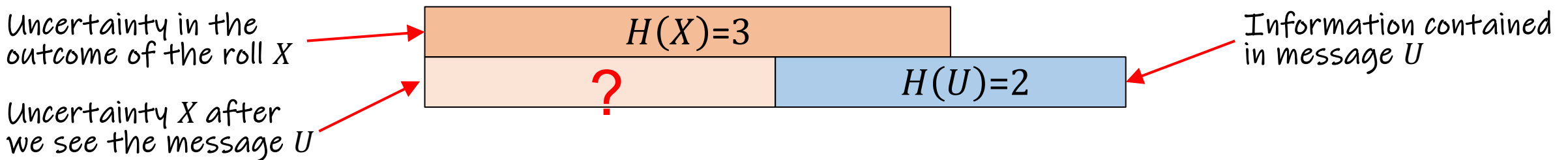


EXAMPLE 5: We again roll the unbiased 8-sided die with sides $\mathcal{X}=\{1,2, \dots,8\}$.

We then get a message U : "The outcome of the roll is even, and by the way, the next president of the US will be ...". Assuming two equally likely outcomes for the election, how much information did we learn?

- We still learn $3-2=1$ bit about the roll of the die X .
 - We also learn 1 bit about the election outcome.
- } We learned 2 bits (U contains 2 bits)

Wasn't information supposed to be additive?



Learning partial information

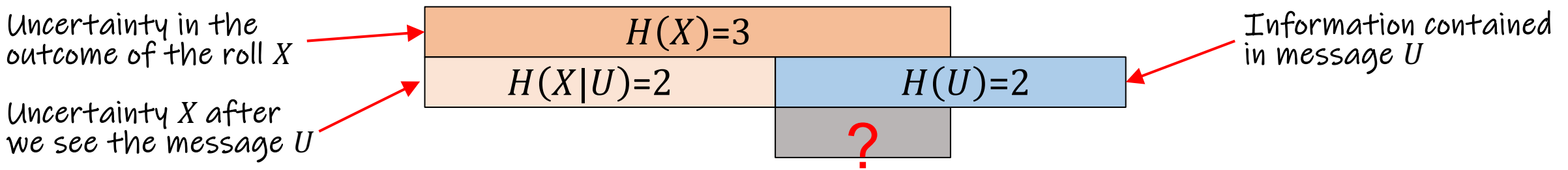


EXAMPLE 5: We again roll the unbiased 8-sided die with sides $\mathcal{X}=\{1,2, \dots,8\}$.

We then get a message U : "The outcome of the roll is even, and by the way, the next president of the US will be ...". Assuming two equally likely outcomes for the election, how much information did we learn?

- We still learn $3-2=1$ bit about the roll of the die X .
 - We also learn 1 bit about the election outcome.
- } We learned 2 bits (U contains 2 bits)

Wasn't information supposed to be additive?



Learning partial information

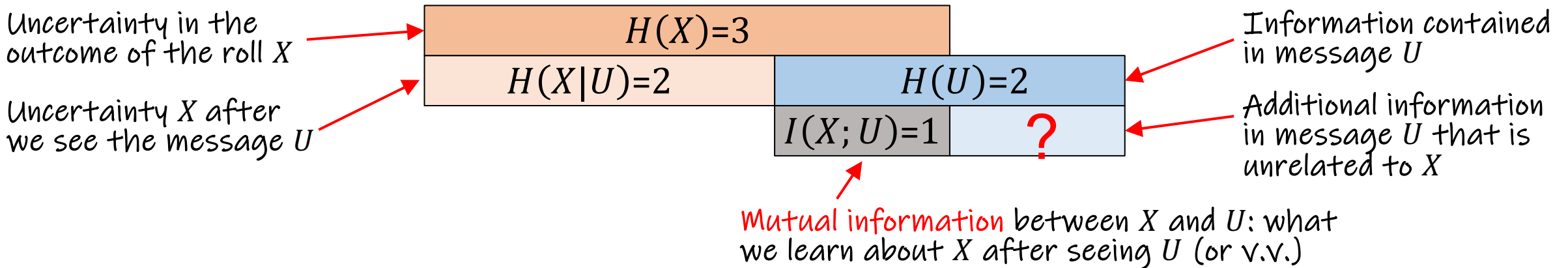


EXAMPLE 5: We again roll the unbiased 8-sided die with sides $\mathcal{X}=\{1,2, \dots,8\}$.

We then get a message U : "The outcome of the roll is even, and by the way, the next president of the US will be ...". Assuming two equally likely outcomes for the election, how much information did we learn?

- We still learn $3-2=1$ bit about the roll of the die X .
 - We also learn 1 bit about the election outcome.
- } We learned 2 bits (U contains 2 bits)

Wasn't information supposed to be additive?



Learning partial information

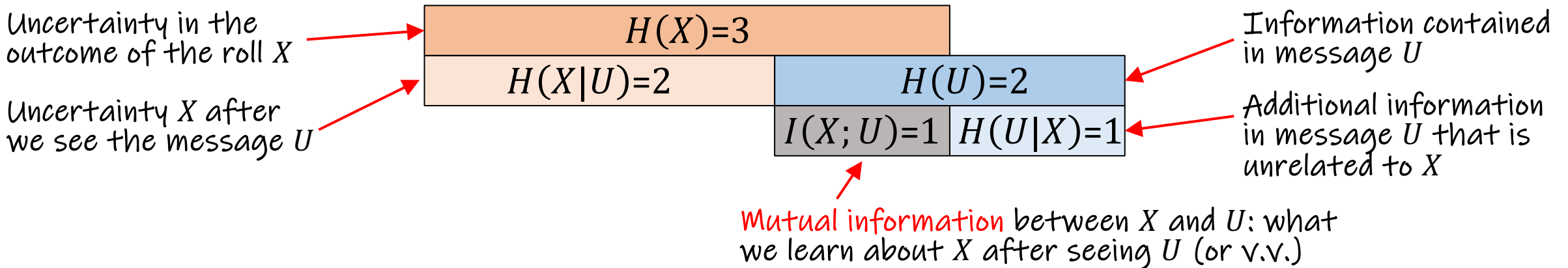


EXAMPLE 5: We again roll the unbiased 8-sided die with sides $\mathcal{X}=\{1,2, \dots,8\}$.

We then get a message U : "The outcome of the roll is even, and by the way, the next president of the US will be ...". Assuming two equally likely outcomes for the election, how much information did we learn?

- We still learn $3-2=1$ bit about the roll of the die X .
 - We also learn 1 bit about the election outcome.
- } We learned 2 bits (U contains 2 bits)

Wasn't information supposed to be additive?



Mutual information

Given two RVs X and Y , **mutual information** is the amount of information that Y provides about X (thus when Y is observed, but X is not).

$$I(X; Y) := H(X) - H(X|Y)$$

Is this function symmetric in X and Y ?

Mutual information

Given two RVs X and Y , **mutual information** is the amount of information that Y provides about X (thus when Y is observed, but X is not).

$$I(X; Y) := H(X) - H(X|Y)$$

$$= H(X) - (H(X, Y) - H(Y))$$

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

$$= H(Y) - H(Y|X)$$

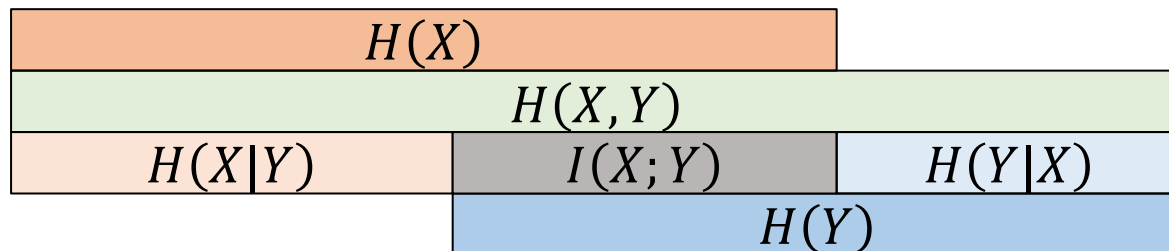
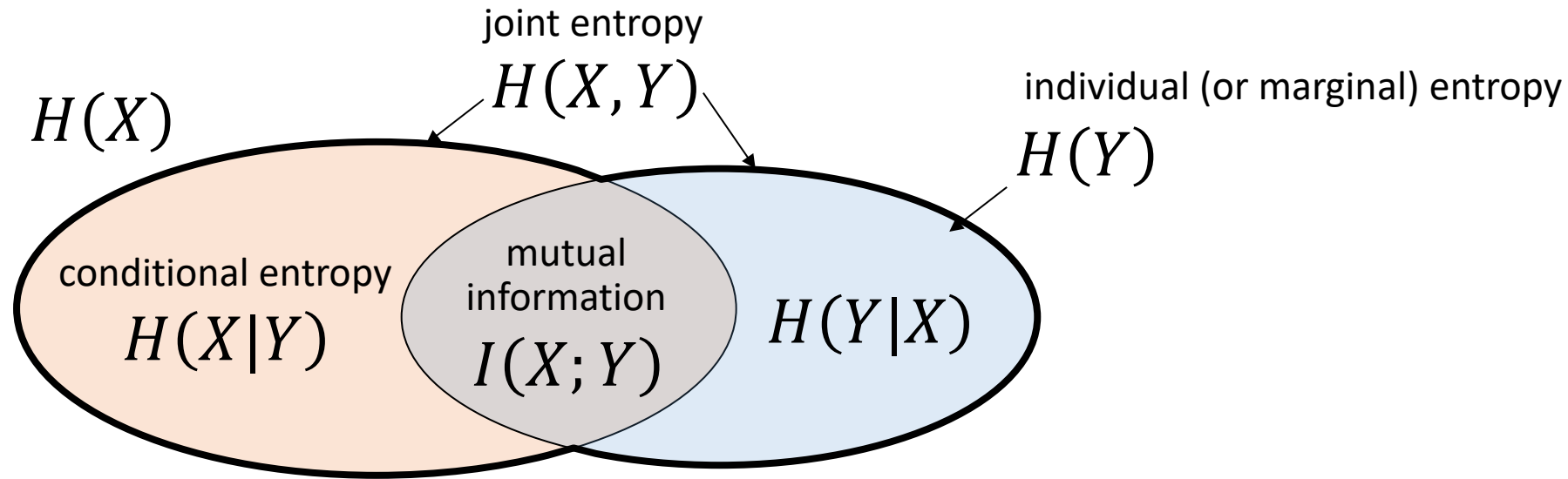
Conditional entropy: the amount of information needed to describe the outcome of RV Y given that we know the value of another RV X .

symmetric in X and Y !

Thus, $I(X; Y) = I(Y; X)$

That's why it is called "**mutual**" information (it does not "prefer" X or Y).
Reduction of the uncertainty of one RV once we observe the other.

Entropy, conditional entropy, mutual information



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

$$H(X, Y) = H(X) + H(Y) - I(X; Y)$$

The bar diagrams are inspired by Fig 8.1 in "MacKay. Information Theory, Inference, and learning Algorithms. Cambridge University Press, 2002." <https://www.inference.org.uk/itprnn/book.pdf>.

In particular, see the Interesting discussion and explanation in the solution to exercise 8.8 for why VENN diagrams (with more than 2 variables) can be misleading.

Gatterbauer, Aslam. Foundations and Applications of Information Theory: <https://northeastern-datalab.github.io/cs7840/>

Entropy, conditional entropy, mutual information

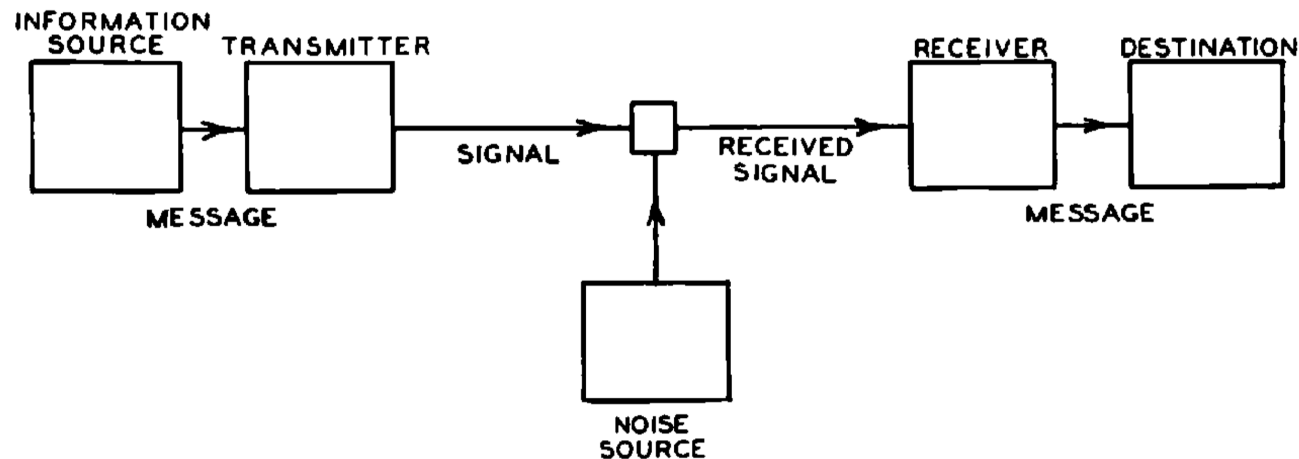
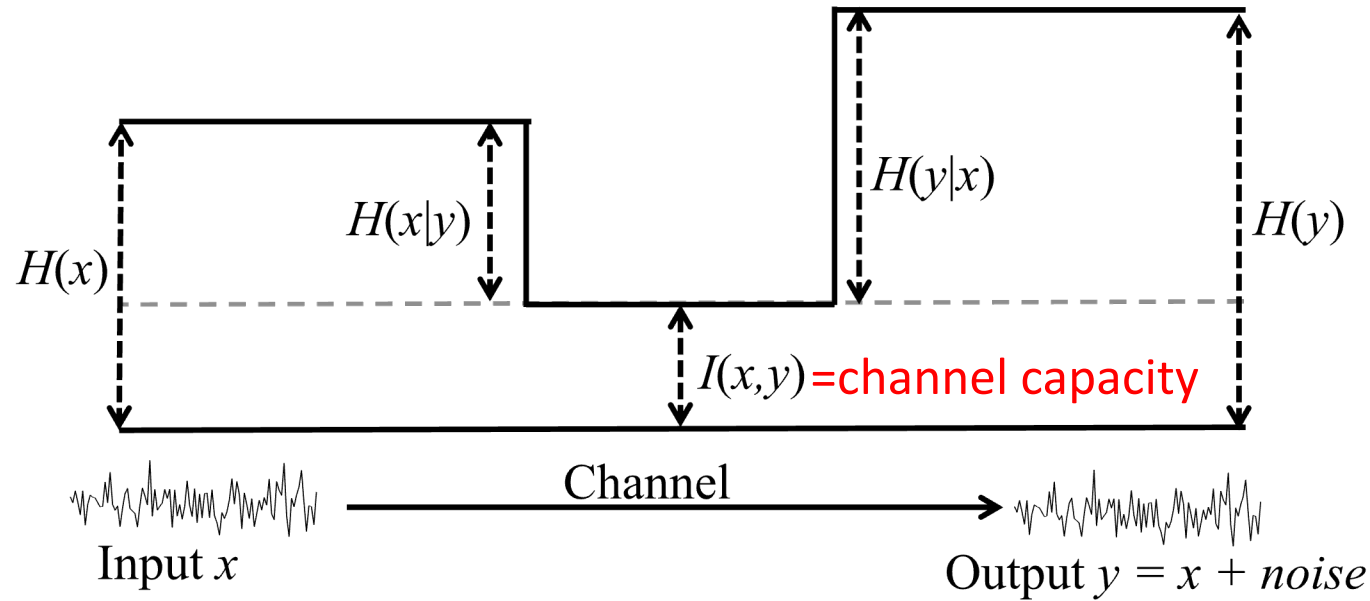


Fig. 1—Schematic diagram of a general communication system.

Self-information

What is $I(X; X)$?

How much does X tell us about itself?

Self-information

What is $I(X; X)$?

How much does X tell us about itself?

$$I(X; X) = H(X) - H(X|X)$$

 $= 0$ no uncertainty (entropy) left)

$$I(X; X) = H(X)$$

We learn from X everything about X
Entropy is "self-information".

Relative entropy
= KL divergence
(\neq Cross-Entropy)

Relative Entropy = KL divergence (\neq Cross-Entropy)

The **relative entropy** (or KL divergence) of a distribution p with respect to a distribution q defined on the alphabet \mathcal{X} of RV X is:

$$D_{\text{KL}}(p||q) = \mathbb{E}_p \left[\lg \left(\frac{p(X)}{q(X)} \right) \right] = \sum_{x \in \mathcal{X}} p(x) \cdot \lg \left(\frac{p(x)}{q(x)} \right)$$

$\mathbb{E}_p[\dots]$ also written as $\mathbb{E}_{X \sim p}[\dots]$ for the expected value operator w.r. to the distribution p

It measures the inefficiency for assuming a distribution q instead of a true distribution p for RV.

If we use q to construct a binary code, the expected message length is called **cross-entropy**:

$$H(p||q) = D_{\text{KL}}(p||q) + H(p)$$

my surprise for seeing x , given my assumption of $q(x)$

$$= \mathbb{E}_p \left[\lg \left(\frac{1}{q(X)} \right) \right] = \sum_{x \in \mathcal{X}} p(x) \cdot \lg \left(\frac{1}{q(x)} \right)$$

my expected surprise given p as the true distribution

$H(p)$	$D(p q)$
$H(p q)$	

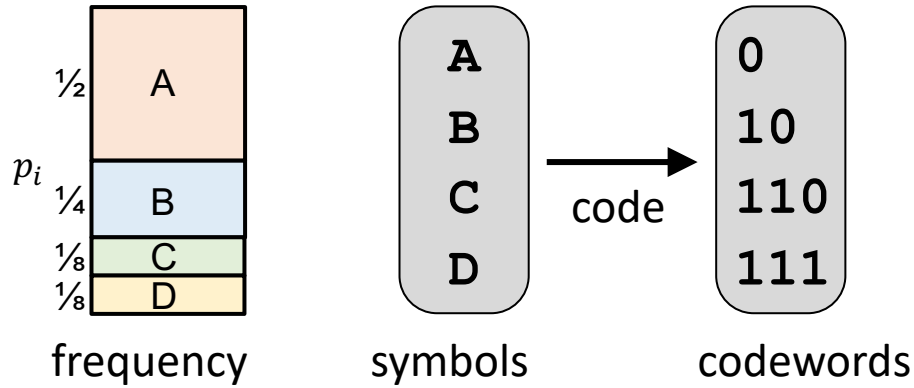
Cross-entropy is usually written as $H(p, q)$, but that notation hides its asymmetry and looks too similar to joint entropy. We prefer the notation $H(p||q)$ which captures the asymmetry with a similar notation as $D_{\text{KL}}(p||q)$. Another non-standard notation is $H_p(q)$ which shows that p is the true distribution, whereas q determines the surprise.

Gatterbauer, Aslam. Foundations and Applications of Information Theory: <https://northeastern-datalab.github.io/cs7840/>

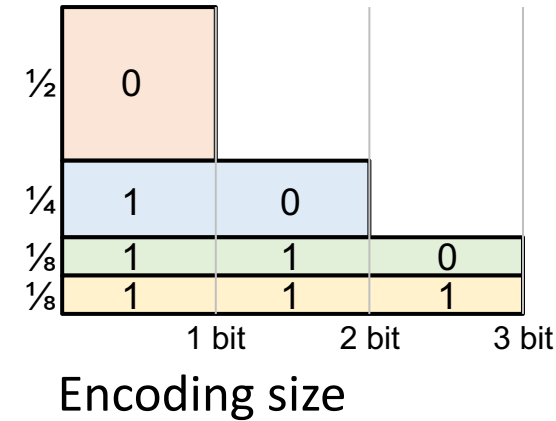


Compressing messages via variable length codes

- Assume we have the following symbol frequency:



New expected length :



$$\lg\left(\frac{1}{2}\right) = -1$$

$$\lg\left(\frac{1}{4}\right) = -2$$

$$\lg\left(\frac{1}{8}\right) = -3$$

$$\frac{1}{2} \cdot 1$$

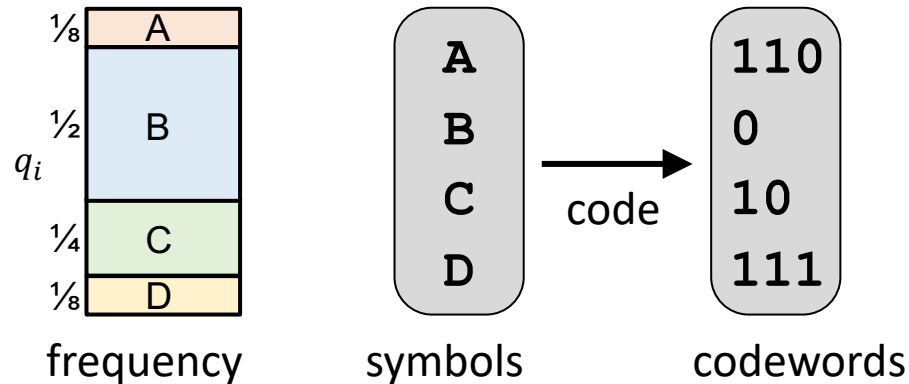
$$\frac{1}{4} \cdot 2 = 1.75 \text{ bits!}$$

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

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

$$\text{Entropy } H(\mathbf{p}) := -\sum_i p_i \cdot \lg(p_i) = 1.75 \text{ bits!}$$

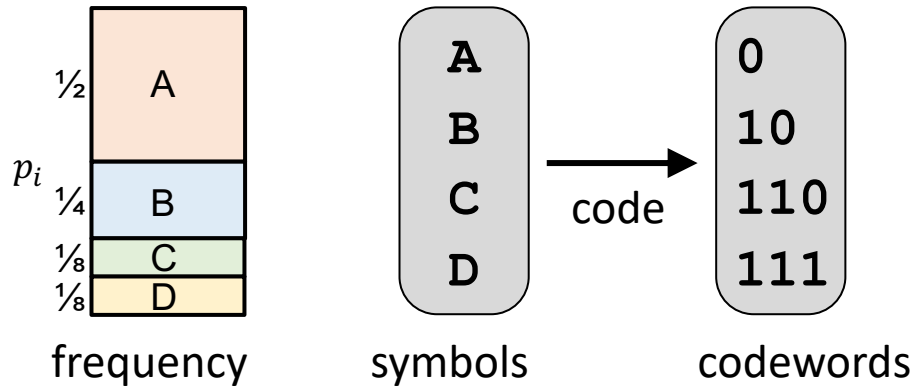
- What if we assume following distribution:



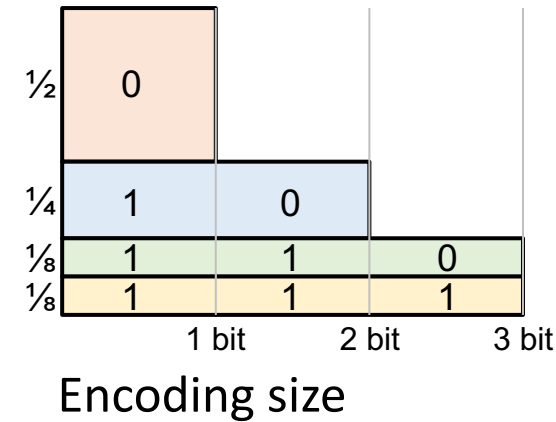
What is our expected message length per symbol if we use that code, but p is the actual distribution ?

Compressing messages via variable length codes

- Assume we have the following symbol frequency:



New expected length :



$$\lg\left(\frac{1}{2}\right) = -1$$

$$\lg\left(\frac{1}{4}\right) = -2$$

$$\lg\left(\frac{1}{8}\right) = -3$$

$$\frac{1}{2} \cdot 1$$

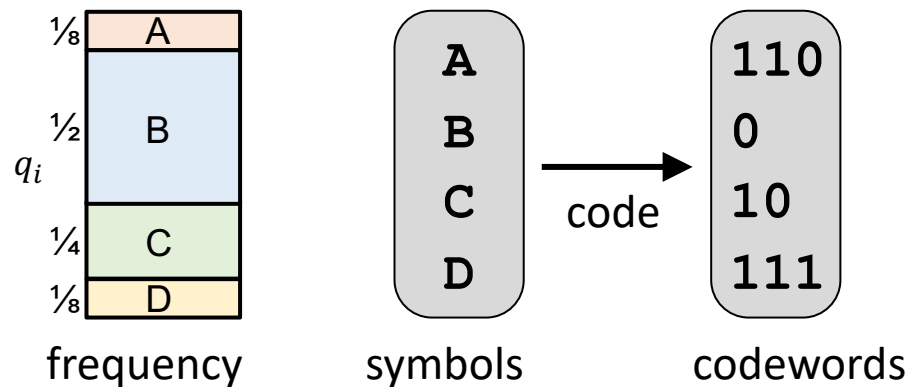
$$\frac{1}{4} \cdot 2 = 1.75 \text{ bits!}$$

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

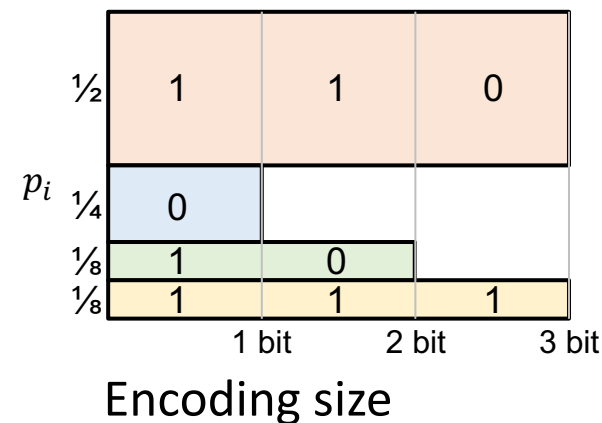
$$\frac{1}{8} \cdot 3$$

$$\text{Entropy } H(\mathbf{p}) := -\sum_i p_i \cdot \lg(p_i) = 1.75 \text{ bits!}$$

- What if we assume following distribution:



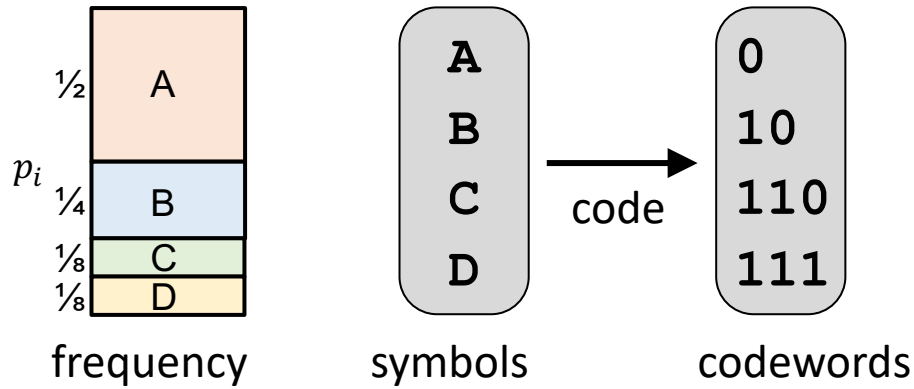
Our new expected message length per symbol:



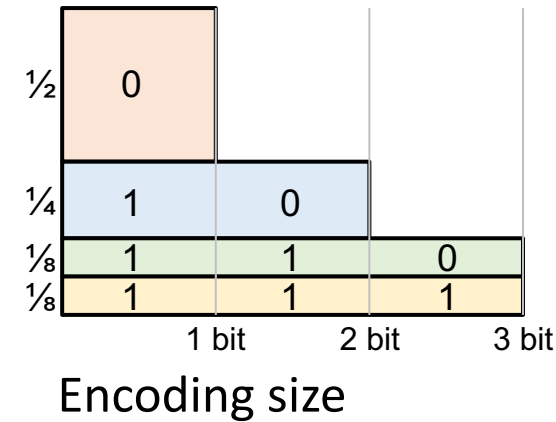


Compressing messages via variable length codes

- Assume we have the following symbol frequency:



New expected length :



$$\lg\left(\frac{1}{2}\right) = -1$$

$$\lg\left(\frac{1}{4}\right) = -2$$

$$\lg\left(\frac{1}{8}\right) = -3$$

$$\frac{1}{2} \cdot 1$$

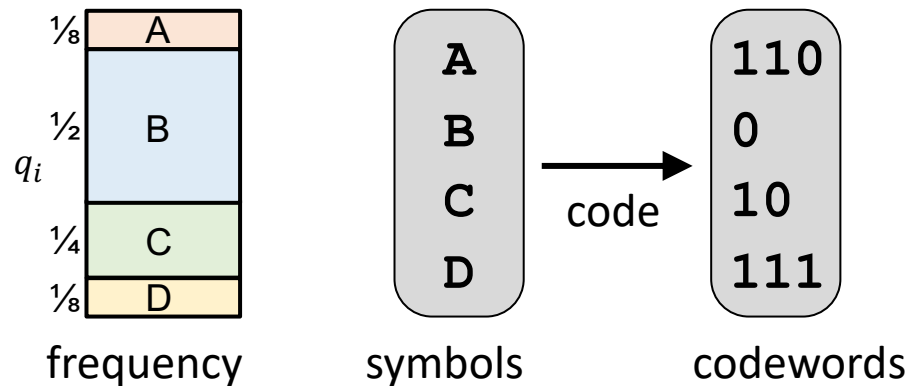
$$\frac{1}{4} \cdot 2 = 1.75 \text{ bits!}$$

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

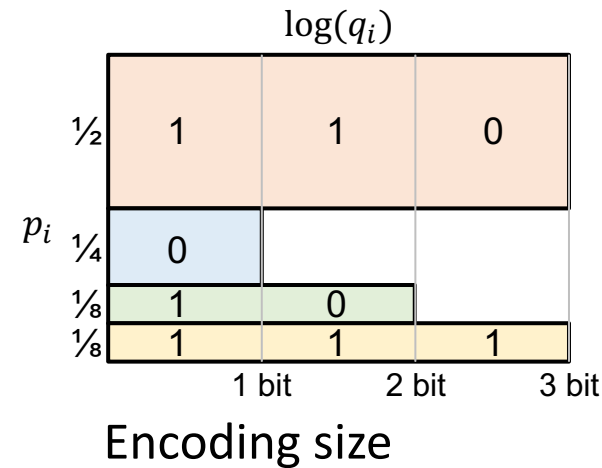
$$\frac{1}{8} \cdot 3$$

$$\text{Entropy } H(\mathbf{p}) := -\sum_i p_i \cdot \lg(p_i) = 1.75 \text{ bits!}$$

- What if we assume following distribution:



Our new expected message length per symbol:

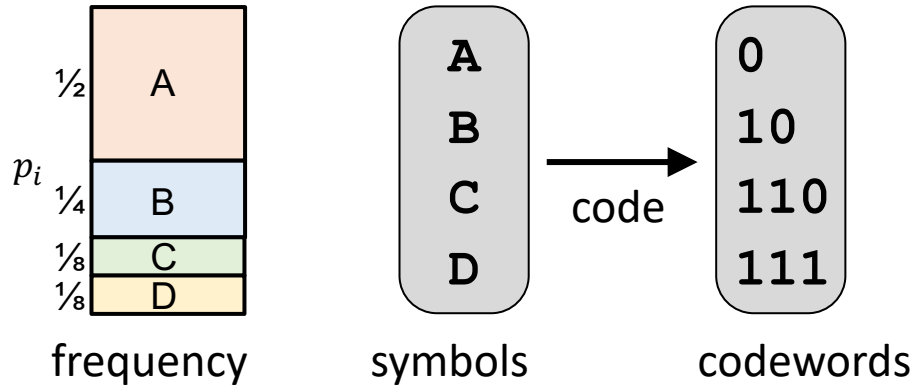


What is the formula we need to evaluate ?

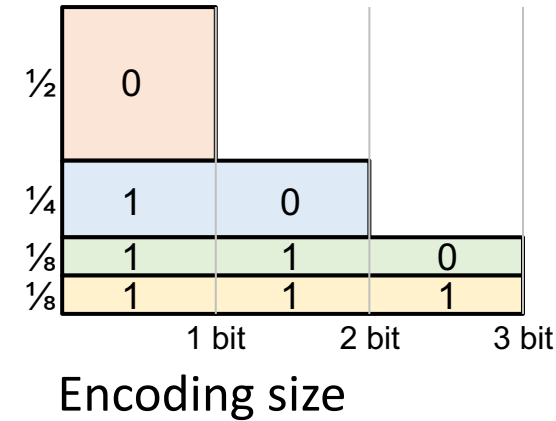
Compressing messages via variable length codes



- Assume we have the following symbol frequency:



New expected length :



$$\lg\left(\frac{1}{2}\right) = -1$$

$$\lg\left(\frac{1}{4}\right) = -2$$

$$\lg\left(\frac{1}{8}\right) = -3$$

$$\frac{1}{2} \cdot 1$$

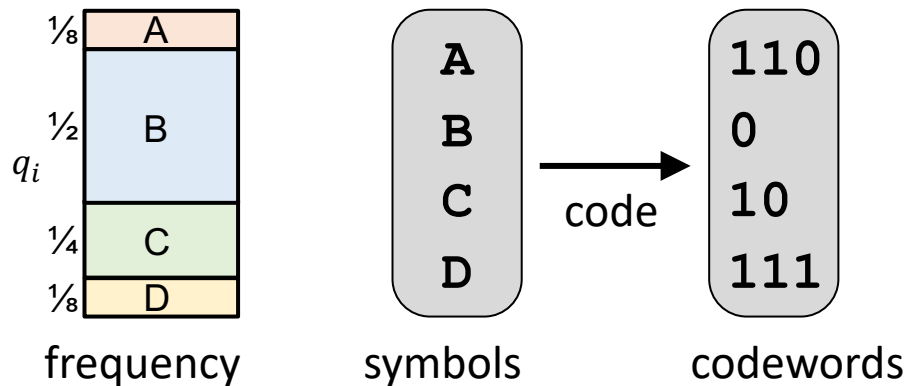
$$\frac{1}{4} \cdot 2 = 1.75 \text{ bits!}$$

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

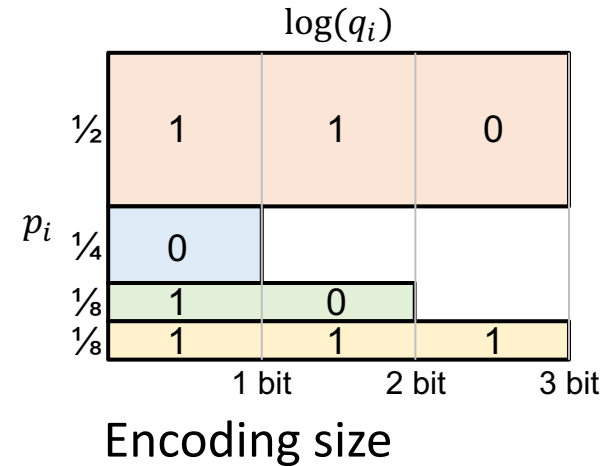
$$\frac{1}{8} \cdot 3$$

$$\text{Entropy } H(\mathbf{p}) := -\sum_i p_i \cdot \lg(p_i) = 1.75 \text{ bits!}$$

- What if we assume following distribution:



Our new expected message length per symbol:



$$= 2.375 \text{ bits!}$$

$$-\sum_i p_i \cdot \lg(q_i)$$

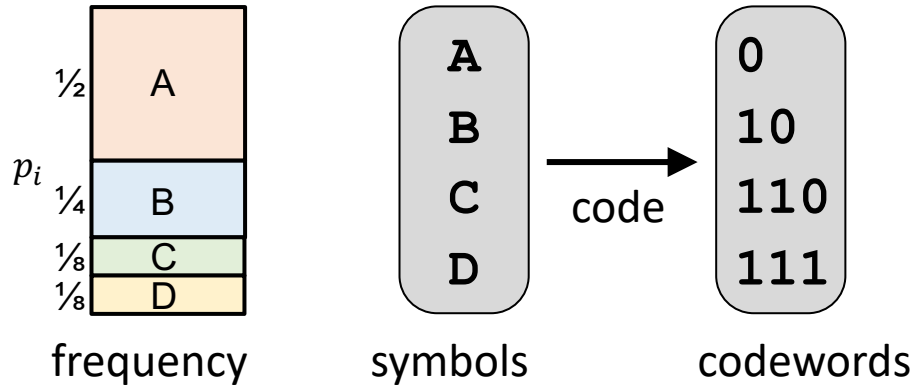
What is this formula called



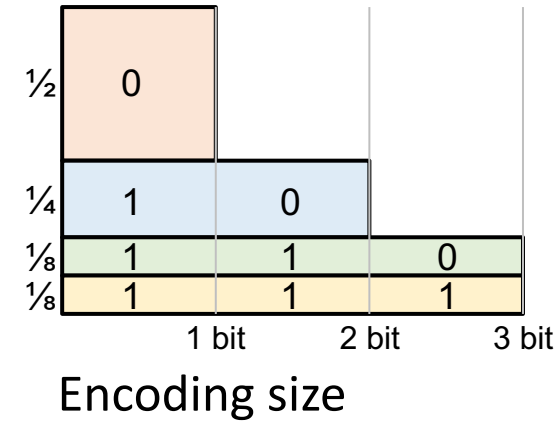
Compressing messages via variable length codes



- Assume we have the following symbol frequency:



New expected length :



$$\lg\left(\frac{1}{2}\right) = -1$$

$$\lg\left(\frac{1}{4}\right) = -2$$

$$\lg\left(\frac{1}{8}\right) = -3$$

$$\frac{1}{2} \cdot 1$$

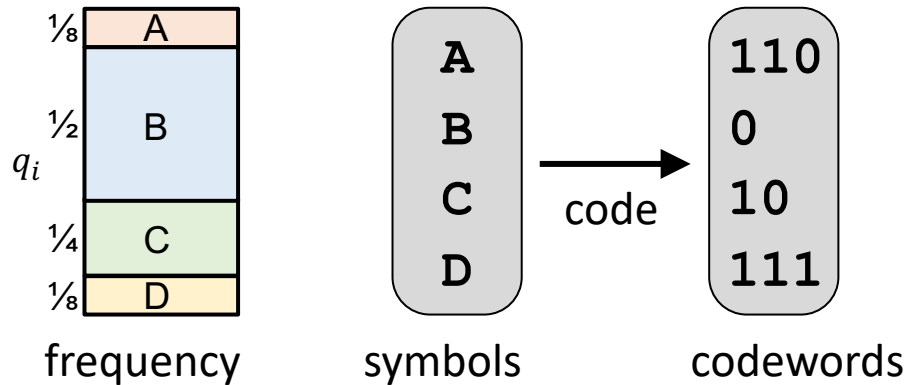
$$\frac{1}{4} \cdot 2 = 1.75 \text{ bits!}$$

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

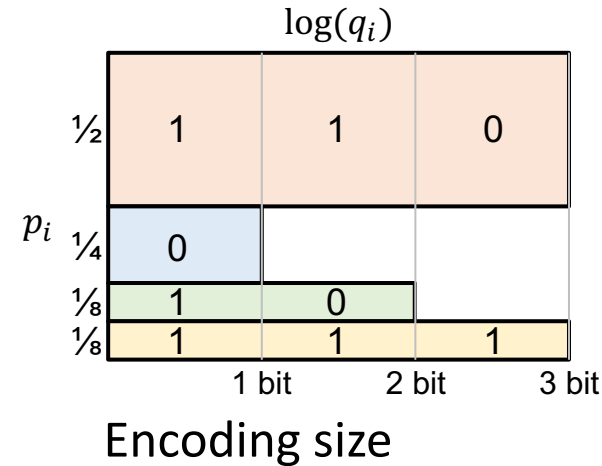
$$\frac{1}{8} \cdot 3$$

Entropy $H(\mathbf{p}) := -\sum_i p_i \cdot \lg(p_i) = 1.75 \text{ bits!}$

- What if we assume following distribution:



Our new expected message length per symbol:



$$= 2.375 \text{ bits!}$$

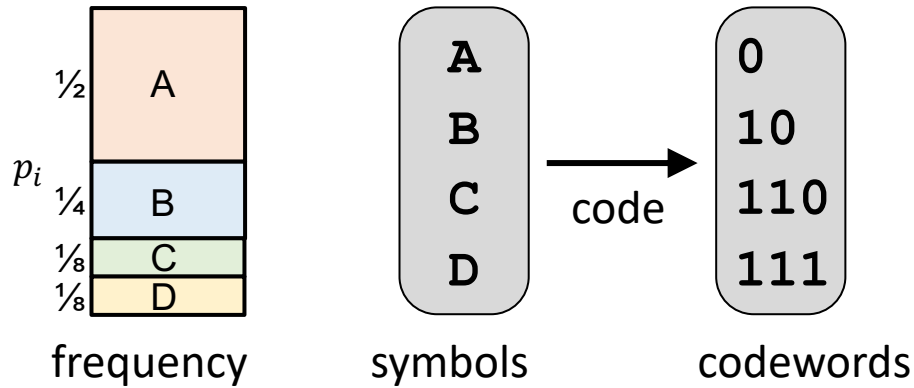
$$-\sum_i p_i \cdot \lg(q_i)$$

Cross entropy $H(p||q)$ ☺

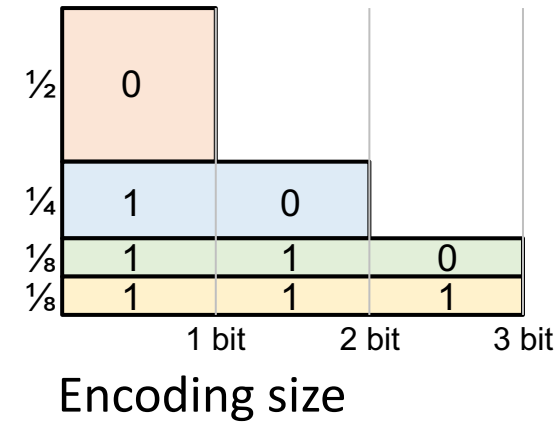
Which distribution q minimizes $H(p||q)$?

Compressing messages via variable length codes

- Assume we have the following symbol frequency:



New expected length :



$$\lg\left(\frac{1}{2}\right) = -1$$

$$\lg\left(\frac{1}{4}\right) = -2$$

$$\lg\left(\frac{1}{8}\right) = -3$$

$$\frac{1}{2} \cdot 1$$

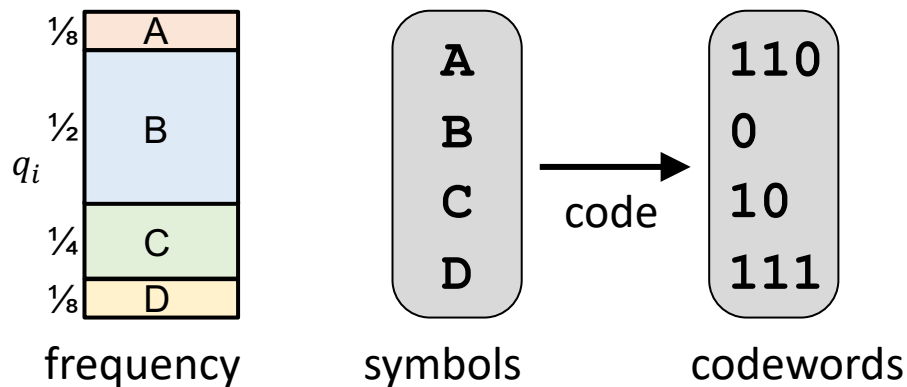
$$\frac{1}{4} \cdot 2 = 1.75 \text{ bits!}$$

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

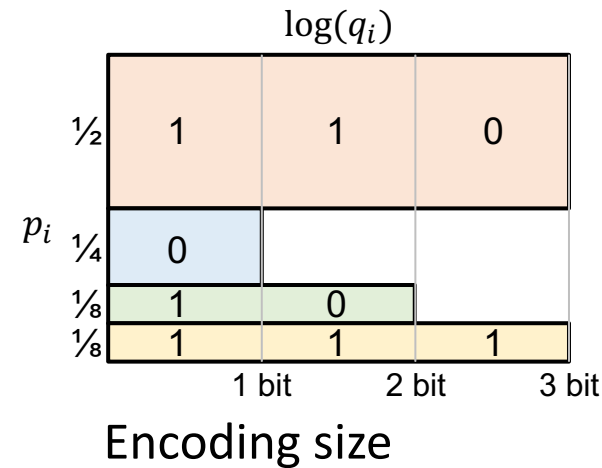
$$\frac{1}{8} \cdot 3$$

$$\text{Entropy } H(\mathbf{p}) := -\sum_i p_i \cdot \lg(p_i) = 1.75 \text{ bits!}$$

- What if we assume following distribution:



Our new expected message length per symbol:



$$= 2.375 \text{ bits!}$$

$$-\sum_i p_i \cdot \lg(q_i)$$

Cross entropy $H(p||q)$ ☺

$$q = p \text{ minimizes } H(p||q)$$

Properties of Relative Entropy = KL divergence

1. Relative entropy is asymmetric (does not satisfy triangle inequality, thus not a metric):

$$D_{\text{KL}}(p||q) \neq D_{\text{KL}}(q||p)$$

EXAMPLE : $\mathbf{u} = \begin{pmatrix} 0.5 \\ 0.5 \end{pmatrix}$ $\mathbf{p} = \begin{pmatrix} p \\ \bar{p} \end{pmatrix}$ $\bar{p} = 1 - p$

	$D_{\text{KL}}(\mathbf{p} \mathbf{u})$	$D_{\text{KL}}(\mathbf{u} \mathbf{p})$
$p = 0.5$?	?
$p = 0$?	?
$p = 0.01$?	?

Properties of Relative Entropy = KL divergence

1. Relative entropy is asymmetric (does not satisfy triangle inequality, thus not a metric):

$$D_{\text{KL}}(p||q) \neq D_{\text{KL}}(q||p)$$

EXAMPLE : $\mathbf{u} = \begin{pmatrix} 0.5 \\ 0.5 \end{pmatrix}$ $\mathbf{p} = \begin{pmatrix} p \\ \bar{p} \end{pmatrix}$ $\bar{p} = 1 - p$

	$D_{\text{KL}}(\mathbf{p} \mathbf{u})$	$D_{\text{KL}}(\mathbf{u} \mathbf{p})$
$p = 0.5$	0	0
$p = 0$?	?
$p = 0.01$?	?

Properties of Relative Entropy = KL divergence

1. Relative entropy is asymmetric (does not satisfy triangle inequality, thus not a metric):

$$D_{\text{KL}}(p||q) \neq D_{\text{KL}}(q||p)$$

EXAMPLE : $\mathbf{u} = \begin{pmatrix} 0.5 \\ 0.5 \end{pmatrix}$ $\mathbf{p} = \begin{pmatrix} p \\ \bar{p} \end{pmatrix}$ $\bar{p} = 1 - p$

	$D_{\text{KL}}(\mathbf{p} \mathbf{u})$	$D_{\text{KL}}(\mathbf{u} \mathbf{p})$
$p = 0.5$	0	0
$p = 0$	1	∞
$p = 0.01$?	?

Properties of Relative Entropy = KL divergence

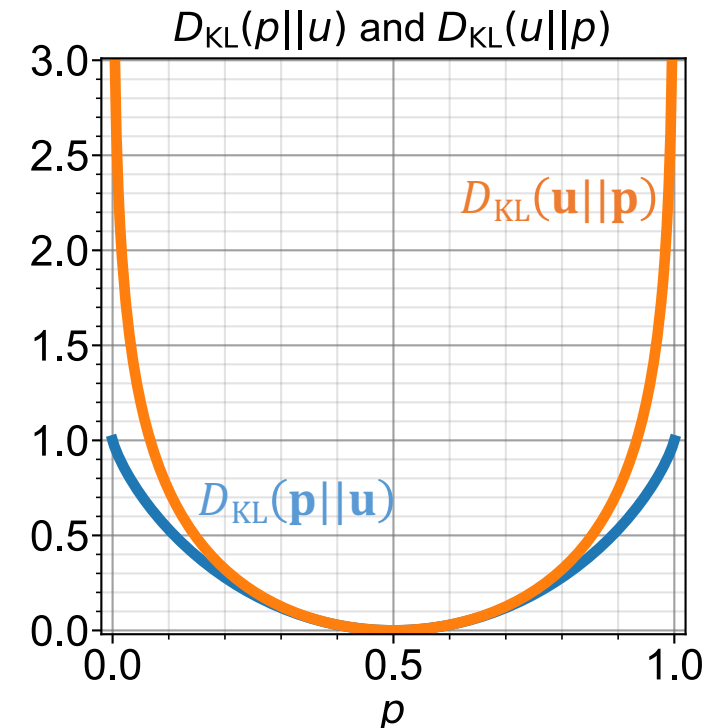
1. Relative entropy is asymmetric (does not satisfy triangle inequality, thus not a metric):

$$D_{\text{KL}}(p||q) \neq D_{\text{KL}}(q||p)$$

EXAMPLE : $\mathbf{u} = \begin{pmatrix} 0.5 \\ 0.5 \end{pmatrix}$ $\mathbf{p} = \begin{pmatrix} p \\ \bar{p} \end{pmatrix}$ $\bar{p} = 1 - p$

	$D_{\text{KL}}(\mathbf{p} \mathbf{u})$	$D_{\text{KL}}(\mathbf{u} \mathbf{p})$
$p = 0.5$	0	0
$p = 0$	1	∞
$p = 0.01$	0.92	2.33
	$\underbrace{.01 \lg\left(\frac{.01}{.5}\right)}_{-0.06} + \underbrace{.99 \lg\left(\frac{.99}{.5}\right)}_{0.96}$	$\underbrace{.5 \lg\left(\frac{.5}{.01}\right)}_{2.82} + \underbrace{.5 \lg\left(\frac{.5}{.99}\right)}_{-0.49}$

What about cross entropies $H(\mathbf{p}||\mathbf{u})$ and $H(\mathbf{u}||\mathbf{p})$?



Properties of Relative Entropy = KL divergence

1. Relative entropy is asymmetric (does not satisfy triangle inequality, thus not a metric):

$$D_{\text{KL}}(p||q) \neq D_{\text{KL}}(q||p)$$

2. $D_{\text{KL}}(p||p) = ?$

Properties of Relative Entropy = KL divergence

1. Relative entropy is asymmetric (does not satisfy triangle inequality, thus not a metric):

$$D_{\text{KL}}(p||q) \neq D_{\text{KL}}(q||p)$$

2. $D_{\text{KL}}(p||p) = 0$

3. $D_{\text{KL}}(p||q) \geq 0$ for all distributions p, q (equality only holds for $p = q$)

We will prove that next (with Jensen's inequality)

Commuting functions: an apparent digression

- Do functions commute with taking the expectation?

$$\mathbb{E}[f(X)] = f(\mathbb{E}[X])$$



Commuting functions: an apparent digression

- Do functions commute with taking the expectation?
- No! This only holds for **linear** functions:
- **Jensen's inequality** for **convex** f :

$$\cancel{\mathbb{E}[f(X)] = f(\mathbb{E}[X])}$$

$$f(x) = ax + b$$

$$\mathbb{E}[ax + b] = a\mathbb{E}[x] + b$$



Commuting functions: an apparent digression

- Do functions commute with taking the expectation?

$$\cancel{\mathbb{E}[f(X)] = f(\mathbb{E}[X])}$$

- No! This only holds for **linear** functions:

$$f(x) = ax + b$$

$$\mathbb{E}[ax + b] = a\mathbb{E}[x] + b$$

- **Jensen's inequality** for **convex** f :

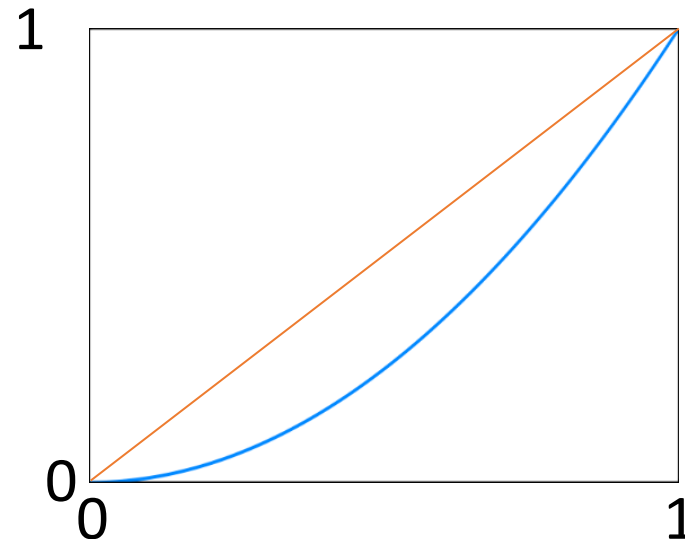
$$\mathbb{E}[f(X)] \geq f(\mathbb{E}[X])$$

- Example $f(x) = x^2$:

Consider the interval $0 \leq x \leq 1$:

$$f(\mathbb{E}[X]) = ?$$

$$\mathbb{E}[f(X)] = ?$$



Commuting functions: an apparent digression

- Do functions commute with taking the expectation?

$$\cancel{\mathbb{E}[f(X)] = f(\mathbb{E}[X])}$$

- No! This only holds for **linear** functions:

$$f(x) = ax + b$$

$$\mathbb{E}[ax + b] = a\mathbb{E}[x] + b$$

- **Jensen's inequality** for **convex** f :

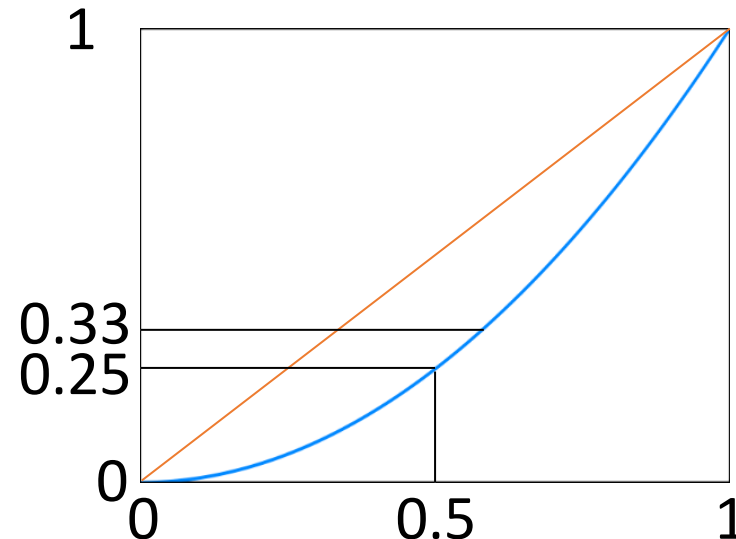
$$\mathbb{E}[f(X)] \geq f(\mathbb{E}[X])$$

- Example $f(x) = x^2$:

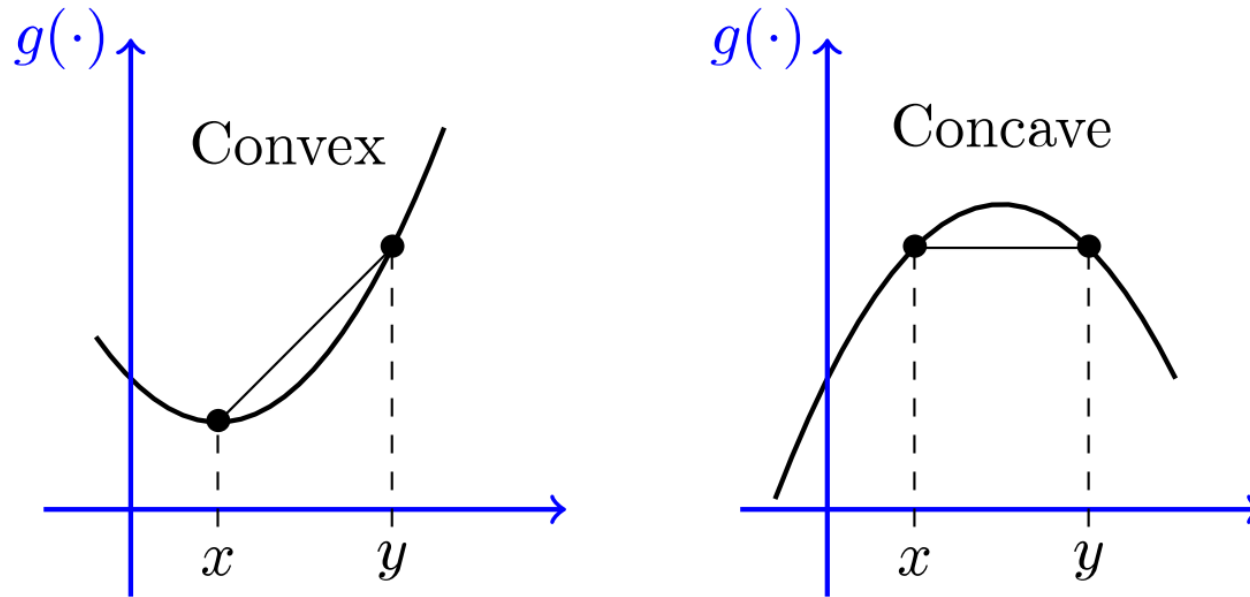
Consider the interval $0 \leq x \leq 1$:

$$f(\mathbb{E}[X]) = f(\mathbb{E}[X]) = f(0.5) = 0.25$$

$$\mathbb{E}[f(X)] = \frac{\int_0^1 f(x)}{1-0} = \frac{x^3}{3} \Big|_0^1 = 0.33$$



Background: Convex / Concave function



Definition 6.3

Consider a function $g : I \rightarrow \mathbb{R}$, where I is an interval in \mathbb{R} . We say that g is a **convex** function if, for any two points x and y in I and any $\alpha \in [0, 1]$, we have

$$g(\alpha x + (1 - \alpha)y) \leq \alpha g(x) + (1 - \alpha)g(y).$$

We say that g is **concave** if

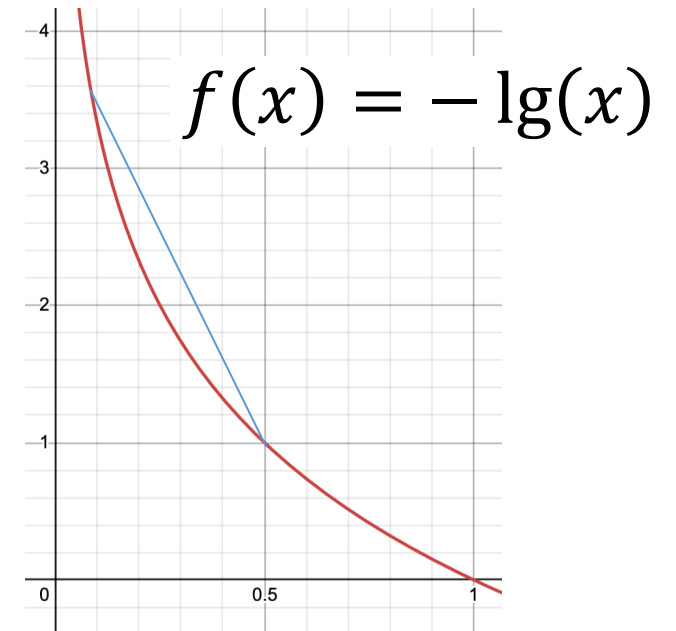
$$g(\alpha x + (1 - \alpha)y) \geq \alpha g(x) + (1 - \alpha)g(y).$$

Information inequality $D_{KL}(p||q) \geq 0$

Ingredients:

1. $-\lg(x)$ is convex
2. Jensen's inequality $\mathbb{E}[f(X)] \geq f(\mathbb{E}[X])$

$$D_{KL}(p||q) = \mathbb{E}_p \left[\lg \left(\frac{p(X)}{q(X)} \right) \right]$$
$$= \text{?}$$



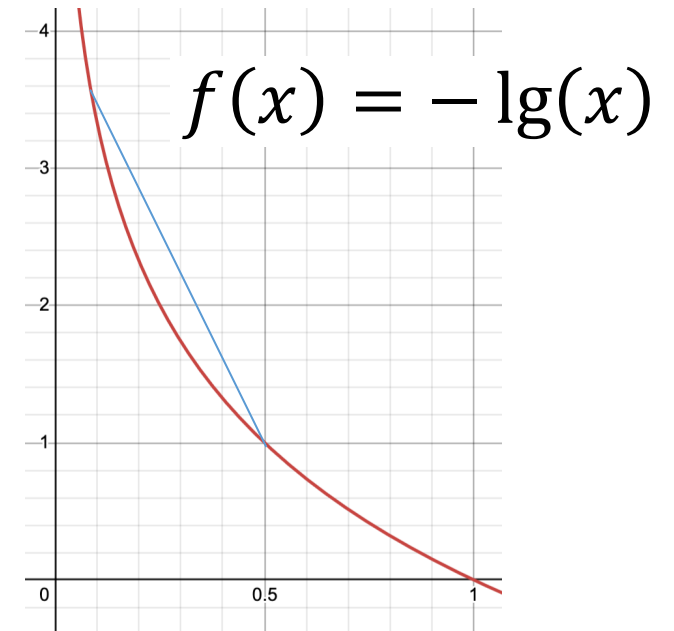
Information inequality $D_{KL}(p||q) \geq 0$

Ingredients:

1. $-\lg(x)$ is convex
2. Jensen's inequality $\mathbb{E}[f(X)] \geq f(\mathbb{E}[X])$

$$\begin{aligned} D_{KL}(p||q) &= \mathbb{E}_p \left[\lg \left(\frac{p(X)}{q(X)} \right) \right] \\ &= \mathbb{E}_p \left[-\lg \left(\frac{q(X)}{p(X)} \right) \right] \\ &\geq -\lg \left(\mathbb{E}_p \left[\frac{q(X)}{p(X)} \right] \right) = -\lg \left(\underbrace{\sum_x p(x) \cdot \frac{q(x)}{p(x)}}_{=1} \right) = 0 \end{aligned}$$

$D_{KL}(p||q) = 0$ iff ?



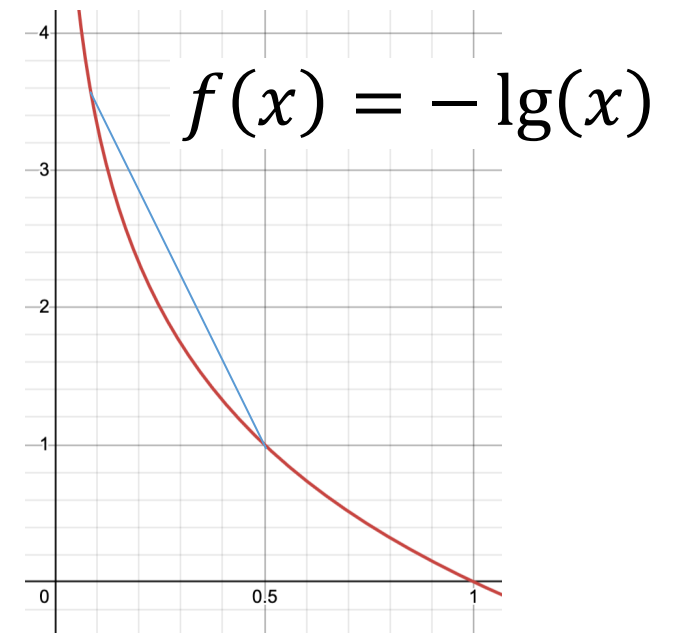
Information inequality $D_{KL}(p||q) \geq 0$

Ingredients:

1. $-\lg(x)$ is convex
2. Jensen's inequality $\mathbb{E}[f(X)] \geq f(\mathbb{E}[X])$

$$\begin{aligned} D_{KL}(p||q) &= \mathbb{E}_p \left[\lg \left(\frac{p(X)}{q(X)} \right) \right] \\ &= \mathbb{E}_p \left[-\lg \left(\frac{q(X)}{p(X)} \right) \right] \\ &\geq -\lg \left(\mathbb{E}_p \left[\frac{q(X)}{p(X)} \right] \right) = -\lg \left(\underbrace{\sum_x \cancel{p(x)} \cdot \frac{q(x)}{\cancel{p(x)}}}_{=1} \right) = 0 \end{aligned}$$

$D_{KL}(p||q) = 0$ iff $q(x) = p(x)$ for all x .



Mutual information as relative entropy and thus ≥ 0

Given two RVs X and Y , **mutual information** is the amount of information that Y provides about X when Y is observed, but X is not.

$$I(X; Y) := H(X) - H(X|Y)$$

≥ 0 ?

notation $x \in \mathcal{X}, y \in \mathcal{Y}$

Mutual information as relative entropy and thus ≥ 0

Given two RVs X and Y , **mutual information** is the amount of information that Y provides about X when Y is observed, but X is not.

$$I(X; Y) := H(X) - H(X|Y)$$

notation $x \in \mathcal{X}, y \in \mathcal{Y}$

$$\begin{aligned} &= \sum_x p(x) \cdot \lg\left(\frac{1}{p(x)}\right) - \sum_{x,y} p(x,y) \cdot \lg\left(\frac{1}{p(x|y)}\right) \\ &= \sum_{x,y} p(x,y) \cdot \lg\left(\frac{1}{p(x)}\right) - \sum_{x,y} p(x,y) \cdot \lg\left(\frac{1}{p(x|y)}\right) \\ &= \sum_{x,y} p(x,y) \cdot \lg\left(\frac{p(x|y)}{p(x)}\right) \\ &= \sum_{x,y} p(x,y) \cdot \lg\left(\frac{p(x,y)}{p(x) \cdot p(y)}\right) = \text{?} \end{aligned}$$

Mutual information as relative entropy and thus ≥ 0

Given two RVs X and Y , **mutual information** is the amount of information that Y provides about X when Y is observed, but X is not.

$$I(X; Y) := H(X) - H(X|Y)$$

notation $x \in \mathcal{X}, y \in \mathcal{Y}$

$$= \sum_x p(x) \cdot \lg\left(\frac{1}{p(x)}\right) - \sum_{x,y} p(x,y) \cdot \lg\left(\frac{1}{p(x|y)}\right)$$

$$= \sum_{x,y} p(x,y) \cdot \lg\left(\frac{1}{p(x)}\right) - \sum_{x,y} p(x,y) \cdot \lg\left(\frac{1}{p(x|y)}\right)$$

$$= \sum_{x,y} p(x,y) \cdot \lg\left(\frac{p(x|y)}{p(x)}\right)$$

$$= \sum_{x,y} p(x,y) \cdot \lg\left(\frac{p(x,y)}{p(x) \cdot p(y)}\right) = D_{\text{KL}}(p(x,y) || p(x) \cdot p(y)) \geq 0$$

When equality?



Mutual information is the **relative entropy** between joint distribution and product of their marginal distributions!

Mutual information as relative entropy and thus ≥ 0

Given two RVs X and Y , **mutual information** is the amount of information that Y provides about X when Y is observed, but X is not.

$$I(X; Y) := H(X) - H(X|Y)$$

notation $x \in \mathcal{X}, y \in \mathcal{Y}$

$$= \sum_x p(x) \cdot \lg\left(\frac{1}{p(x)}\right) - \sum_{x,y} p(x,y) \cdot \lg\left(\frac{1}{p(x|y)}\right)$$

$$= \sum_{x,y} p(x,y) \cdot \lg\left(\frac{1}{p(x)}\right) - \sum_{x,y} p(x,y) \cdot \lg\left(\frac{1}{p(x|y)}\right)$$

$$= \sum_{x,y} p(x,y) \cdot \lg\left(\frac{p(x|y)}{p(x)}\right)$$

$$= \sum_{x,y} p(x,y) \cdot \lg\left(\frac{p(x,y)}{p(x) \cdot p(y)}\right) = D_{\text{KL}}(p(x,y) || p(x) \cdot p(y)) \geq 0$$

equality when X and Y
are independent!

alternative notation:
 $D_{\text{KL}}(p_{X,Y} || p_X \cdot p_Y)$

Mutual information is the **relative entropy** between joint distribution and product of their marginal distributions!

Conditioning reduces entropy, in expectation

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

(follows from $I(X; Y) = H(X) - H(X|Y) \geq 0$)

The **nonnegativity of mutual information** implies that on average the entropy of X conditioned on the observation $Y = y$ is \leq than the entropy of X (which intuitively makes sense: getting more information only reduces uncertainty, in expectation).

But importantly, the inequality is applied to averaged quantities. It is still possible that there is new rare evidence y for which:



$$H(X) < H(X|Y = y)$$

Example: in a court case, specific new evidence might increase uncertainty, but on the average evidence decreases uncertainty.

But new concrete evidence may increase entropy



EXAMPLE 6: Consider the joint ensemble (X, Y) with Boolean domains $\mathcal{X} = \mathcal{Y} = \{0, 1\}$ and following joint distribution.

		y	
		0	1
x	0	$\frac{1}{2}$	$\frac{1}{4}$
	1	0	$\frac{1}{4}$

		y	
		0	1
x	0	■	■
	1		■

$$H(X) = \quad ?$$

$$H(X|y = 0) = \quad ?$$

$$H(X|y = 1) = \quad ?$$

$$H(X|Y) = \quad ?$$

But new concrete evidence may increase entropy



EXAMPLE 6: Consider the joint ensemble (X, Y) with Boolean domains $\mathcal{X} = \mathcal{Y} = \{0,1\}$ and following joint distribution.

$p(x, y)$	y	Σ
	0	1
0	$\frac{1}{2}$	$\frac{1}{4}$
1	0	$\frac{1}{4}$
Σ	$\frac{1}{2}$	$\frac{1}{2}$

	y
	0
0	■
1	■

$$H(X) = \frac{3}{4} \lg\left(\frac{4}{3}\right) + \frac{1}{4} \lg(4) = 0.811$$

$$H(X|y = 0) = 0$$

$$H(X|y = 1) = 1$$

$$H(X|Y) = \frac{1}{2} \underbrace{H(X|y = 0)}_0 + \frac{1}{2} \underbrace{H(X|y = 1)}_1 = 0.5$$

$$H(X|Y) \leq H(X) < H(X|y = 1)$$

$$0.5 \quad 0.811 \quad 1$$

Part 1: Theory

L08: Basics of entropy (4/6)

[multivariate entropies]

Wolfgang Gatterbauer, Javed Aslam

cs7840 Foundations and Applications of Information Theory (fa24)

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

9/30/2024

Pre-class conversations

- Last class recapitulation
- Slide decks: more overall consistent updates coming
- Next scribe correct towards end of week
- Python scripts also coming soon

- Today so far: compression
- Today next:
 - Multi-variate entropies
 - Markov Chains & Data Processing inequality

Three-term (multivariate) entropies,
conditional mutual information,
interaction information

Conditioning & chain rules

$$H(Y|X) = \sum_x p(x) \cdot H(Y|X = x)$$

Conditional entropy $H(Y|X)$: the expected amount of information needed to describe the outcome of RV Y , given that the value of another RV X is known

$$H(X, Y) = ?$$

Two-variable chain rule

Conditioning & chain rules

$$H(Y|X) = \sum_x p(x) \cdot H(Y|X = x)$$

Conditional entropy $H(Y|X)$: the expected amount of information needed to describe the outcome of RV Y , given that the value of another RV X is known

$$H(X, Y) = H(X) + H(Y|X) \quad \text{Two-variable chain rule}$$

$$H(X, Y|Z) = ?$$

Conditional chain rule.

$\mathbb{E}[H(X, Y)|Z]$ ← Notice the implied precedence rule

Conditional joint entropy $H(X, Y|Z)$: expected joint entropy of X and Y together, given that Z is known

Conditioning & chain rules

$$H(Y|X) = \sum_x p(x) \cdot H(Y|X = x)$$

Conditional entropy $H(Y|X)$: the expected amount of information needed to describe the outcome of RV Y , given that the value of another RV X is known

$$H(X, Y) = H(X) + H(Y|X) \quad \text{Two-variable chain rule}$$

$$\underbrace{H(X, Y|Z)} = H(X|Z) + H(Y|X, Z)$$

Conditional chain rule.

$\mathbb{E}[H(X, Y)|Z]$ ← Notice the implied precedence rule

Conditioning on an event creates a new probability space where the same probability concepts apply.

Conditional joint entropy $H(X, Y|Z)$: expected joint entropy of X and Y together, given that Z is known

$$H(X, Y|Z) \quad ? \quad H(X|Z) + H(Y|Z)$$

\leq or \geq

Conditioning & chain rules

$$H(Y|X) = \sum_x p(x) \cdot H(Y|X = x)$$

Conditional entropy $H(Y|X)$: the expected amount of information needed to describe the outcome of RV Y , given that the value of another RV X is known

$$H(X, Y) = H(X) + H(Y|X) \quad \text{Two-variable chain rule}$$

$$\underbrace{H(X, Y|Z)} = H(X|Z) + H(Y|X, Z)$$

Conditional chain rule.

$\mathbb{E}[H(X, Y)|Z]$ ← Notice the implied precedence rule

Conditioning on an event creates a new probability space where the same probability concepts apply.

Conditional joint entropy $H(X, Y|Z)$: expected joint entropy of X and Y together, given that Z is known

$$H(X, Y|Z) \leq H(X|Z) + H(Y|Z)$$

Equality holds if X and Y are conditionally independent, given Z (Proof similar to unconditional case).

$$H(X, Y, Z) = ?$$

Three-variable chain rule

Conditioning & chain rules

$$H(Y|X) = \sum_x p(x) \cdot H(Y|X = x)$$

Conditional entropy $H(Y|X)$: the expected amount of information needed to describe the outcome of RV Y , given that the value of another RV X is known

$$H(X, Y) = H(X) + H(Y|X) \quad \text{Two-variable chain rule}$$

$$\underbrace{H(X, Y|Z)} = H(X|Z) + H(Y|X, Z)$$

Conditional chain rule.

$\mathbb{E}[H(X, Y)|Z]$ ← Notice the implied precedence rule

Conditioning on an event creates a new probability space where the same probability concepts apply.

Conditional joint entropy $H(X, Y|Z)$: expected joint entropy of X and Y together, given that Z is known

$$H(X, Y|Z) \leq H(X|Z) + H(Y|Z)$$

Equality holds if X and Y are conditionally independent, given Z (Proof similar to unconditional case).

$$H(X, Y, Z) = H(X) + H(Y|X) + H(Z|X, Y) \quad \text{Three-variable chain rule}$$

Conditional mutual information & interaction information

$$\underbrace{I(X; Y|Z)}_{\mathbb{E}[I(X; Y)|Z]} = ?$$

Conditional mutual information $I(X; Y|Z)$:
expected mutual information of X and Y ,
given Z is known

Conditional mutual information & interaction information

$$\underbrace{I(X; Y|Z)}_{\mathbb{E}[I(X; Y)|Z]} = H(X|Z) + H(Y|Z) - \underbrace{H(X, Y|Z)}_{H(Y|Z) + H(X|Y, Z)}$$

Conditional mutual information $I(X; Y|Z)$:
expected mutual information of X and Y ,
given Z is known

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

$$J(X; Y; Z) = ?$$

Interaction information (often also called "mutual information*"): measures the influence of a variable Z on the amount of information shared between X and Y .

* Alternative notations include $J(X; Y; Z)$ and $R(X; Y; Z)$. We don't recommend calling it "mutual information" and thus also replace the more common notation $I(X; Y; Z)$ with $J(X; Y; Z)$.
Some sources prefer not to even define that measure at all (we will discuss the reasons in a bit) https://en.wikipedia.org/wiki/Interaction_information.
Gatterbauer, Aslam. Foundations and Applications of Information Theory: <https://northeastern-datalab.github.io/cs7840/>

Part 1: Theory

L09: Basics of entropy (5/6)

[multivariate entropies, interaction information, Markov chains, data processing inequality]

Wolfgang Gatterbauer, Javed Aslam

cs7840 Foundations and Applications of Information Theory (fa24)

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

10/2/2024

Pre-class conversations

- Last class recapitulation
- Slide decks: please continue checking for errors / inconsistencies / unclear details
- Your own Python scripts could be part of your next scribes!

- Today:
 - Multi-variate entropies
 - Markov Chains, Data Processing inequality, sufficient statistics
 - Possibly starting with Part 2: axioms

Three-term (multivariate) entropies,
conditional mutual information,
interaction information
(continued)

Conditional mutual information & interaction information

$$\underbrace{I(X; Y|Z)}_{\mathbb{E}[I(X; Y)|Z]} = ?$$

Conditional mutual information $I(X; Y|Z)$:
expected mutual information of X and Y ,
given Z is known

Conditional mutual information & interaction information

$$\underbrace{I(X; Y|Z)}_{\mathbb{E}[I(X; Y)|Z]} = H(X|Z) + H(Y|Z) - \underbrace{H(X, Y|Z)}_{H(Y|Z) + \underbrace{H(X|Y, Z)}_{H(X|(Y, Z))}}$$

Conditional mutual information $I(X; Y|Z)$: expected mutual information of X and Y , given Z is known

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

$$J(X; Y; Z) = ?$$

Interaction information (often also called "mutual information*"): measures the influence of a variable Z on the amount of information shared between X and Y .

* Alternative notations include $J(X; Y; Z)$ and $R(X; Y; Z)$. We don't recommend calling it "mutual information" and thus also replace the more common notation $I(X; Y; Z)$ with $J(X; Y; Z)$. Some sources prefer not to even define that measure at all (we will discuss the reasons in a bit) https://en.wikipedia.org/wiki/Interaction_information. Gatterbauer, Aslam. Foundations and Applications of Information Theory: <https://northeastern-datalab.github.io/cs7840/>

Conditional mutual information & interaction information

$$\underbrace{I(X; Y|Z)}_{\mathbb{E}[I(X; Y)|Z]} = H(X|Z) + H(Y|Z) - \underbrace{H(X, Y|Z)}_{H(Y|Z) + H(X|Y, Z)}$$

Conditional mutual information $I(X; Y|Z)$: expected mutual information of X and Y , given Z is known

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

$$J(X; Y; Z) = I(X; Y) - I(X; Y|Z)$$

Interaction information (often also called "mutual information*"): measures the influence of a variable Z on the amount of information shared between X and Y .

Is it symmetric in all the variables ?

* Alternative notations include $J(X; Y; Z)$ and $R(X; Y; Z)$. We don't recommend calling it "mutual information" and thus also replace the more common notation $I(X; Y; Z)$ with $J(X; Y; Z)$. Some sources prefer not to even define that measure at all (we will discuss the reasons in a bit) https://en.wikipedia.org/wiki/Interaction_information. Gatterbauer, Aslam. Foundations and Applications of Information Theory: <https://northeastern-datalab.github.io/cs7840/>

Conditional mutual information & interaction information

$$\underbrace{I(X; Y|Z)}_{\mathbb{E}[I(X; Y)|Z]} = H(X|Z) + H(Y|Z) - \underbrace{H(X, Y|Z)}_{H(Y|Z) + H(X|Y, Z)}$$

Conditional mutual information $I(X; Y|Z)$: expected mutual information of X and Y , given Z is known

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

$$J(X; Y; Z) = I(X; Y) - I(X; Y|Z)$$

Interaction information (often also called "mutual information*"): measures the influence of a variable Z on the amount of information shared between X and Y .

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

(...) thus symmetric in all 3 variables!

* Alternative notations include $J(X; Y; Z)$ and $R(X; Y; Z)$. We don't recommend calling it "mutual information" and thus also replace the more common notation $I(X; Y; Z)$ with $J(X; Y; Z)$.

Some sources prefer not to even define that measure at all (we will discuss the reasons in a bit) https://en.wikipedia.org/wiki/Interaction_information.

Gatterbauer, Aslam. Foundations and Applications of Information Theory: <https://northeastern-datalab.github.io/cs7840/>

Interaction information example



EXAMPLE: Consider the joint ensemble (X, Y, Z) with Boolean domains $\mathcal{X} = \mathcal{Y} = \mathcal{Z} = \{0,1\}$. X and Y are independent uniform binary variables. And let Z be the XOR of X and Y : $z = \text{XOR}(x, y)$, or equally, $z = x + y \pmod{2}$.

x	y	z
0	0	0
0	1	1
1	0	1
1	1	0

Interaction information example



EXAMPLE: Consider the joint ensemble (X, Y, Z) with Boolean domains $\mathcal{X} = \mathcal{Y} = \mathcal{Z} = \{0,1\}$. X and Y are independent uniform binary variables. And let Z be the XOR of X and Y : $z = \text{XOR}(x, y)$, or equally, $z = x + y \pmod{2}$.

Thus any 2 variables functionally determine the 3rd, e.g. $(x, z) \rightarrow y$!

x	y	z	p
0	0	0	$\frac{1}{4}$
0	1	1	$\frac{1}{4}$
1	0	1	$\frac{1}{4}$
1	1	0	$\frac{1}{4}$
0	0	1	0
...	0

Interaction information example



EXAMPLE: Consider the joint ensemble (X, Y, Z) with Boolean domains $\mathcal{X} = \mathcal{Y} = \mathcal{Z} = \{0,1\}$. X and Y are independent uniform binary variables. And let Z be the XOR of X and Y : $z = \text{XOR}(x, y)$, or equally, $z = x + y \pmod{2}$.

x	y	z	p
0	0	0	$\frac{1}{4}$
0	1	1	$\frac{1}{4}$
1	0	1	$\frac{1}{4}$
1	1	0	$\frac{1}{4}$
0	0	1	0
...	0

$$H(X) = ?$$

Thus any 2 variables functionally determine the 3rd, e.g. $(x, z) \rightarrow y$!

$$H(X|Y) = ?$$
$$I(X; Y) = ?$$

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

$$J(X; Y; Z) = ?$$

Interaction information example



EXAMPLE: Consider the joint ensemble (X, Y, Z) with Boolean domains $\mathcal{X} = \mathcal{Y} = \mathcal{Z} = \{0,1\}$. X and Y are independent uniform binary variables. And let Z be the XOR of X and Y : $z = \text{XOR}(x, y)$, or equally, $z = x + y \bmod 2$.

x	y	z	p
0	0	0	$\frac{1}{4}$
0	1	1	$\frac{1}{4}$
1	0	1	$\frac{1}{4}$
1	1	0	$\frac{1}{4}$

$$H(X) = 1$$

Similarly, $H(Y) = 1$ and $H(Z) = 1$

$$H(X|Y) = H(X) = 1$$

Similarly, all variables are pairwise independent

$$I(X; Y) = 0$$

$$H(X|Y, Z) = 0$$

$$I(X; Y|Z) = 1$$

Thus, if Z is observed, then X and Y become dependent:
(knowing $X = x$ and $Z = z$, tells you what Y is: $y = z - x \bmod 2$)

Thus the **conditional mutual information** is bigger than the **unconditional mutual information**: $I(X; Y|Z) > I(X; Y)$

$$J(X; Y; Z) = I(X; Y) - I(X; Y|Z) = -1$$

When VENN diagrams confuse more than help (?)

EXAMPLE: Consider the joint ensemble (X, Y, Z) with Boolean domains $\mathcal{X} = \mathcal{Y} = \mathcal{Z} = \{0,1\}$. X and Y are independent uniform binary variables. And let Z be the XOR of X and Y : $z = \text{XOR}(x, y)$, or equally, $z = x + y \text{ mod } 2$.

x	y	z	p
0	0	0	$\frac{1}{4}$
0	1	1	$\frac{1}{4}$
1	0	1	$\frac{1}{4}$
1	1	0	$\frac{1}{4}$

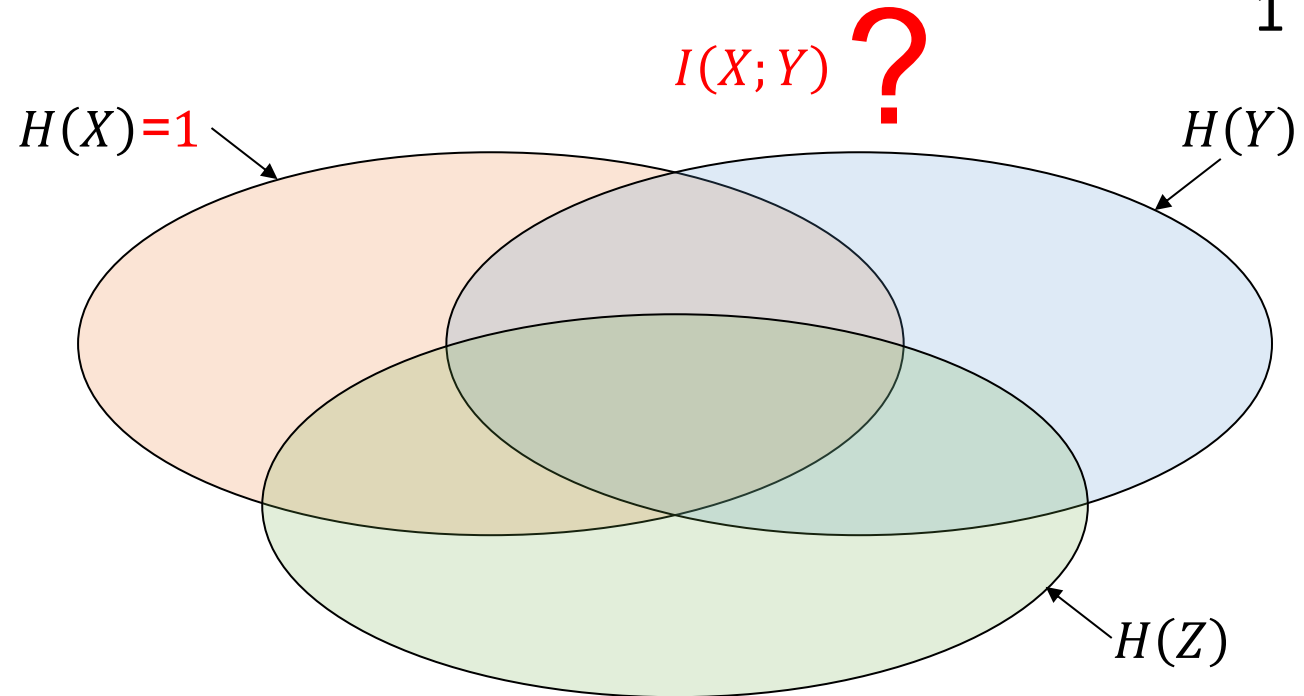
$$H(X) = 1$$

$$H(X|Y) = H(X) = 1$$

$$I(X;Y) = 0$$

$$H(X|Y, Z) = 0$$

$$I(X;Y|Z) = 1$$



$$J(X;Y;Z) = I(X;Y) - I(X;Y|Z) = -1$$

When VENN diagrams confuse more than help (?)

EXAMPLE: Consider the joint ensemble (X, Y, Z) with Boolean domains $\mathcal{X} = \mathcal{Y} = \mathcal{Z} = \{0,1\}$. X and Y are independent uniform binary variables. And let Z be the XOR of X and Y : $z = \text{XOR}(x, y)$, or equally, $z = x + y \pmod 2$.

x	y	z	p
0	0	0	$\frac{1}{4}$
0	1	1	$\frac{1}{4}$
1	0	1	$\frac{1}{4}$
1	1	0	$\frac{1}{4}$

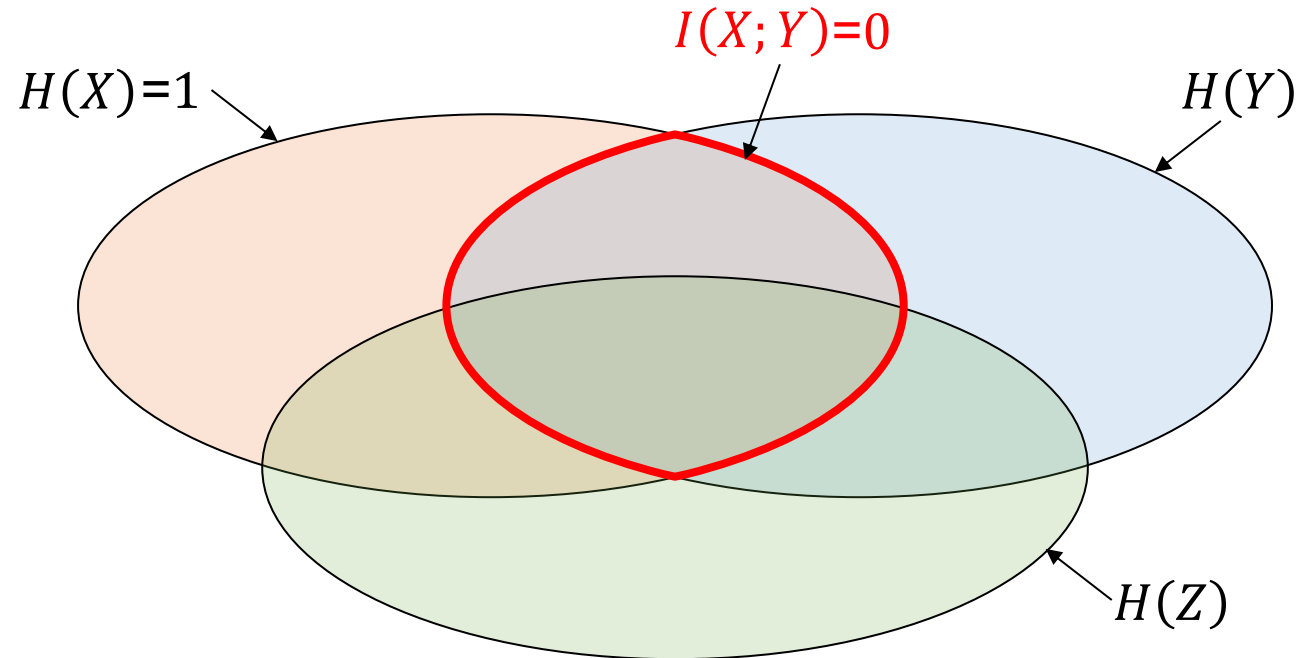
$$H(X) = 1$$

$$H(X|Y) = H(X) = 1$$

$$I(X;Y) = 0$$

$$H(X|Y, Z) = 0$$

$$I(X;Y|Z) = 1$$



$$J(X;Y;Z) = I(X;Y) - I(X;Y|Z) = -1$$

When VENN diagrams confuse more than help (?)

EXAMPLE: Consider the joint ensemble (X, Y, Z) with Boolean domains $\mathcal{X} = \mathcal{Y} = \mathcal{Z} = \{0,1\}$. X and Y are independent uniform binary variables. And let Z be the XOR of X and Y : $z = \text{XOR}(x, y)$, or equally, $z = x + y \text{ mod } 2$.

x	y	z	p
0	0	0	$\frac{1}{4}$
0	1	1	$\frac{1}{4}$
1	0	1	$\frac{1}{4}$
1	1	0	$\frac{1}{4}$

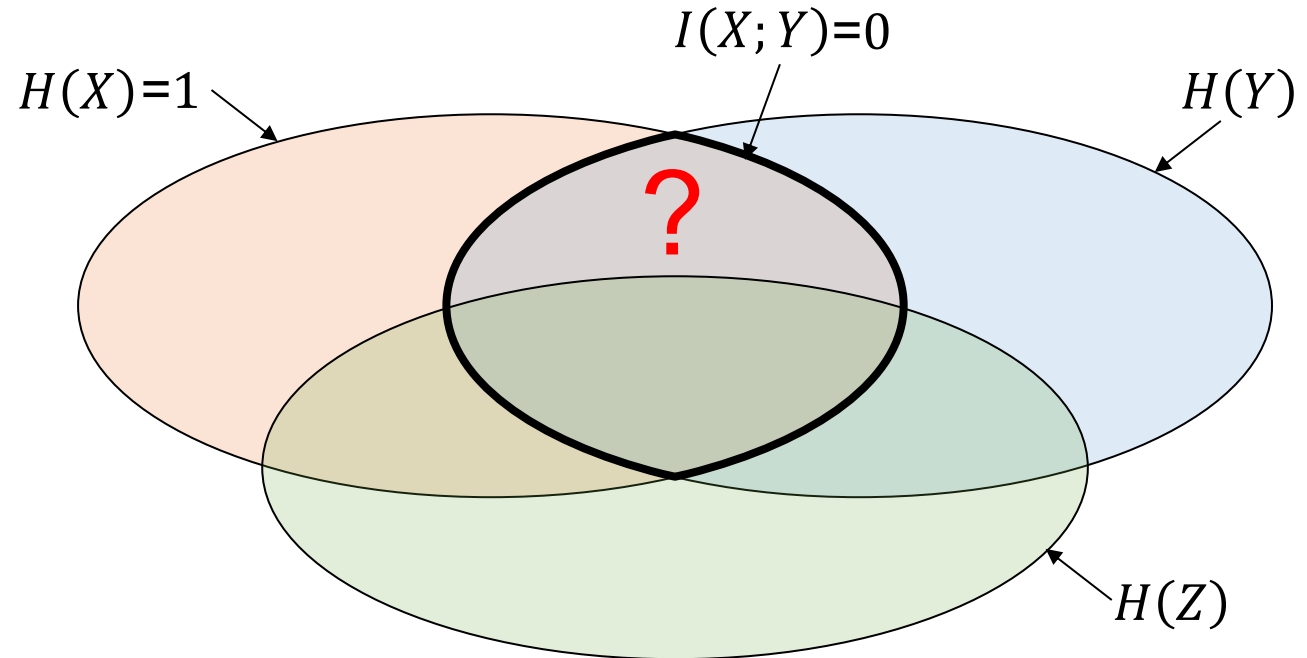
$$H(X) = 1$$

$$H(X|Y) = H(X) = 1$$

$$I(X;Y) = 0$$

$$H(X|Y, Z) = 0$$

$$I(X;Y|Z) = 1$$



$$J(X;Y;Z) = I(X;Y) - I(X;Y|Z) = -1$$

When VENN diagrams confuse more than help (?)

EXAMPLE: Consider the joint ensemble (X, Y, Z) with Boolean domains $\mathcal{X} = \mathcal{Y} = \mathcal{Z} = \{0,1\}$. X and Y are independent uniform binary variables. And let Z be the XOR of X and Y : $z = \text{XOR}(x, y)$, or equally, $z = x + y \pmod 2$.

x	y	z	p
0	0	0	$\frac{1}{4}$
0	1	1	$\frac{1}{4}$
1	0	1	$\frac{1}{4}$
1	1	0	$\frac{1}{4}$

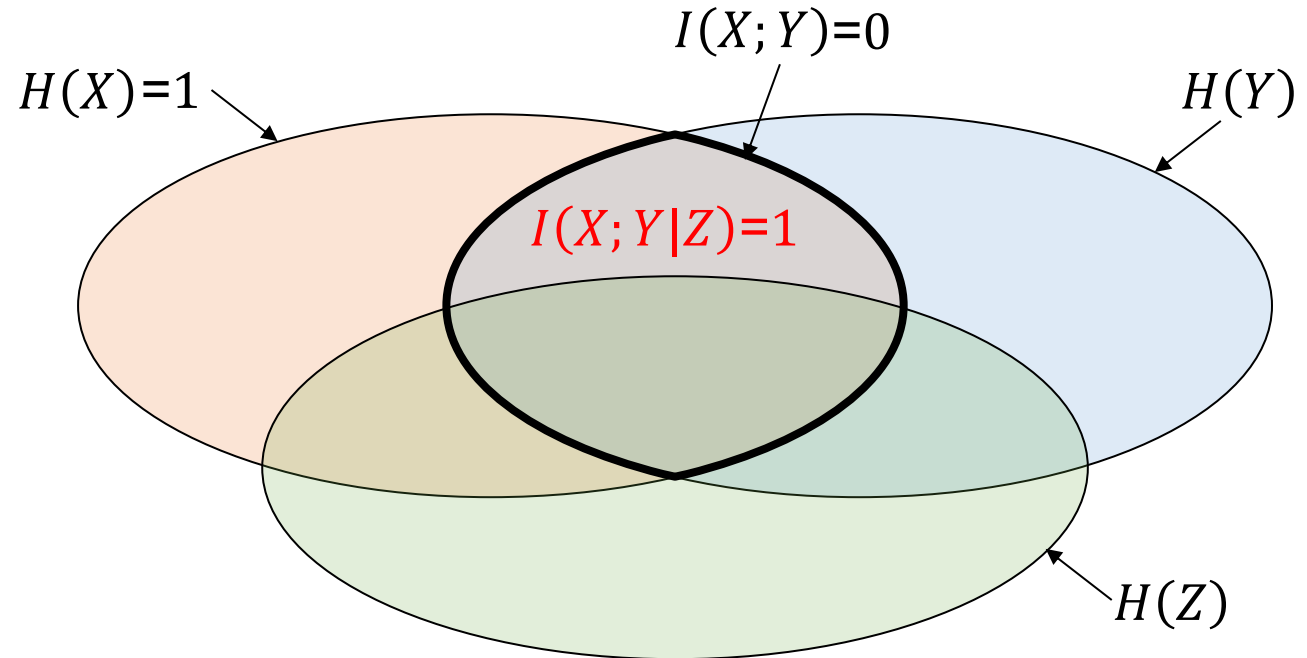
$$H(X) = 1$$

$$H(X|Y) = H(X) = 1$$

$$I(X;Y) = 0$$

$$H(X|Y, Z) = 0$$

$$I(X;Y|Z) = 1$$



$$J(X;Y;Z) = I(X;Y) - I(X;Y|Z) = -1$$

When VENN diagrams confuse more than help (?)

EXAMPLE: Consider the joint ensemble (X, Y, Z) with Boolean domains $\mathcal{X} = \mathcal{Y} = \mathcal{Z} = \{0,1\}$. X and Y are independent uniform binary variables. And let Z be the XOR of X and Y : $z = \text{XOR}(x, y)$, or equally, $z = x + y \pmod{2}$.

x	y	z	p
0	0	0	$\frac{1}{4}$
0	1	1	$\frac{1}{4}$
1	0	1	$\frac{1}{4}$
1	1	0	$\frac{1}{4}$

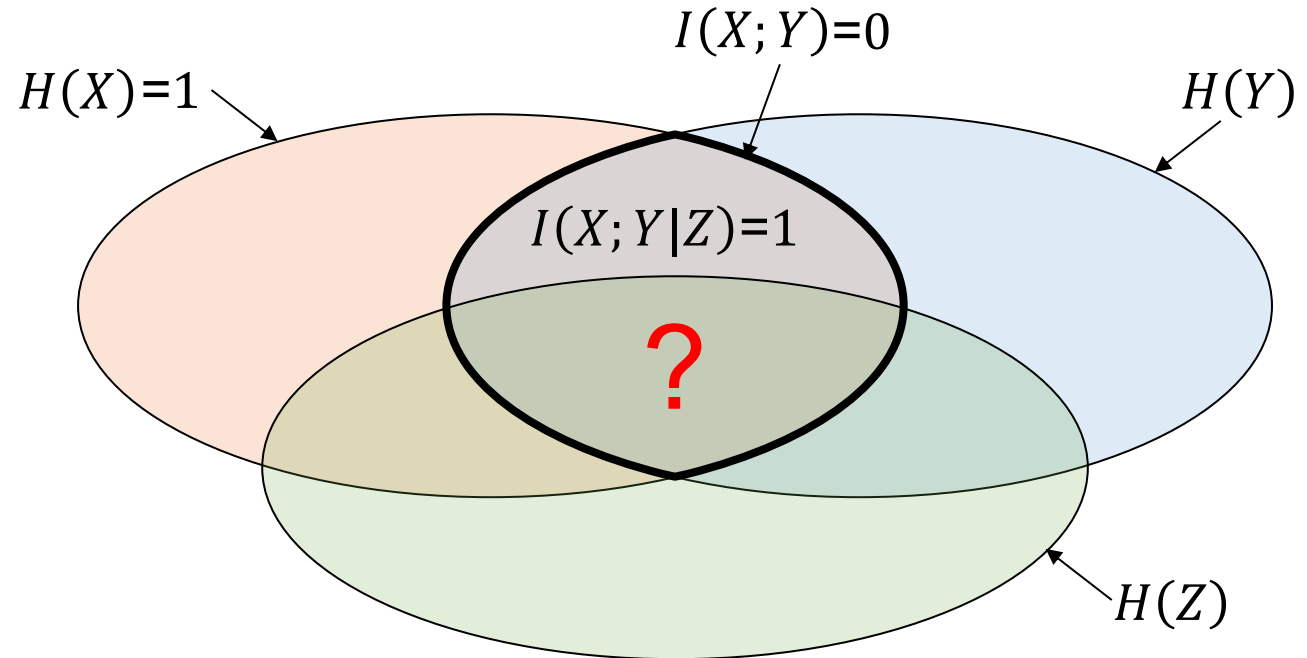
$$H(X) = 1$$

$$H(X|Y) = H(X) = 1$$

$$I(X;Y) = 0$$

$$H(X|Y, Z) = 0$$

$$I(X;Y|Z) = 1$$



$$J(X;Y;Z) = I(X;Y) - I(X;Y|Z) = -1$$

When VENN diagrams confuse more than help (?)

EXAMPLE: Consider the joint ensemble (X, Y, Z) with Boolean domains $\mathcal{X} = \mathcal{Y} = \mathcal{Z} = \{0,1\}$. X and Y are independent uniform binary variables. And let Z be the XOR of X and Y : $z = \text{XOR}(x, y)$, or equally, $z = x + y \pmod 2$.

x	y	z	p
0	0	0	$\frac{1}{4}$
0	1	1	$\frac{1}{4}$
1	0	1	$\frac{1}{4}$
1	1	0	$\frac{1}{4}$

$$H(X) = 1$$

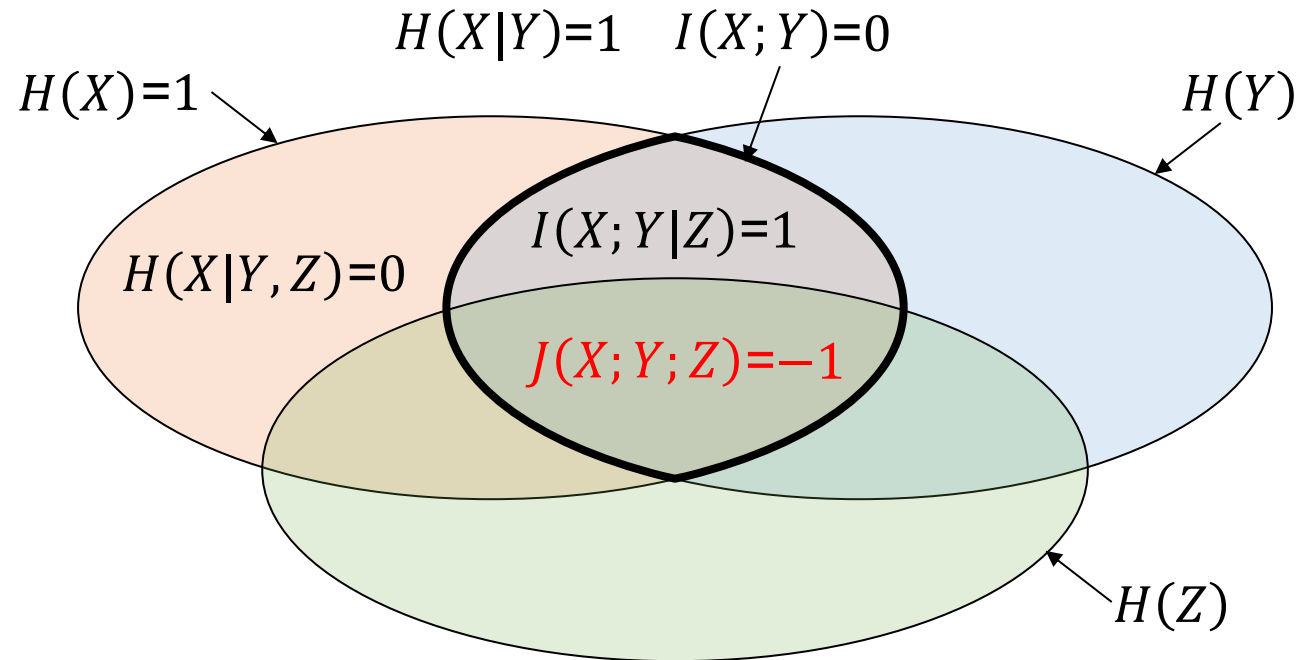
$$H(X|Y) = H(X) = 1$$

$$I(X;Y) = 0$$

$$H(X|Y, Z) = 0$$

$$I(X;Y|Z) = 1$$

$$J(X;Y;Z) = I(X;Y) - I(X;Y|Z) = -1$$



\Rightarrow VENN diagrams applied to joint entropies with ≥ 2 variables can mislead

[MacKay'02] on VENN diagrams and three-term entropies

Figure 8.1. The relationship between joint information, marginal entropy, conditional entropy and mutual entropy.

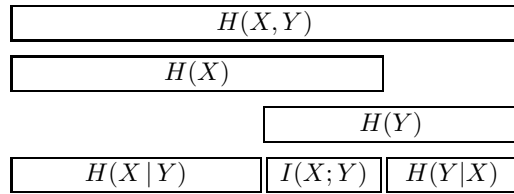
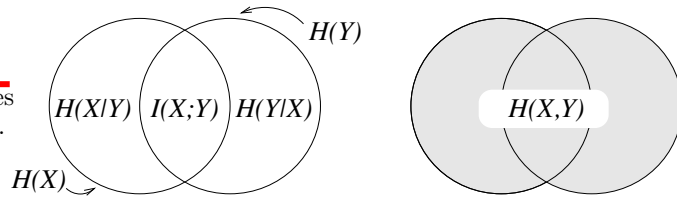


Figure 8.2. A misleading representation of entropies (contrast with figure 8.1).



Exercise 8.8. [3, p.143] Many texts draw figure 8.1 in the form of a Venn diagram (figure 8.2). Discuss why this diagram is a misleading representation of entropies. Hint: consider the three-variable ensemble XYZ in which $x \in \{0, 1\}$ and $y \in \{0, 1\}$ are independent binary variables and $z \in \{0, 1\}$ is defined to be $z = x + y \text{ mod } 2$.

[MacKay'02] on VENN diagrams and three-term entropies

Figure 8.1. The relationship between joint information, marginal entropy, conditional entropy and mutual entropy.

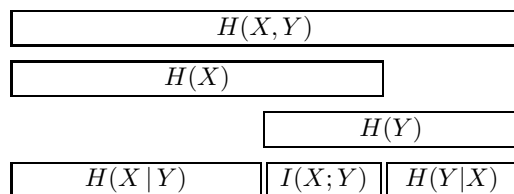
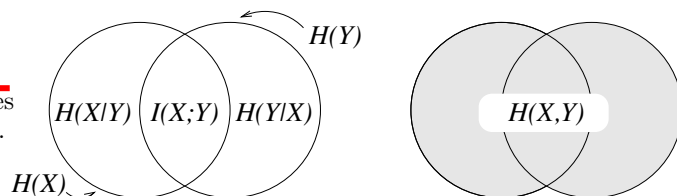
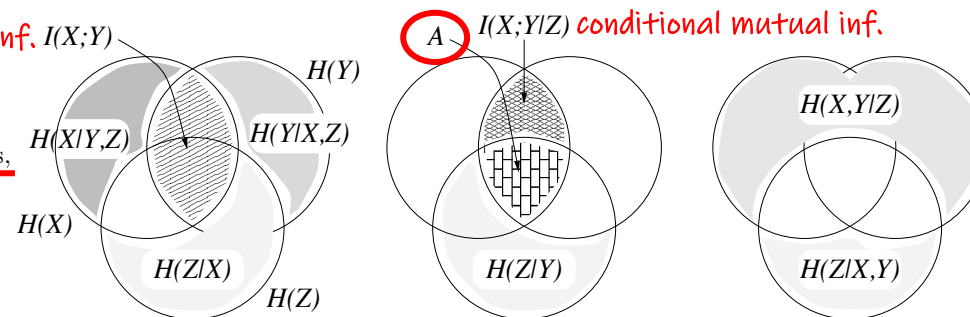


Figure 8.2. A misleading representation of entropies (contrast with figure 8.1).



unconditional mutual inf. $I(X; Y)$

Figure 8.3. A misleading representation of entropies, continued.



Secondly, the depiction in terms of Venn diagrams encourages one to believe that all the areas correspond to positive quantities. In the special case of two random variables it is indeed true that $H(X|Y)$, $I(X; Y)$ and $H(Y|X)$ are positive quantities. But as soon as we progress to three-variable ensembles, we obtain a diagram with positive-looking areas that may actually correspond to negative quantities. Figure 8.3 correctly shows relationships such as

$$H(X) + H(Z|X) + H(Y|X, Z) = H(X, Y, Z). \quad (8.31)$$

But it gives the misleading impression that the conditional mutual information $I(X; Y|Z)$ is less than the mutual information $I(X; Y)$. In fact the area labelled A can correspond to a negative quantity. Consider the joint ensemble (X, Y, Z) in which $x \in \{0, 1\}$ and $y \in \{0, 1\}$ are independent binary variables and $z \in \{0, 1\}$ is defined to be $z = x + y \text{ mod } 2$. Then clearly $H(X) = H(Y) = 1$ bit. Also $H(Z) = 1$ bit. And $H(Y|X) = H(Y) = 1$ since the two variables are independent. So the mutual information between X and Y is zero. $I(X; Y) = 0$. However, if z is observed, X and Y become dependent — knowing x , given z , tells you what y is: $y = z - x \text{ mod } 2$. So $I(X; Y|Z) = 1$ bit. Thus the area labelled A must correspond to -1 bits for the figure to give the correct answers.

The above example is not at all a capricious or exceptional illustration. The binary symmetric channel with input X , noise Y , and output Z is a situation in which $I(X; Y) = 0$ (input and noise are independent) but $I(X; Y|Z) > 0$ (once you see the output, the unknown input and the unknown noise are intimately related!).

The Venn diagram representation is therefore valid only if one is aware that positive areas may represent negative quantities. With this proviso kept in mind, the interpretation of entropies in terms of sets can be helpful (Yeung, 1991).

Exercise 8.8. [3, p.143] Many texts draw figure 8.1 in the form of a Venn diagram (figure 8.2). Discuss why this diagram is a misleading representation of entropies. Hint: consider the three-variable ensemble XYZ in which $x \in \{0, 1\}$ and $y \in \{0, 1\}$ are independent binary variables and $z \in \{0, 1\}$ is defined to be $z = x + y \text{ mod } 2$.

...

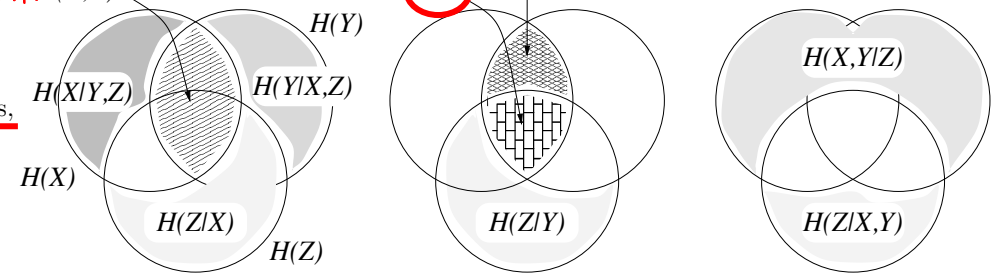
Solution to exercise 8.8 (p.141). The depiction of entropies in terms of Venn diagrams is misleading for at least two reasons.

First, one is used to thinking of Venn diagrams as depicting sets; but what are the 'sets' $H(X)$ and $H(Y)$ depicted in figure 8.2, and what are the objects that are members of those sets? I think this diagram encourages the novice student to make inappropriate analogies. For example, some students imagine that the random outcome (x, y) might correspond to a point in the diagram, and thus confuse entropies with probabilities.

[Cover,Thomas'06] & [MacKay'02] on three-term entropies

unconditional mutual inf. $I(X;Y)$ $I(X;Y|Z)$ conditional mutual inf.

Figure 8.3. A misleading representation of entropies, continued.



2.25 Venn diagrams. There isn't really a notion of mutual information common to three random variables. Here is one attempt at a definition: Using Venn diagrams, we can see that the mutual information common to three random variables X , Y , and Z can be defined by

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

This quantity is symmetric in X , Y , and Z , despite the preceding asymmetric definition. Unfortunately, $I(X; Y; Z)$ is not necessarily nonnegative. Find X , Y , and Z such that $I(X; Y; Z) < 0$, and prove the following two identities:

- (a) $I(X; Y; Z) = H(X, Y, Z) - H(X) - H(Y) - H(Z) + I(X; Y) + I(Y; Z) + I(Z; X).$
- (b) $I(X; Y; Z) = H(X, Y, Z) - H(X, Y) - H(Y, Z) - H(Z, X) + H(X) + H(Y) + H(Z).$

The first identity can be understood using the Venn diagram analogy for entropy and mutual information. The second identity follows easily from the first.

The conditional mutual information between X and Y given Z is the average over z of the above conditional mutual information.

$$I(X; Y | Z) = H(X | Z) - H(X | Y, Z). \quad (8.10)$$

No other 'three-term entropies' will be defined. For example, expressions such as $I(X; Y; Z)$ and $I(X | Y; Z)$ are illegal. But you may put conjunctions of arbitrary numbers of variables in each of the three spots in the expression $I(X; Y | Z)$ – for example, $I(A, B; C, D | E, F)$ is fine: it measures how much information on average c and d convey about a and b , assuming e and f are known.

[Yeung'08] disagrees and heavily uses "information diagrams"

3.5 Information Diagrams

We have established in Section 3.3 a one-to-one correspondence between Shannon's information measures and set theory. Therefore, it is valid to use an information diagram, which is a variation of a Venn diagram, to represent the relationship between Shannon's information measures.

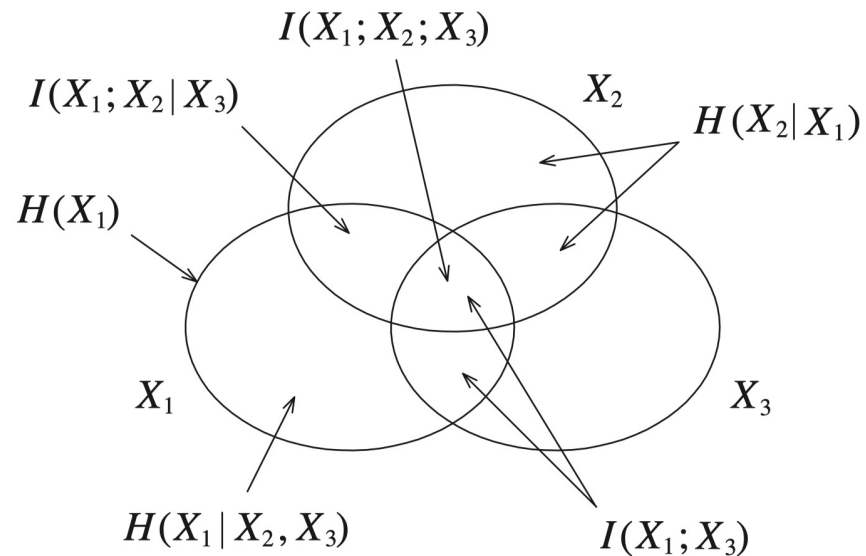


Fig. 3.4. The generic information diagram for X_1 , X_2 , and X_3 .

Interaction information

$J(X; Y; Z)$

- measures the influence of a variable Z on the amount of information shared between X and Y .* (And it is symmetric)
- It is positive when Z inhibits (i.e., accounts for or explains some of) the correlation between X and Y (that happens here).
- It is negative when Z facilitates or enhances the correlation (e.g., when X and Y are independent but not conditionally independent given Z , that's our last example).

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

* Alternative notations include $J(X; Y; Z)$ and $R(X; Y; Z)$. We don't recommend calling it "mutual information" and thus also replace the more common notation $I(X; Y; Z)$ with $J(X; Y; Z)$.

For more details, see https://en.wikipedia.org/wiki/Interaction_information.

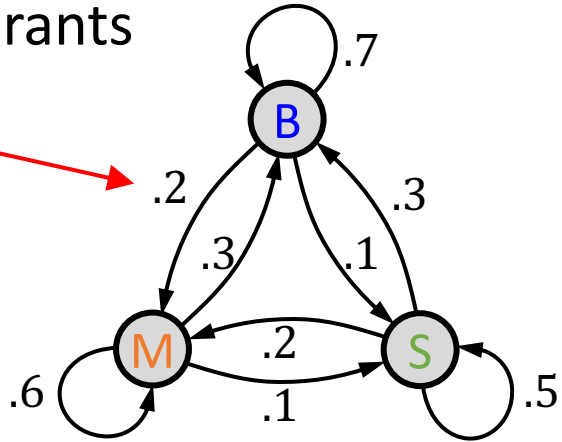
Gatterbauer, Aslam. Foundations and Applications of Information Theory: <https://northeastern-datalab.github.io/cs7840/>

Entropy rates of Markov Chains

Markov Chain

EXAMPLE: restaurants

$$\mathbb{P}[M|B] = 0.2$$



How to find the **stationary distribution** μ ?



State transition matrix \mathbf{P} :

$$\mathbf{P} = \begin{array}{c|ccc|c} & \text{B} & \text{M} & \text{S} & \Sigma \\ \hline \text{B} & .7 & .2 & .1 & 1 \\ \text{M} & .3 & .6 & .1 & 1 \\ \text{S} & .3 & .2 & .5 & 1 \\ \hline \Sigma & 1.3 & 1.0 & .7 & \end{array} \text{ row-stochastic}$$

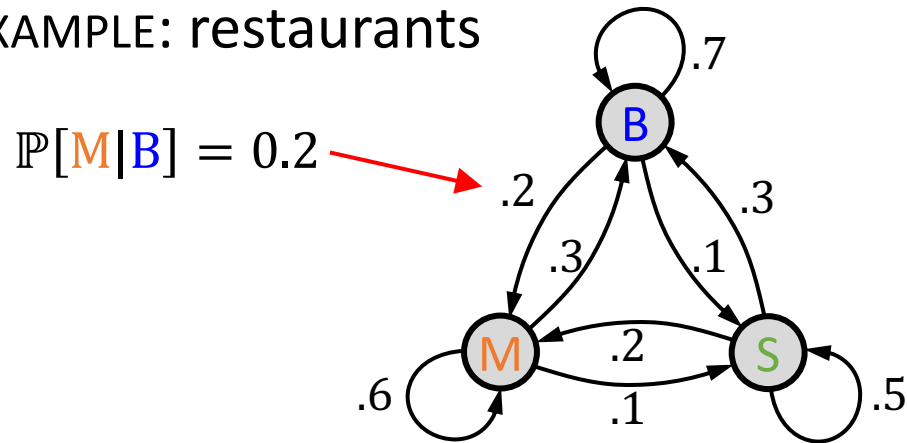
$$P_{i,j} = \mathbb{P}[X_{n+1} = j | X_n = i]:$$

probability of choosing j after i

$\mathbf{P}_{i,:}$: row vector (probability distribution)

Markov Chain

EXAMPLE: restaurants



How to find the stationary distribution μ ?
By finding the largest eigenvector of \mathbf{P} ,
i.e. solving an equation system:

$$\mu = \mathbf{P}^T \mu \quad \mu_j = \sum_i \mu_i P_{i,j} \text{ for all } j$$

transpose

$$\mu = \begin{pmatrix} 1/2 \\ 1/3 \\ 1/6 \end{pmatrix} \quad H(\mu) = 1.460$$

State transition matrix \mathbf{P} :

$$\mathbf{P} = \begin{array}{c} \mathbf{B} \\ \mathbf{M} \\ \mathbf{S} \\ \Sigma \end{array} \begin{array}{ccc} \mathbf{B} & \mathbf{M} & \mathbf{S} \\ \left(\begin{array}{ccc} .7 & .2 & .1 \\ .3 & .6 & .1 \\ .3 & .2 & .5 \end{array} \right) & & \\ \Sigma & 1.3 & 1.0 & .7 \end{array} \begin{array}{l} 1 \\ 1 \\ 1 \end{array} \text{ row-stochastic}$$

$$P_{i,j} = \mathbb{P}[X_{n+1} = j | X_n = i]:$$

probability of choosing j after i

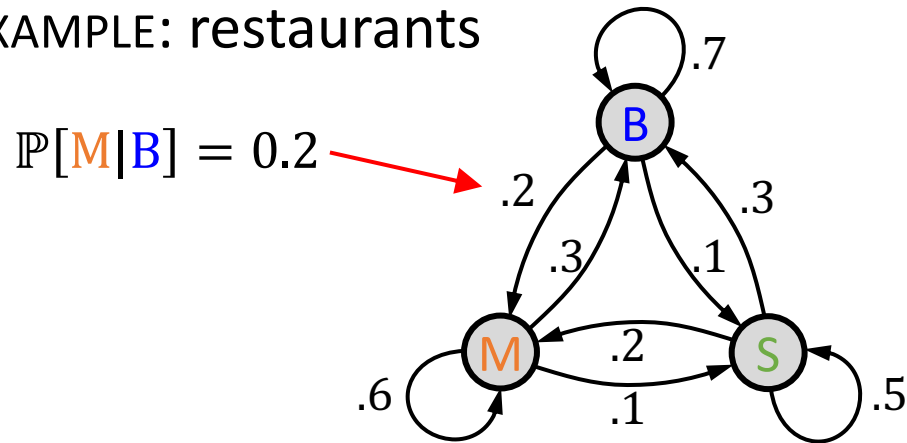
$\mathbf{P}_{i,:}$: row vector (probability distribution)

What would be the state transition matrix \mathbf{P}' with
same stationary distribution μ if there was no
memory: $\mathbb{P}[X_{n+1} = j | X_n = i] = \mathbb{P}[X_{n+1} = j]$



Markov Chain

EXAMPLE: restaurants



How to find the stationary distribution μ ?
By finding the largest eigenvector of \mathbf{P} ,
i.e. solving an equation system:

$$\mu = \mathbf{P}^T \mu \quad \mu_j = \sum_i \mu_i P_{i,j} \text{ for all } j$$

transpose (with arrow pointing to \mathbf{P}^T)

$$\mu = \begin{pmatrix} 1/2 \\ 1/3 \\ 1/6 \end{pmatrix} \quad H(\mu) = 1.460$$

State transition matrix \mathbf{P} :

$$\mathbf{P} = \begin{matrix} & \begin{matrix} \text{B} & \text{M} & \text{S} & \Sigma \end{matrix} \\ \begin{matrix} \text{B} \\ \text{M} \\ \text{S} \\ \Sigma \end{matrix} & \begin{pmatrix} .7 & .2 & .1 \\ .3 & .6 & .1 \\ .3 & .2 & .5 \\ 1.3 & 1.0 & .7 \end{pmatrix} \end{matrix} \quad \text{row-stochastic}$$

$P_{i,j} = \mathbb{P}[X_{n+1} = j | X_n = i]$:
probability of choosing j after i

$\mathbf{P}_{i,:}$: row vector (probability distribution)

What would be the state transition matrix \mathbf{P}' with
same stationary distribution μ if there was no
memory: $\mathbb{P}[X_{n+1} = j | X_n = i] = \mathbb{P}[X_{n+1} = j]$

$$\mathbf{P}' = \begin{matrix} & \begin{matrix} \text{B} & \text{M} & \text{S} & \Sigma \end{matrix} \\ \begin{matrix} \text{B} \\ \text{M} \\ \text{S} \\ \Sigma \end{matrix} & \begin{pmatrix} 1/2 & 1/3 & 1/6 \\ 1/2 & 1/3 & 1/6 \\ 1/2 & 1/3 & 1/6 \\ 1.5 & 1.0 & .5 \end{pmatrix} \end{matrix}$$

$$P_{i,j}' = \mu_j$$

Markov Chains and information measures

$X \rightarrow Y \rightarrow Z$ is a Markov chain if ?

$$p(x, y, z) =$$

Markov Chains and information measures

$X \rightarrow Y \rightarrow Z$ is a Markov chain if $X \perp Z | Y$, and thus: *The future depends only on the current state (not the previous ones)*

$$p(x, y, z) = ?$$

Markov Chains and information measures

$X \rightarrow Y \rightarrow Z$ is a **Markov chain** if $X \perp Z | Y$, and thus: *The future depends only on the current state (not the previous ones)*

$$p(x, y, z) = p(x) \cdot p(y|x) \cdot p(z|y)$$

In general, $p(x, y, z) = p(x) \cdot p(y|x) \cdot p(z|x, y)$

$$p(x, z|y) = ?$$

Markov Chains and information measures

$X \rightarrow Y \rightarrow Z$ is a **Markov chain** if $X \perp Z|Y$, and thus: *The future depends only on the current state (not the previous ones)*

$$p(x, y, z) = p(x) \cdot p(y|x) \cdot p(z|y)$$

$$\text{In general, } p(x, y, z) = p(x) \cdot p(y|x) \cdot p(z|x, y)$$

$$p(x, z|y) = p(x|y) \cdot p(z|y)$$

$$\text{In general, } p(x, z|y) = p(x|y) \cdot p(z|x, y)$$

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

Markov Chains and information measures

$X \rightarrow Y \rightarrow Z$ is a **Markov chain** if $X \perp Z|Y$, and thus: *The future depends only on the current state (not the previous ones)*

$$p(x, y, z) = p(x) \cdot p(y|x) \cdot p(z|y)$$

$$\text{In general, } p(x, y, z) = p(x) \cdot p(y|x) \cdot p(z|x, y)$$

$$p(x, z|y) = p(x|y) \cdot p(z|y)$$

$$\text{In general, } p(x, z|y) = p(x|y) \cdot p(z|x, y)$$

$$I(X; Z|Y) = 0$$

What does this mean for the interaction information $J(X; Y; Z)$?

?

Markov Chains and information measures

$X \rightarrow Y \rightarrow Z$ is a **Markov chain** if $X \perp Z | Y$, and thus: *The future depends only on the current state (not the previous ones)*

$$p(x, y, z) = p(x) \cdot p(y|x) \cdot p(z|y)$$

$$\text{In general, } p(x, y, z) = p(x) \cdot p(y|x) \cdot p(z|x, y)$$

$$p(x, z|y) = p(x|y) \cdot p(z|y)$$

$$\text{In general, } p(x, z|y) = p(x|y) \cdot p(z|x, y)$$

$$I(X; Z|Y) = 0$$

What does this mean for the interaction information $J(X; Y; Z)$?

$$J(X; Z; Y) = I(X; Z) - \underbrace{I(X; Z|Y)}_{= 0} = I(X; Z) \geq 0$$

- Recall: $J(X; Z; Y)$ measures the influence of a variable Y on the amount of information shared between X and Z .
- It is **positive when Y inhibits** (i.e., accounts for or explains some of) the correlation between X and Z (that happens here).
- It is **negative** when Y facilitates or enhances the correlation (e.g., when X and Y are independent yet not conditionally independent given Z , see earlier example).

Markov Chains and stationary stochastic processes

$X \rightarrow Y \rightarrow Z$ is a **Markov chain** if $X \perp Z | Y$, and thus:

$$p(x, y, z) = p(x) \cdot p(y|x) \cdot p(z|y)$$

$$\text{In general, } p(x, y, z) = p(x) \cdot p(y|x) \cdot p(z|x, y)$$

$$p(x, z|y) = p(x|y) \cdot p(z|y)$$

$$\text{In general, } p(x, z|y) = p(x|y) \cdot p(z|x, y)$$

$$I(X; Z|Y) = 0$$

A discrete stochastic process (X_1, X_2, \dots) is a **Markov chain** if

?

Markov Chains and stationary stochastic processes

$X \rightarrow Y \rightarrow Z$ is a **Markov chain** if $X \perp Z|Y$, and thus:

$$p(x, y, z) = p(x) \cdot p(y|x) \cdot p(z|y)$$

$$\text{In general, } p(x, y, z) = p(x) \cdot p(y|x) \cdot p(z|x, y)$$

$$p(x, z|y) = p(x|y) \cdot p(z|y)$$

$$\text{In general, } p(x, z|y) = p(x|y) \cdot p(z|x, y)$$

$$I(X; Z|Y) = 0$$

A discrete stochastic process (X_1, X_2, \dots) is a **Markov chain** if each RV depends only on the one preceding it and is conditionally independent of all the other preceding RVs

$$\mathbb{P}[x_{n+1}|x_n, x_{n-1}, \dots, x_1] = \mathbb{P}[x_{n+1}|x_n,]$$

A stochastic process $\{X_i\} = (X_1, X_2, \dots)$ is **stationary** if ...



Markov Chains and stationary stochastic processes

$X \rightarrow Y \rightarrow Z$ is a **Markov chain** if $X \perp Z|Y$, and thus:

$$p(x, y, z) = p(x) \cdot p(y|x) \cdot p(z|y)$$

$$\text{In general, } p(x, y, z) = p(x) \cdot p(y|x) \cdot p(z|x, y)$$

$$p(x, z|y) = p(x|y) \cdot p(z|y)$$

$$\text{In general, } p(x, z|y) = p(x|y) \cdot p(z|x, y)$$

$$I(X; Z|Y) = 0$$

A discrete stochastic process (X_1, X_2, \dots) is a **Markov chain** if each RV depends only on the one preceding it and is conditionally independent of all the other preceding RVs

$$\mathbb{P}[x_{n+1}|x_n, x_{n-1}, \dots, x_1] = \mathbb{P}[x_{n+1}|x_n,]$$

A stochastic process $\{X_i\} = (X_1, X_2, \dots)$ is **stationary** if the joint distribution of any subsequence is invariant w.r.t. shifts in the time index

$$\mathbb{P}[(x_1, x_2, \dots, x_k)] = \mathbb{P}[(x_{1+\ell}, x_{2+\ell}, \dots, x_{k+\ell})]$$

Entropy rate for stationary Markov Chain

The **entropy rate** of a stochastic process $\{X_i\}$ is the average entropy per symbol:

$$H(\mathcal{X}) = \lim_{n \rightarrow \infty} \frac{1}{n} \cdot H(X_1, X_2, \dots, X_n) \quad H(X_1, X_2, \dots, X_n) \rightarrow n \cdot H(\mathcal{X})$$

For a **stationary** stochastic process, this is equal to the rate of information innovation

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

For **stationary Markov Chain**, the entropy rate is

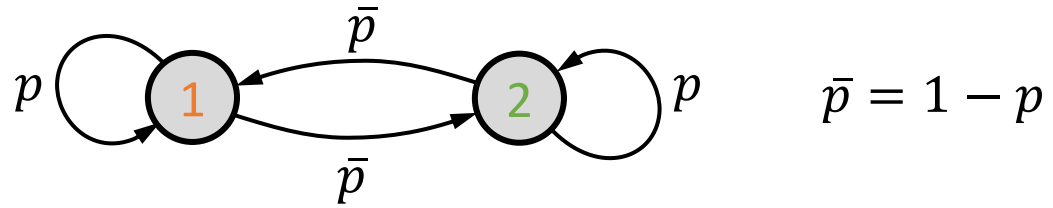
$$H(\mathcal{X}) = H(X_2 | X_1) \quad \text{where the conditional entropy is calculated using the } \underline{\text{stationary distribution}} \text{ (!)}$$

$$= \sum_i \mu_i \cdot H(X_2 | X_1 = i) = \sum_i \mu_i \cdot H(\mathbf{P}_{i:\cdot}) = \mathbb{E}_{i \sim \mu} [H(\mathbf{P}_{i:\cdot})]$$

$$= - \sum_i \mu_i \cdot \sum_j P_{ij} \cdot \lg(P_{ij})$$

Markov Chain

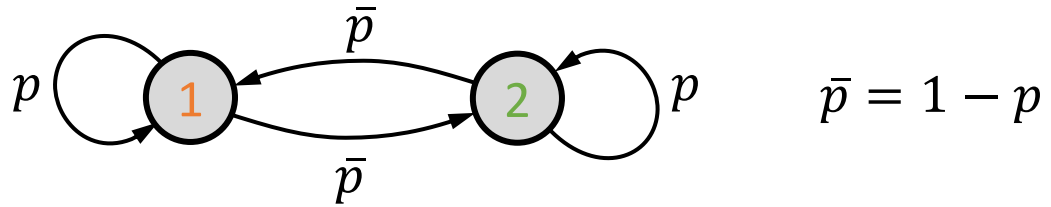
EXAMPLE: A simple two-state Markov Chain



P = ?

Markov Chain

EXAMPLE: A simple two-state Markov Chain

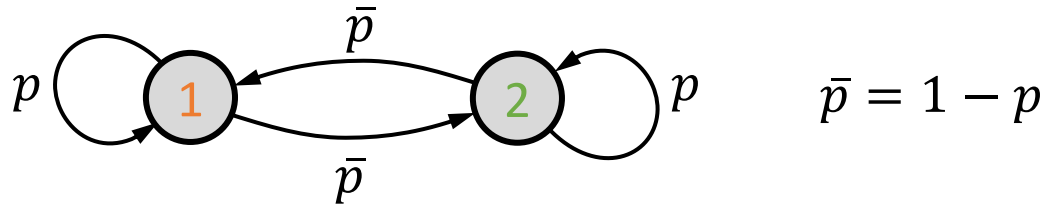


$$\mathbf{P} = \begin{pmatrix} p & \bar{p} \\ \bar{p} & p \end{pmatrix}$$

$$\boldsymbol{\mu} = ?$$

Markov Chain

EXAMPLE: A simple two-state Markov Chain



$$\mathbf{P} = \begin{pmatrix} p & \bar{p} \\ \bar{p} & p \end{pmatrix}$$

$$\boldsymbol{\mu} = \begin{pmatrix} 0.5 \\ 0.5 \end{pmatrix}$$

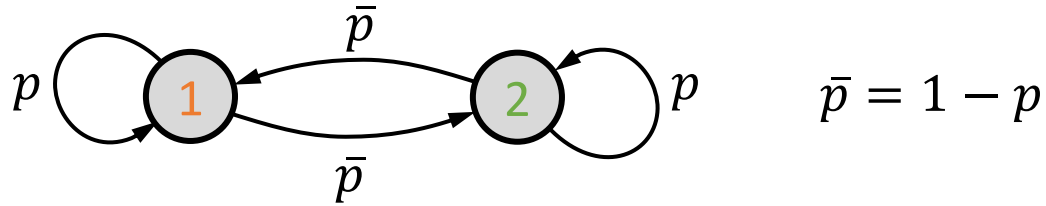
$$H(\boldsymbol{\mu}) = ?$$

$p = 0.95$: A random sample:
 AAAAAAAAABBBBBBBAAAAAAAABBBAAAAAAAAAAAA
 AAAAAAAAAAAAAAABBBBBBBBBBAAAAAAAAAAAAABBBBBBABBVBVVBBBAAAABAAAAABBBB
 BBBBVBVVBBBBBBBBBBBBBBBBBBBBBAA
 AAA
 AAAAAABAAABBBBBBAAAAAAAAAAAAAAAAABBBBA
 ABBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBAAAAAAAAAAAAAAAAAAAAAAAAAAAAAABBB
 BBBBVBVVBBBBBBBBBBBBBBBBBBBBBAAAAAAAAAAAAAAAABBBBABBVBVVBBBBBBBBBBBB
 BBBBVBVVBBBAAAAAAA
 AAAAAABBAAAAAAA
 AABBBBBBBBBBBBBBBBBBBBBBBBBBBBBB
 BBBBAA
 AAAAAAAAAAAAAABBBBBBBBBBBBBBBBBBAAAAAABAAAAAAAAAAAAABBBBVBVVBBBBBBBBBB
 BBBAAAAAAAAAAAAABBBBBBAA
 AAAAAAAAAAAAAABBBBBBBBBBBBBBBBBBAAAAAAAAAAAAAAAAAAAAABBBBVBVVBBBBBBBBBB
 BBBAAAAAAAAAAAAABBBBBBAA
 AAAAAAAAAAAAAABBBBVBVVBBBBBBBBBBBBBAAAAAAAAAAAAAAAAAAAAABBBBVBVVBBBBBB
 AAAAAAAAAAAAAABBBBVBVVBBBBBBBBBBBBBAAAAAAAAAAAAAAAAAAAAABBBBVBVVBBBBBB
 AAAAAAAAAAAAAABBBBVBVVBBBBBBBBBBBBBAAAAAAAAAAAAAAAAAAAAABBBBVBVVBBBBBB

$p = 0.05$: A random sample:
 AB
 AAAB
 AB
 BAB
 ABABBAB
 AB
 BAB
 AB
 AAB
 BBAB
 ABABABAAB
 ABAABABABBAAB
 BAB
 AB
 BABABABBAB
 AB
 AB

Markov Chain

EXAMPLE: A simple two-state Markov Chain



$$\mathbf{P} = \begin{pmatrix} p & \bar{p} \\ \bar{p} & p \end{pmatrix}$$

$$\boldsymbol{\mu} = \begin{pmatrix} 0.5 \\ 0.5 \end{pmatrix}$$

$$H(\boldsymbol{\mu}) = 1$$

$$H(\mathbf{P}) = ?$$

$p = 0.95$:

A random sample:

```

AAAAAAAAABBBBBBAAAAAAAAABBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBAAAAAAAAAAAA
AAAAAAAAAAAAABBBBBBBBBBAAAAAAAAAAAAABBBBBBABBVBVVVBAABAAAAABBBBB
BBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBB
AAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAA
AAAAAABAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAABBBBBBAAAAAAAAAAAAAAAAABBBBA
AABBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBAAAAAAAAAAAAAAAAAAAAAAAAAABB
BBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBAAAAAAAAAAAAAAAAABBBBBABBVBVVVBB
BBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBB
AAAAAABBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBB
AAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAABBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBB
BBBBBAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAA
AAAAAAAAAAAAABBBBBBBBBBBBBBBBBBBBAAABAAAAAAAAABBBBBBBBBBBBBBBBBBB
BBBAAAAAAAAABBBBBBAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAA
AAAAAAAAAAAAABBBBBBBBBBBBBBBBBBBBAAAAAAAAAAAAAAAAABBBBBBBBBBBBBBBBB
BBBAAAAAAAAABBBBBBAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAA
AAAAAAAAAAAAABBBBBBBBBBBBBBBBBBBBAAAAAAAAAAAAAAAAABBBBBBBBBBBBBBBBB
AAAAAAAAAAAAABBBBBBBBBBBBBBBBBBBBAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAA
AAAAAAAAAAAAABBBBBBBBBBBBBBBBBBBBAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAA
AAAAAAAAAAAAABBBBBBBBBBBBBBBBBBBBAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAA
    
```

$p = 0.05$:

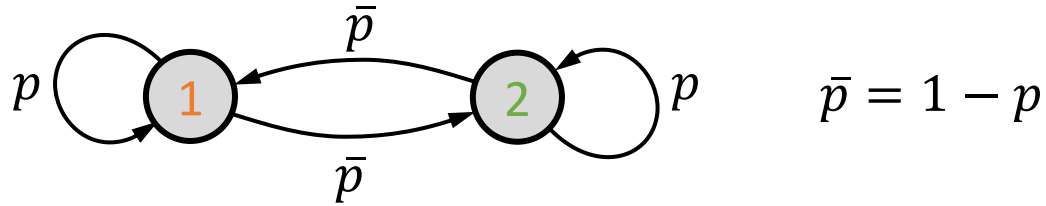
A random sample:

```

ABABABABABABABABABABABABABABABABABABABABABABABABABABABABABABABABAB
AAABABABABABABABABABABABABABABABABABABABABABABABABABABABABABABABAB
ABABABABABABABABABABABABABABABABABABABABABABABABABABABABABABABABAB
BABABABABABABABABABABABABABABABABABABABABABABABABABABABABABABABABAB
ABBBABABABABABABABABABABABABABABABABABABABABABABABABABABABABABABAB
ABABABABABABABABABABABABABABABABABABABABABABABABABABABABABABABABAB
BABABABABABABABABABABABABABABABABABABABABABABABABABABABABABABABABAB
ABABABABABABABABABABABABABABABABABABABABABABABABABABABABABABABABAB
AABABABABABABABABABABABABABABABABABABABABABABABABABABABABABABABAB
BBABABABABABABABABABABABABABABABABABABABABABABABABABABABABABABABAB
ABABABAABABABABABABABABABABABABABABABABABABABABABABABABABABABABAB
ABAABABABBAABABABABABABABABABABABABABABABABABABABABABABABABABABAB
BABABABABABABABABABABABABABABABABABABABABABABABABABABABABABABABAB
ABABABABABABABABABABABABABABABABABABABABABABABABABABABABABABABAB
BABABABABABABABABABABABABABABABABABABABABABABABABABABABABABABABAB
ABABABABABABABABABABABABABABABABABABABABABABABABABABABABABABABAB
ABABABABABABABABABABABABABABABABABABABABABABABABABABABABABABABAB
    
```

Markov Chain

EXAMPLE: A simple two-state Markov Chain

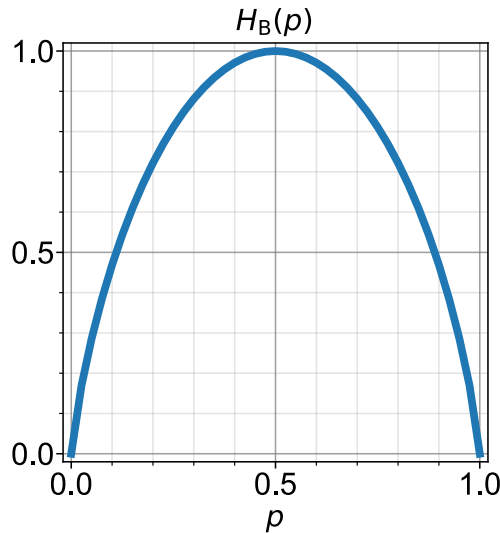


$$\mathbf{P} = \begin{pmatrix} p & \bar{p} \\ \bar{p} & p \end{pmatrix}$$

$$\boldsymbol{\mu} = \begin{pmatrix} 0.5 \\ 0.5 \end{pmatrix}$$

$$H(\boldsymbol{\mu}) = 1$$

$$H(\mathbf{P}) = \mathbb{E}_{i \sim \boldsymbol{\mu}}[H(\mathbf{P}_i)] = H_B(p)$$



$$H_B(p) = 0.286$$

$p = 0.95$:

A random sample:

```

AAAAAAAAABBBBBBAAAAAAAAABBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBAAAAAAAAAAAA
AAAAAAAAAAAAAAAAABBBBBBBBBBAAAAAAAAAAAAAAAAABBBBBBABBVBVVVBAABAAAAABBBBB
BBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBB
AAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAA
AAAAAAAAABAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAABBBBBBBBAAAAAAAAAAAAAAAAABBBBA
AABBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBAAAAAAAAAAAAAAAAAAAAAAAAAABB
BBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBAAAAAAAAAAAAAABBBBBBABBVBVVVBBBB
BBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBB
AAAAAAAAABBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBAAAA
AAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAABBBBBBBBBBBBBBBBBBBBBBBBBBBBBBB
BBBBBAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAA
AAAAAAAAAAAAABBBBBBBBBBBBBBBBBBBBAAAABBAAAAAAAAAAABBBBBBBBBBBBBBBBBBB
BBBAAAAAAAAAABBBBBBBBAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAA
AAAAAAAAAAAAABBBBBBBBBBBBBBBBBBBBAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAA
AAAAAAAAAAAAABBBBBBBBBBBBBBBBBBBBAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAA
AAAAAAAAAAAAABBBBBBBBBBBBBBBBBBBBAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAA
AAAAAAAAAAAAABBBBBBBBBBBBBBBBBBBBAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAAA

```

$p = 0.05$:

A random sample:

```

ABABABABABABABABABABABABABABABABABABABABABABAABABAABABABABABABABABAB
AAABABABABABABABABABABABABABABABABABABABABABAABABABABABABABABABAB
ABABABABABABABABABABABABABABABAAAABABABABABABABBABABABABAABABABA
BABABABABABABABAABABABABABABABABABABABABABABBABABABABABABABABAB
ABABBABABABAABAABABABABABABABABABABABABABABBAABABABABABABABABAB
ABABABABABBABABABABABABABABABABABABABABABABABABABBABABABABABABABA
BABABABABABABABABABBABABAABABABAABABBABABABABABABABABABABABABAB
ABABABABABABABABABABABABABABABABABABBABABABABABABABABABABABABAB
AABABABABABBABABABABABABABABABABABABABABABAABABABBABABABABABABA
BBABABABABABABABABABABABABABABABABABABABABABABABABABABABABABABB
ABABABAABABABABABABABABABABABABABABABABABBABABABABABABABABBAB
ABAABABABBAABABABABABABABABABABABABABABABABABABABABABABABABABBABA
BABABABABABABABABAABABABABABABABAABABABABABABABAABBBABAB
ABABABABABABABAABABABABABABABABABABABABABABABABAABABABABABA
BABABABBABABABABABABABABABABABABABABBABABABABABABABABABABABAB
ABABABABABAABABBABABABABABABAABABABABABABABABABABABABABABABAB

```

Data Processing Inequality

Data Processing Inequality for $X \rightarrow Y \rightarrow Z$

Intuitively, the **data processing inequality** states that no clever transformation of a received representation Y can increase the information about the original information X .

THEOREM: Suppose we have a **Markov chain** $X \rightarrow Y \rightarrow Z$ (and thus $X \perp Z|Y$), then

$$I(X; Y) \stackrel{?}{\leq \text{ or } \geq} I(X; Z)$$

Data Processing Inequality for $X \rightarrow Y \rightarrow Z$

Intuitively, the **data processing inequality** states that no clever transformation of a received representation Y can increase the information about the original information X .

THEOREM: Suppose we have a **Markov chain** $X \rightarrow Y \rightarrow Z$ (and thus $X \perp Z|Y$), then

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

COROLLARY: If $Z = f(Y)$, then $I(X; Y) \geq I(X; f(Y))$. Thus functions of Y cannot increase the information about X . In other words, no processing of Y , deterministic or random, can increase the information that Y contains about X (unless you add additional outside information).

This follows from $X \rightarrow Y \rightarrow f(Y)$ forming a Markov chain.

Data Processing Inequality for $X \rightarrow Y \rightarrow Z$

Intuitively, the **data processing inequality** states that no clever transformation of a received representation Y can increase the information about the original information X .

THEOREM: Suppose we have a **Markov chain** $X \rightarrow Y \rightarrow Z$ (and thus $X \perp Z|Y$), then

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

PROOF:

$$\underbrace{I(X; Y, Z)}_{I(X; (Y, Z))} = ?$$

Data Processing Inequality for $X \rightarrow Y \rightarrow Z$

Intuitively, the **data processing inequality** states that no clever transformation of a received representation Y can increase the information about the original information X .

THEOREM: Suppose we have a **Markov chain** $X \rightarrow Y \rightarrow Z$ (and thus $X \perp Z|Y$), then

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

PROOF:

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

Data Processing Inequality for $X \rightarrow Y \rightarrow Z$

Intuitively, the **data processing inequality** states that no clever transformation of a received representation Y can increase the information about the original information X .

THEOREM: Suppose we have a **Markov chain** $X \rightarrow Y \rightarrow Z$ (and thus $X \perp Z|Y$), then

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

PROOF:

$$\begin{aligned} \underbrace{I(X; Y, Z)}_{I(X; (Y, Z))} &= H(X) - \underbrace{H(X|Y, Z)}_{I(X|(Y, Z))} = \underbrace{H(X) + (-H(X|Z) + H(X|Z))}_{I(X; Z)} - \underbrace{H(X|Y, Z)}_{I(X; Y|Z)} \\ &= I(X; Z) + I(X; Y|Z) \\ &= I(X; Y) + \underbrace{I(X; Z|Y)}_{?} \quad (\text{similarly, from symmetry}) \end{aligned}$$

Data Processing Inequality for $X \rightarrow Y \rightarrow Z$

Intuitively, the **data processing inequality** states that no clever transformation of a received representation Y can increase the information about the original information X .

THEOREM: Suppose we have a **Markov chain** $X \rightarrow Y \rightarrow Z$ (and thus $X \perp Z|Y$), then

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

PROOF:

$$\begin{aligned} \underbrace{I(X; Y, Z)}_{I(X; (Y, Z))} &= H(X) - \underbrace{H(X|Y, Z)}_{I(X|(Y, Z))} = \underbrace{H(X) + (-H(X|Z) + H(X|Z))}_{I(X; Z)} - \underbrace{H(X|Y, Z)}_{I(X; Y|Z)} \\ &= I(X; Z) + I(X; Y|Z) \\ &= I(X; Y) + \underbrace{I(X; Z|Y)}_{= 0} \quad (\text{similarly, from symmetry}) \end{aligned}$$

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

Data Processing Inequality for $X \rightarrow Y \rightarrow Z$

Intuitively, the **data processing inequality** states that no clever transformation of a received representation Y can increase the information about the original information X .

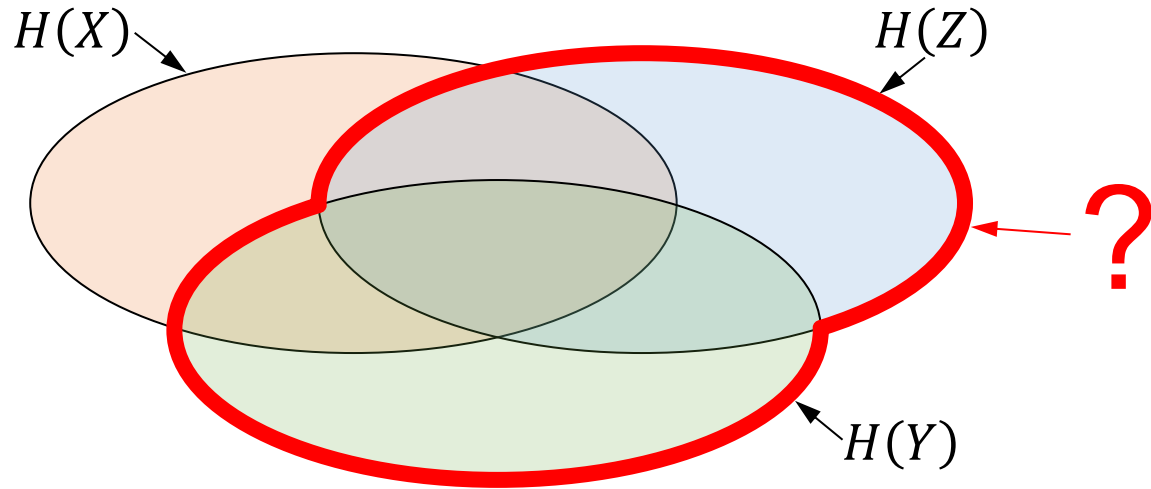
THEOREM: Suppose we have a **Markov chain** $X \rightarrow Y \rightarrow Z$ (and thus $X \perp Z|Y$), then

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

PROOF:

$$\begin{aligned} \underbrace{I(X; Y, Z)}_{I(X; (Y, Z))} &= H(X) - \underbrace{H(X|Y, Z)}_{I(X|(Y, Z))} = \underbrace{H(X) + (-H(X|Z) + H(X|Z))}_{I(X; Z)} - \underbrace{H(X|Y, Z)}_{I(X; Y|Z)} \\ &= I(X; Z) + I(X; Y|Z) \\ &= I(X; Y) + \underbrace{I(X; Z|Y)}_{= 0} \quad (\text{similarly, from symmetry}) \\ &\Rightarrow I(X; Y) = I(X; Z) + \underbrace{I(X; Y|Z)}_{\geq 0} \quad \text{since mutual information is always non-negative} \\ &\Rightarrow \boxed{I(X; Y) \geq I(X; Z)} \end{aligned}$$

Data Processing Inequality for $X \rightarrow Y \rightarrow Z$

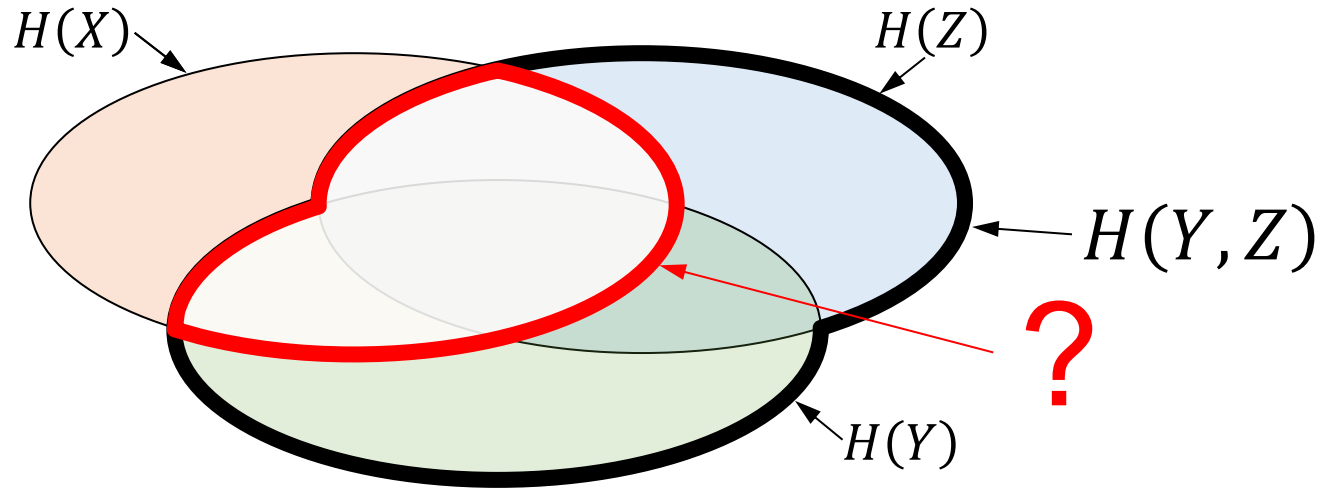


$$\begin{aligned}
 \underbrace{I(X; Y, Z)}_{I(X; (Y, Z))} &= H(X) - \underbrace{H(X|Y, Z)}_{I(X|Y, Z)} = \underbrace{H(X) + (-H(X|Z) + H(X|Z))}_{I(X; Z)} - \underbrace{H(X|Y, Z)}_{I(X; Y|Z)} \\
 &= I(X; Z) + I(X; Y|Z) \\
 &= I(X; Y) + \underbrace{I(X; Z|Y)}_{= 0} \quad (\text{similarly, from symmetry})
 \end{aligned}$$

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

$$\Rightarrow \boxed{I(X; Y) \geq I(X; Z)} \quad \underbrace{\geq 0}_{\text{since mutual information is always non-negative}}$$

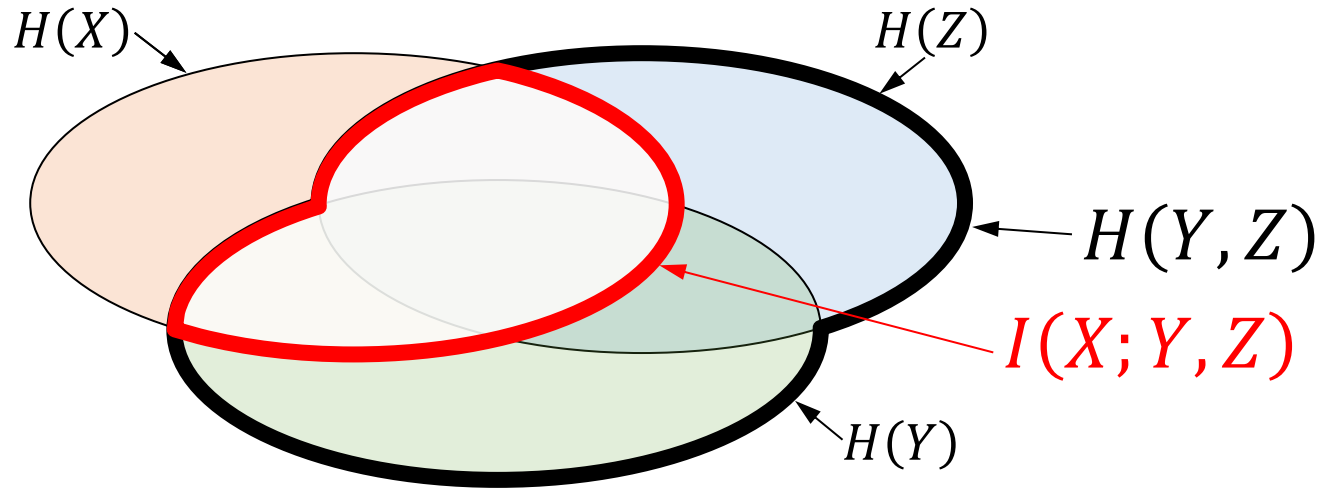
Data Processing Inequality for $X \rightarrow Y \rightarrow Z$



$$\begin{aligned}
 \underbrace{I(X; Y, Z)}_{I(X; (Y, Z))} &= H(X) - \underbrace{H(X|Y, Z)}_{I(X|(Y, Z))} = \underbrace{H(X) + (-H(X|Z) + H(X|Z))}_{I(X; Z)} - \underbrace{H(X|Y, Z)}_{I(X; Y|Z)} \\
 &= I(X; Z) + I(X; Y|Z) \\
 &= I(X; Y) + \underbrace{I(X; Z|Y)}_{= 0} \quad (\text{similarly, from symmetry})
 \end{aligned}$$

$$\begin{aligned}
 &\Rightarrow I(X; Y) = I(X; Z) + I(X; Y|Z) \\
 &\Rightarrow \boxed{I(X; Y) \geq I(X; Z)} \quad \underbrace{\geq 0}_{\text{since mutual information is always non-negative}}
 \end{aligned}$$

Data Processing Inequality for $X \rightarrow Y \rightarrow Z$



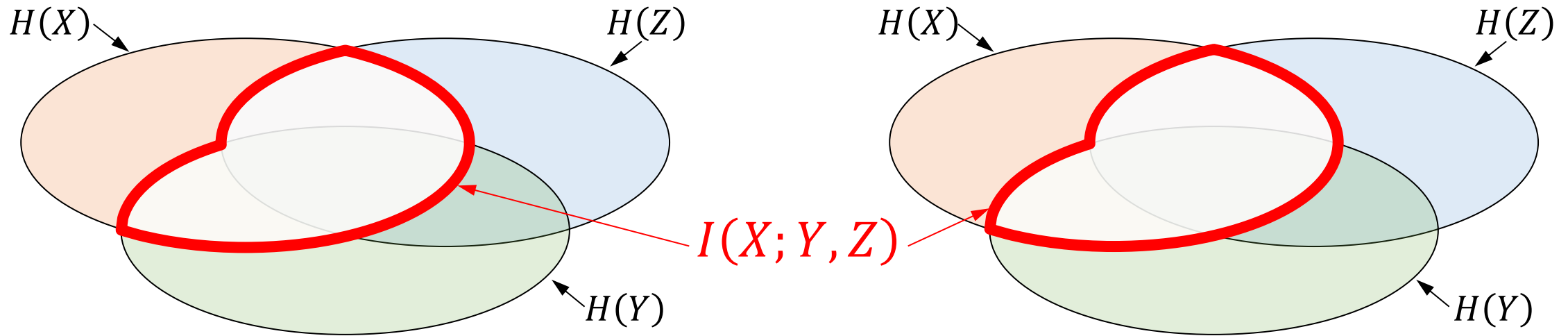
$$\begin{aligned}
 \underbrace{I(X; Y, Z)}_{I(X; (Y, Z))} &= H(X) - \underbrace{H(X|Y, Z)}_{I(X|(Y, Z))} = \underbrace{H(X) + (-H(X|Z) + H(X|Z))}_{I(X; Z)} - \underbrace{H(X|Y, Z)}_{I(X; Y|Z)} \\
 &= I(X; Y) + \underbrace{I(X; Z|Y)}_{= 0} \quad (\text{similarly, from symmetry})
 \end{aligned}$$

$$\Rightarrow I(X; Y) = I(X; Z) + \underbrace{I(X; Y|Z)}_{\geq 0}$$

$$\Rightarrow \boxed{I(X; Y) \geq I(X; Z)}$$

since mutual information is always non-negative

Data Processing Inequality for $X \rightarrow Y \rightarrow Z$



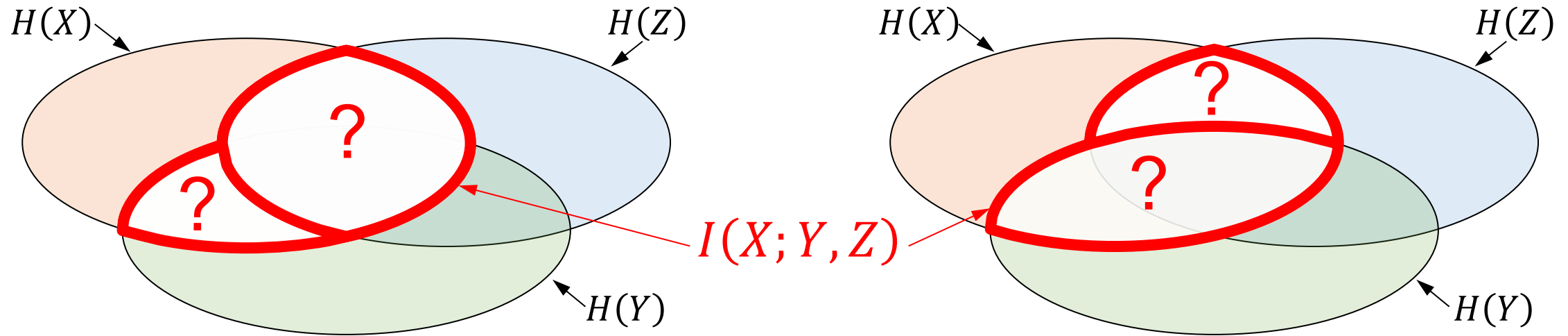
$$\begin{aligned}
 \underbrace{I(X; Y, Z)}_{I(X; (Y, Z))} &= H(X) - \underbrace{H(X|Y, Z)}_{I(X|(Y, Z))} = \underbrace{H(X) + (-H(X|Z) + H(X|Z))}_{I(X; Z)} - \underbrace{H(X|Y, Z)}_{I(X; Y|Z)} \\
 &= I(X; Y) + \underbrace{I(X; Z|Y)}_{= 0} \quad (\text{similarly, from symmetry})
 \end{aligned}$$

$$\Rightarrow I(X; Y) = I(X; Z) + \underbrace{I(X; Y|Z)}_{\geq 0}$$

$$\Rightarrow \boxed{I(X; Y) \geq I(X; Z)}$$

since mutual information is always non-negative

Data Processing Inequality for $X \rightarrow Y \rightarrow Z$

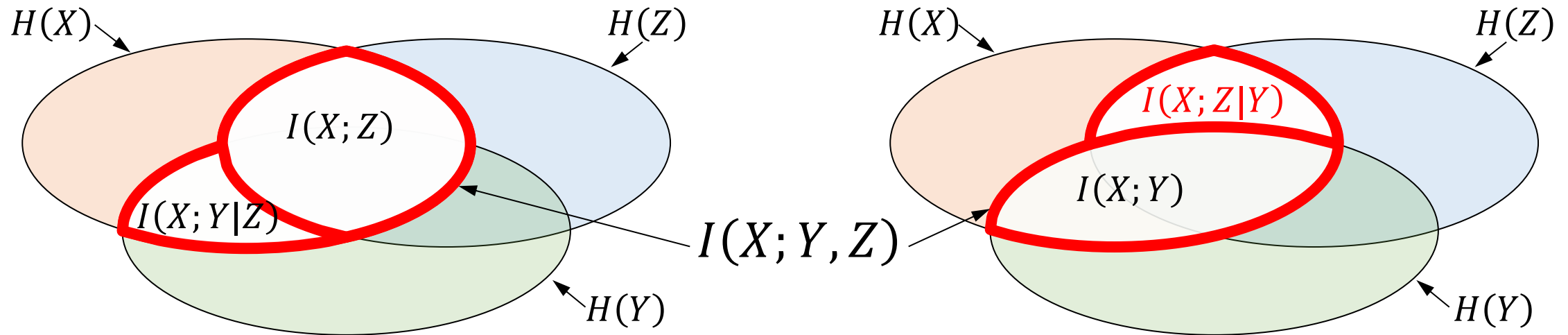


$$\begin{aligned}
 \underbrace{I(X; Y, Z)}_{I(X; (Y, Z))} &= H(X) - \underbrace{H(X|Y, Z)}_{I(X|(Y, Z))} = \underbrace{H(X) + (-H(X|Z) + H(X|Z))}_{I(X; Z)} - \underbrace{H(X|Y, Z)}_{I(X; Y|Z)} \\
 &= I(X; Y) + \underbrace{I(X; Z|Y)}_{= 0} \quad (\text{similarly, from symmetry})
 \end{aligned}$$

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

$$\Rightarrow \boxed{I(X; Y) \geq I(X; Z)} \quad \underbrace{\geq 0}_{\text{since mutual information is always non-negative}}$$

Data Processing Inequality for $X \rightarrow Y \rightarrow Z$



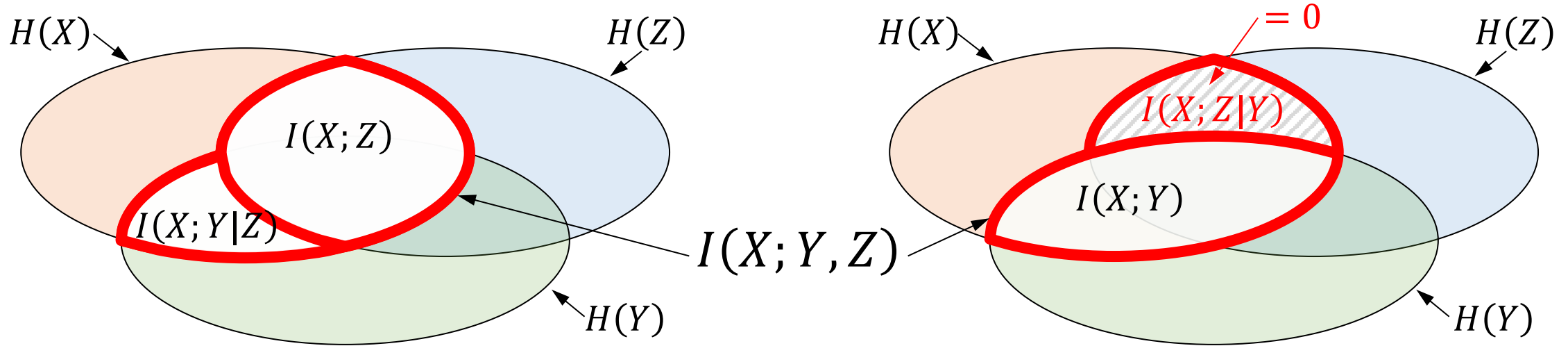
$$\begin{aligned}
 \underbrace{I(X; Y, Z)}_{I(X; (Y, Z))} &= H(X) - \underbrace{H(X|Y, Z)}_{I(X|(Y, Z))} = \underbrace{H(X) + (-H(X|Z) + H(X|Z))}_{I(X; Z)} - \underbrace{H(X|Y, Z)}_{I(X; Y|Z)} \\
 &= I(X; Y) + \underbrace{I(X; Z|Y)}_{= 0} \quad (\text{similarly, from symmetry})
 \end{aligned}$$

$$\Rightarrow I(X; Y) = I(X; Z) + \underbrace{I(X; Y|Z)}_{\geq 0}$$

$$\Rightarrow \boxed{I(X; Y) \geq I(X; Z)}$$

since mutual information is always non-negative

Data Processing Inequality for $X \rightarrow Y \rightarrow Z$



$$\begin{aligned}
 \underbrace{I(X; Y, Z)}_{I(X; (Y, Z))} &= H(X) - \underbrace{H(X|Y, Z)}_{I(X|(Y, Z))} = \underbrace{H(X) + (-H(X|Z) + H(X|Z))}_{I(X; Z)} - \underbrace{H(X|Y, Z)}_{I(X; Y|Z)} \\
 &= I(X; Y) + \underbrace{I(X; Z|Y)}_{= 0} \quad (\text{similarly, from symmetry})
 \end{aligned}$$

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

$$\Rightarrow \boxed{I(X; Y) \geq I(X; Z)}$$

$$\geq 0$$

since mutual information is always non-negative

Part 1: Theory

L10: Basics of entropy (6/6)

[data processing inequality, sufficient statistics, information inequalities]

Wolfgang Gatterbauer, Javed Aslam

cs7840 Foundations and Applications of Information Theory (fa24)

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

10/7/2024

Pre-class conversations

- Last class recapitulation
 - Please use our anonymous feedback to let us know which parts were too fast or unclear
 - Web page & readings
 - Today:
 - Sufficient statistics,
 - Information inequalities
 - (3 project ideas)
 - Next time:
 - We skip forward from part 1 to part 3: practical applications for a bit, before later coming back to more theory and the axiomatic approach
- **Lecture 9 (Wed 10/2): Basics of information theory (5/6)** [multivariate entropies, interaction information, Markov chains, data processing inequality]
 - **Lecture 10 (Mon 10/7): Basics of information theory (6/6)** [data processing inequality, sufficient statistics, information inequalities]
 - [Casella,Berger'24] **Statistical inference (2nd ed)**, CRC press, 2024: Ch 6 Principles of data reduction, Ch 6.2.1 Sufficient statistics.
 - [Fithian'24] **Statistics 210a: Theoretical Statistics (Lecture 4 sufficiency)**, Berkeley, 2014.
 - [Yeung'08] **Information Theory and Network Coding**. 2008: Ch 2.6 The basic inequalities, Ch 2.7 Some Useful Information Inequalities, Ch 3.5 Information Diagrams, Ch 13 Information inequalities, Ch 14 Shannon-type inequalities, Ch 15 Beyond Shannon-type inequalities

Knowledge Distillation from NN

Concept of distillation [\[edit \]](#)

Knowledge transfer from a large model to a small one somehow needs to teach the latter without loss of validity. If both models are trained on the same data, the smaller model may have insufficient capacity to learn a **concise knowledge representation** compared to the large model. However, some information about a concise knowledge representation is encoded in the **pseudolikelihoods** assigned to its output: when a model correctly predicts a class, it assigns a large value to the output variable corresponding to such class, and smaller values to the other output variables. The distribution of values among the outputs for a record provides information on how the large model represents knowledge. Therefore, the goal of economical deployment of a valid model can be achieved by training only the large model on the data, exploiting its better ability to learn concise knowledge representations, and then distilling such knowledge into the smaller model, by training it to learn the **soft output** of the large model.^[1]

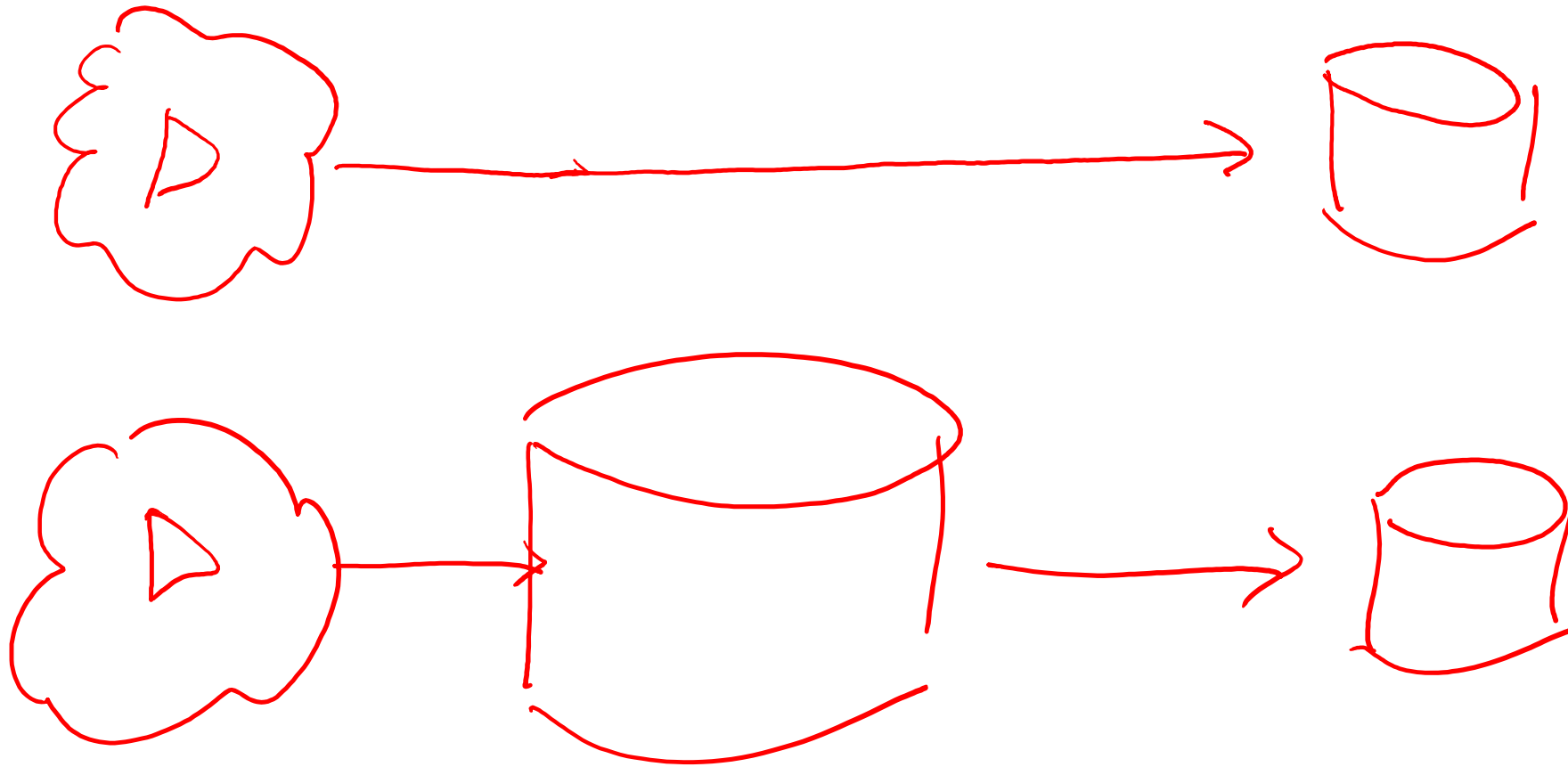
What is Knowledge Distillation?

Knowledge distillation is a powerful technique in machine learning that allows us to transfer knowledge from a large, complex model to a smaller, simpler one. By doing so, we can reduce the memory footprint and computational requirements of the model without significant performance loss.

The fundamental idea behind knowledge distillation is to leverage the soft probabilities or logits of a larger "teacher network" along with the available class labels to train a smaller "student network". These soft probabilities provide more information than just the class labels, enabling the student network to learn more effectively.

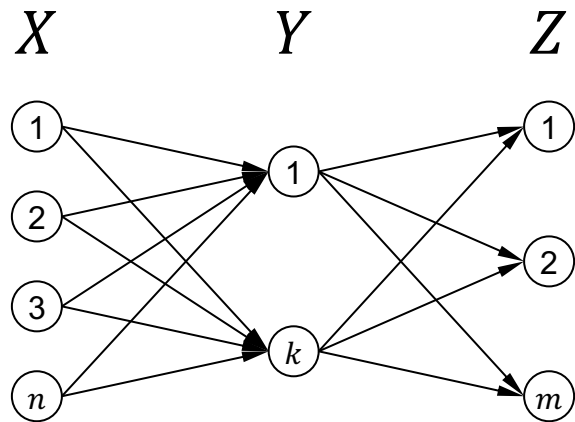
1. **Offline Distillation:** Imagine an aspiring author learning from an already published, successful book. The published book (the teacher model) is complete and fixed. The new writer (the student model) learns from this book, attempting to write their own based on the insights gained. *In the context of neural networks*, this is like using a fully trained, sophisticated neural network to train a simpler, more efficient network. The student network learns from the established knowledge of the teacher without modifying it.

Knowledge Distillation from NN



Bottleneck $X \rightarrow Y \rightarrow Z$

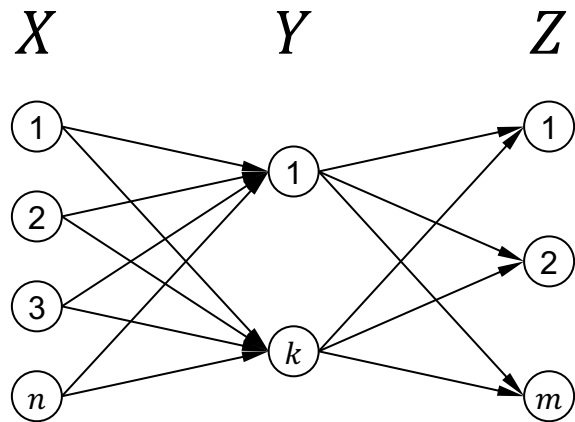
EXAMPLE: suppose a (non-stationary) Markov chain starts in one of n states, necks down to $k < n$ states, and then fans back to $m > k$ states. In other words, $X \rightarrow Y \rightarrow Z$ with $p(x, y, z) = p(x) \cdot p(y|x) \cdot p(z|y)$, and $x \in [n], y \in [k], z \in [m]$.



How can we upper bound $I(X; Z)$?

Bottleneck $X \rightarrow Y \rightarrow Z$

EXAMPLE: suppose a (non-stationary) Markov chain starts in one of n states, necks down to $k < n$ states, and then fans back to $m > k$ states. In other words, $X \rightarrow Y \rightarrow Z$ with $p(x, y, z) = p(x) \cdot p(y|x) \cdot p(z|y)$, and $x \in [n], y \in [k], z \in [m]$.



How can we upper bound $I(X; Z)$?

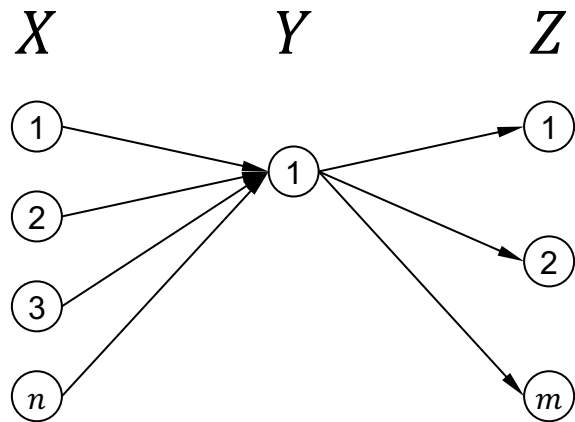
$$\begin{aligned} I(X; Z) &\leq I(X; Y) = H(Y) - H(Y|X) \\ &\leq H(Y) \\ &\leq \lg(k) \end{aligned}$$

\Rightarrow The dependence between X and Z is limited by the size k of the bottleneck.

What if $k = 1$?

Bottleneck $X \rightarrow Y \rightarrow Z$

EXAMPLE: suppose a (non-stationary) Markov chain starts in one of n states, necks down to $k < n$ states, and then fans back to $m > k$ states. In other words, $X \rightarrow Y \rightarrow Z$ with $p(x, y, z) = p(x) \cdot p(y|x) \cdot p(z|y)$, and $x \in [n], y \in [k], z \in [m]$.



How can we upper bound $I(X; Z)$?

$$\begin{aligned} I(X; Z) &\leq I(X; Y) = H(Y) - H(Y|X) \\ &\leq H(Y) \\ &\leq \lg(k) \end{aligned}$$

\Rightarrow The dependence between X and Z is limited by the size k of the bottleneck.

What if $k = 1$? $\Rightarrow I(X; Z) \leq \lg 1 = 0$. $\Rightarrow X$ and Z are **independent**.

Sufficient statistics

Following part builds on text, notation and examples from several sources, in particular:

[Casella,Berger'24] Statistical inference (2nd ed), 2024: Ch 6 Principles of Data Reduction. <https://doi.org/10.1201/9781003456285>

[Fithian'24] Statistics 210a: Theoretical Statistics, Berkeley, 2014: Lecture 4 sufficiency. <https://stat210a.berkeley.edu/fall-2024/reader/sufficiency.html>

[Scott'11] EECS 564: Estimation, Filtering, and Detection, University of Michigan, 2011: Lecture 5 Sufficient statistics. https://web.eecs.umich.edu/~cscott/past_courses/eecs564w11/index.html

[Cover,Thomas'06] Elements of Information Theory (2nd ed), 2006: Ch 2.9 Sufficient Statistics. <https://www.doi.org/10.1002/047174882X>

Parameter estimation

Suppose the probability distribution of a random variable X is determined by a parameter θ :

$$X \sim f_{\theta}(x) \quad \text{Think of this as a conditional distribution: } f_{\theta}(x) = p(x|\theta)$$

EXAMPLE: If X is a discrete Bernoulli RV, then its pmf (probability mass function) is parameterized by p :

$$f_p(x) = ?$$

Parameter estimation

Suppose the probability distribution of a random variable X is determined by a parameter θ :

$$X \sim f_{\theta}(x) \quad \text{Think of this as a conditional distribution: } f_{\theta}(x) = p(x|\theta)$$

EXAMPLE: If X is a discrete Bernoulli RV, then its pmf (probability mass function) is parameterized by p :

$$f_p(x) = \begin{cases} p & \text{if } x = 1 \\ \bar{p} & \text{if } x = 0 \end{cases} \quad \bar{p} := 1 - p$$



EXAMPLE: If X is a continuous Normal RV, then its pdf (probability density function) is parameterized by (μ, σ^2) :

$$f_{(\mu, \sigma^2)}(x) = ?$$

The parameter can also be a vector

Parameter estimation

Suppose the probability distribution of a random variable X is determined by a parameter θ :

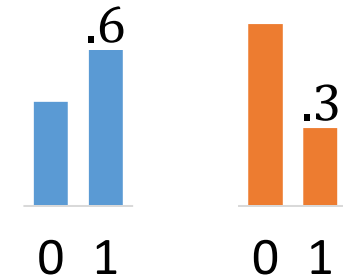
$$X \sim f_{\theta}(x) \quad \text{Think of this as a conditional distribution: } f_{\theta}(x) = p(x|\theta)$$

EXAMPLE: If X is a discrete Bernoulli RV, then its pmf (probability mass function) is parameterized by p :

$$f_p(x) = \begin{cases} p & \text{if } x = 1 \\ \bar{p} & \text{if } x = 0 \end{cases}$$



$\mathbf{x} = (1,1,0,1,1,1,0,0,1,1)$

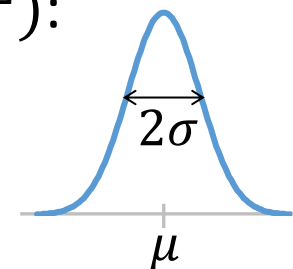


?

EXAMPLE: If X is a continuous Normal RV, then its pdf (probability density function) is parameterized by (μ, σ^2) :

$$f_{(\mu, \sigma^2)}(x) = \frac{1}{\sqrt{2\pi\sigma^2}} e^{-\frac{(x-\mu)^2}{2\sigma^2}}$$

The parameter can also be a vector



In **statistical inference**, we assume the functional form of f is known, but θ is hidden. We then observe a realization (a sample) \mathbf{x} of iid RV's \mathbf{X} and want to guess θ ("estimate θ ").

Independent and Identically Distributed

Parameter estimation

Suppose the probability distribution of a random variable X is determined by a parameter θ :

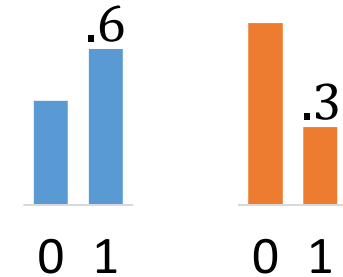
$$X \sim f_{\theta}(x) \quad \text{Think of this as a conditional distribution: } f_{\theta}(x) = p(x|\theta)$$

EXAMPLE: If X is a discrete Bernoulli RV, then its pmf (probability mass function) is parameterized by p :

$$f_p(x) = \begin{cases} p & \text{if } x = 1 \\ \bar{p} & \text{if } x = 0 \end{cases}$$



$\mathbf{x} = (1,1,0,1,1,1,0,0,1,1)$

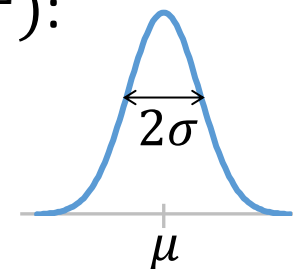


?

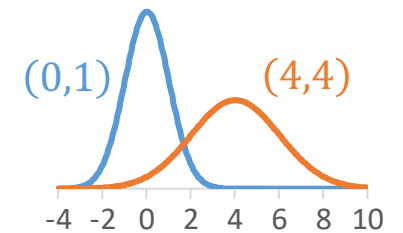
EXAMPLE: If X is a continuous Normal RV, then its pdf (probability density function) is parameterized by (μ, σ^2) :

$$f_{(\mu, \sigma^2)}(x) = \frac{1}{\sqrt{2\pi\sigma^2}} e^{-\frac{(x-\mu)^2}{2\sigma^2}}$$

The parameter can also be a vector



$\mathbf{x} = (5.2, 2.5, 0.3, 4.2)$



?

In **statistical inference**, we assume the functional form of f is known, but θ is hidden. We then observe a realization (a sample) \mathbf{x} of iid RV's \mathbf{X} and want to guess θ ("estimate θ ").

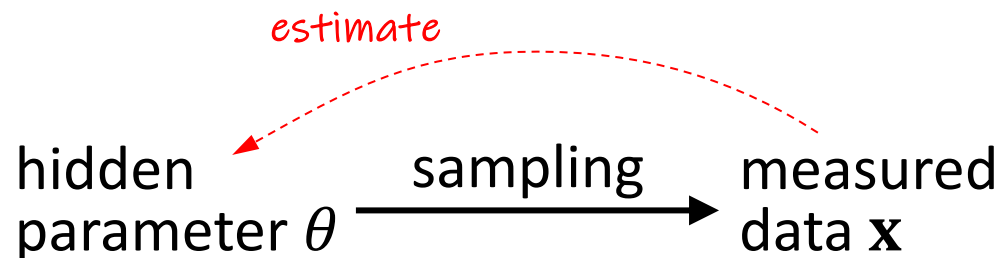
Independent and Identically Distributed

Sufficient statistics

If the sample $\mathbf{x} = (x_1, \dots, x_n)$ and unknown parameter θ , we would like to compress the measurements \mathbf{x} into a low-dimensional statistic without affecting the quality of the possible inference about θ (i.e. we do not want to lose relevant information about θ).

In other words, we are interested in whether there exists a **sufficient statistic** $T(\mathbf{X})$ where the dimension of $\mathbf{t} = T(\mathbf{x})$ is $m < n$, s.t. **\mathbf{t} carries all the useful information from \mathbf{x} about θ .**

If such a sufficient statistic exists, then for the purpose of studying θ , we could discard the raw measurement \mathbf{x} and retain only the compressed statistic \mathbf{t} .

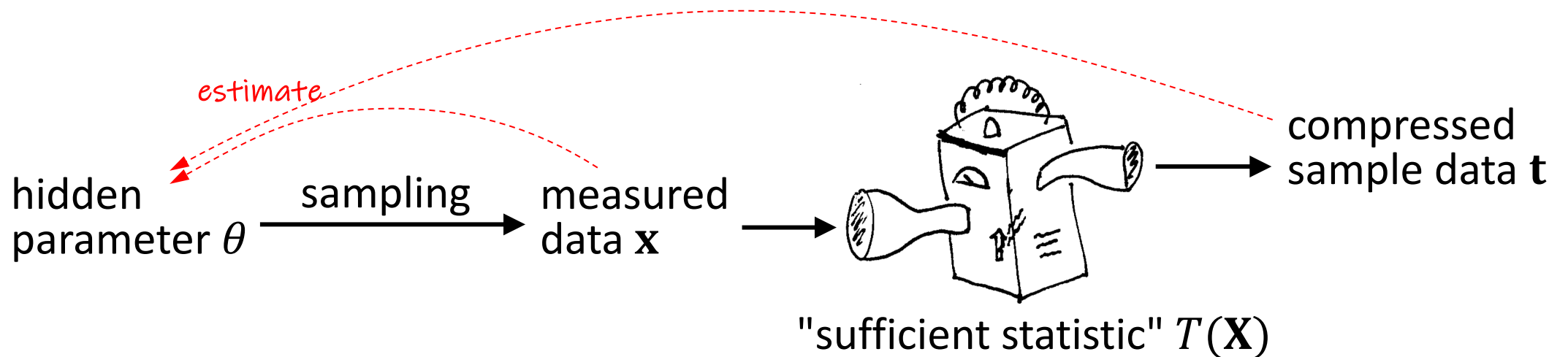


Sufficient statistics

If the sample $\mathbf{x} = (x_1, \dots, x_n)$ and unknown parameter θ , we would like to compress the measurements \mathbf{x} into a low-dimensional statistic without affecting the quality of the possible inference about θ (i.e. we do not want to lose relevant information about θ).

In other words, we are interested in whether there exists a **sufficient statistic** $T(\mathbf{X})$ where the dimension of $\mathbf{t} = T(\mathbf{x})$ is $m < n$, s.t. **\mathbf{t} carries all the useful information from \mathbf{x} about θ** .

If such a sufficient statistic exists, then for the purpose of studying θ , we could discard the raw measurement \mathbf{x} and retain only the compressed statistic \mathbf{t} .



Sufficient statistics in the eyes of information theory

Given a family of distributions $\{f_\theta(x)\}$ indexed a parameter θ . Let $\mathbf{X} = (X_1, \dots, X_n)$ be an iid sample from f_θ , and $T(\mathbf{X})$ be a **statistic** (a quantity computed from the values in the sample). *can also be a vector*

Then $\theta \rightarrow \mathbf{X} \rightarrow T(\mathbf{X})$ forms a Markov chain

From the data processing inequality, we thus know

?

Sufficient statistics in the eyes of information theory

Given a family of distributions $\{f_\theta(x)\}$ indexed a parameter θ . Let $\mathbf{X} = (X_1, \dots, X_n)$ be an iid sample from f_θ , and $T(\mathbf{X})$ be a **statistic** (a quantity computed from the values in the sample). *can also be a vector*

Then $\theta \rightarrow \mathbf{X} \rightarrow T(\mathbf{X})$ forms a Markov chain

From the data processing inequality, we thus know

$$I(\theta; T(\mathbf{X})) \leq I(\theta; \mathbf{X})$$

A statistic is **sufficient** for θ if it preserves all the information in \mathbf{X} about θ :

?

Sufficient statistics in the eyes of information theory

Given a family of distributions $\{f_\theta(x)\}$ indexed a parameter θ . Let $\mathbf{X} = (X_1, \dots, X_n)$ be an iid sample from f_θ , and $T(\mathbf{X})$ be a **statistic** (a quantity computed from the values in the sample). *can also be a vector*

Then $\theta \rightarrow \mathbf{X} \rightarrow T(\mathbf{X})$ forms a Markov chain

From the data processing inequality, we thus know

$$I(\theta; T(\mathbf{X})) \leq I(\theta; \mathbf{X})$$

A statistic is **sufficient** for θ if it preserves all the information in \mathbf{X} about θ :

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

PRACTICAL DEFINITION: A function $T(\mathbf{X})$ is said to be a **sufficient statistic** relative to the family $\{f_\theta(x)\}$ if the conditional distribution of \mathbf{X} given $T(\mathbf{X})$ is independent of θ :

$\theta \perp \mathbf{X} | T(\mathbf{X})$ In other words, $\theta \rightarrow T(\mathbf{X}) \rightarrow \mathbf{X}$ also forms a Markov chain

A possibly helpful way to think about this process is to use a new sample variable: $\theta \rightarrow \mathbf{X} \rightarrow T(\mathbf{X}) = T(\mathbf{X}') \rightarrow \mathbf{X}'$

Example Sufficient statistics

This is the parameter θ

EXAMPLE: Given a sample \mathbf{x} of n iid Bernoulli RVs X_1, \dots, X_n with unknown $\mathbb{P}[X_i = 1] = p$.

Then, given a fixed n , what could be a sufficient statistic $T(\mathbf{X})$ for p ?



Example Sufficient statistics

This is the parameter θ

EXAMPLE: Given a sample \mathbf{x} of n iid Bernoulli RVs X_1, \dots, X_n with unknown $\mathbb{P}[X_i = 1] = p$. Then $k = T(\mathbf{X}) = \sum_i X_i$ is a **sufficient statistic** for θ (assuming n is fixed).



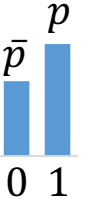
PROOF:

Example Sufficient statistics

This is the parameter θ



EXAMPLE: Given a sample \mathbf{x} of n iid Bernoulli RVs X_1, \dots, X_n with unknown $\mathbb{P}[X_i = 1] = p$. Then $k = T(\mathbf{X}) = \sum_i X_i$ is a **sufficient statistic** for θ (assuming n is fixed).



PROOF: We know that $p \rightarrow \mathbf{X} \rightarrow k$ forms a Markov chain from the fact that k is calculated from \mathbf{X} . To prove that k is a sufficient statistic for p , it is enough to show that $p \rightarrow k \rightarrow \mathbf{X}$ also forms a Markov chain.

We prove that by showing that the conditional distribution of \mathbf{X} given $T(\mathbf{X}) = k$ is independent of θ .

This is a particular sample e.g. $\mathbf{x}=(1, 0, 0, 1, 0, 1, 1)$

$$\mathbb{P}_p[\mathbf{X} = \mathbf{x}] =$$

?

$$\mathbb{P}_p[\mathbf{X} = \mathbf{x} \wedge T(\mathbf{X}) = k] =$$

joint probability

?

$$\mathbb{P}_p[T(\mathbf{X}) = k] =$$

?

$$\mathbb{P}_p[\mathbf{X} = \mathbf{x} | T(\mathbf{X}) = k] =$$

?

Example Sufficient statistics

This is the parameter θ

EXAMPLE: Given a sample \mathbf{x} of n iid Bernoulli RVs X_1, \dots, X_n with unknown $\mathbb{P}[X_i = 1] = p$. Then $k = T(\mathbf{X}) = \sum_i X_i$ is a **sufficient statistic** for θ (assuming n is fixed).



PROOF: We know that $p \rightarrow \mathbf{X} \rightarrow k$ forms a Markov chain from the fact that k is calculated from \mathbf{X} . To prove that k is a sufficient statistic for p , it is enough to show that $p \rightarrow k \rightarrow \mathbf{X}$ also forms a Markov chain.

We prove that by showing that the conditional distribution of \mathbf{X} given $T(\mathbf{X}) = k$ is independent of θ .

$$\mathbb{P}_p[\mathbf{X} = \mathbf{x}] = \prod_i^n p^{x_i} (1-p)^{\bar{x}_i} = p^k (1-p)^{n-k}$$

Very important later: Notice that the density $\mathbb{P}_p[\mathbf{X} = \mathbf{x}]$ depends on \mathbf{x} only through $k = T(\mathbf{X})$. Thus, $\mathbb{P}_p[\mathbf{X} = \mathbf{x}]$ could be written as some function $g(T(\mathbf{x}), \theta)$, which is key to what happens next.

$$\mathbb{P}_p[\mathbf{X} = \mathbf{x} \wedge T(\mathbf{X}) = k] =$$

joint probability

$$\mathbb{P}_p[T(\mathbf{X}) = k] =$$

$$\mathbb{P}_p[\mathbf{X} = \mathbf{x} | T(\mathbf{X}) = k] =$$

Example Sufficient statistics

This is the parameter θ

EXAMPLE: Given a sample \mathbf{x} of n iid Bernoulli RVs X_1, \dots, X_n with unknown $\mathbb{P}[X_i = 1] = p$. Then $k = T(\mathbf{X}) = \sum_i X_i$ is a **sufficient statistic** for θ (assuming n is fixed).



PROOF: We know that $p \rightarrow \mathbf{X} \rightarrow k$ forms a Markov chain from the fact that k is calculated from \mathbf{X} . To prove that k is a sufficient statistic for p , it is enough to show that $p \rightarrow k \rightarrow \mathbf{X}$ also forms a Markov chain.

We prove that by showing that the conditional distribution of \mathbf{X} given $T(\mathbf{X}) = k$ is independent of θ .

$$\mathbb{P}_p[\mathbf{X} = \mathbf{x}] = \prod_i^n p^{x_i} (1-p)^{\bar{x}_i} = p^k (1-p)^{n-k}$$

Very important later: Notice that the density $\mathbb{P}_p[\mathbf{X} = \mathbf{x}]$ depends on \mathbf{x} only through $k = T(\mathbf{X})$. Thus, $\mathbb{P}_p[\mathbf{X} = \mathbf{x}]$ could be written as some function $g(T(\mathbf{x}), \theta)$, which is key to what happens next.

$$\mathbb{P}_p[\mathbf{X} = \mathbf{x} \wedge T(\mathbf{X}) = k] = \begin{cases} \mathbb{P}_p[\mathbf{X} = \mathbf{x}] & \text{if } \sum_i^n x_i = k \\ 0 & \text{otherwise} \end{cases}$$

joint probability

$$\mathbb{P}_p[T(\mathbf{X}) = k] = \binom{n}{k} \cdot p^k (1-p)^{n-k} \quad \text{binomial distribution}$$

$$\mathbb{P}_p[\mathbf{X} = \mathbf{x} | T(\mathbf{X}) = k] = \frac{\mathbb{P}_p[\mathbf{X}=\mathbf{x}, T(\mathbf{X})=k]}{\mathbb{P}_p[k]} = \frac{p^k (1-p)^{n-k}}{\binom{n}{k} \cdot p^k (1-p)^{n-k}} = \begin{cases} \binom{n}{k}^{-1} & \text{if } \sum_i^n x_i = k \\ 0 & \text{otherwise} \end{cases}$$

Thus, we have shown that $\mathbb{P}_p[X|k] = \mathbb{P}[X|k]$ is independent of p . Concretely, all sequences \mathbf{x} with k 1's (and $n-k$ 0's) are equally likely.

Factorization Theorem

In the previous example, we had to guess the sufficient statistic and work out the conditional pmf $\mathbb{P}[\mathbf{X}|T(\mathbf{X}) = T(\mathbf{x})]$ by hand. This can become quite difficult in general.

As we will see next, we didn't really need to go to the trouble of calculating the conditional distribution. Once we noticed that the density $\mathbb{P}_\theta[\mathbf{X} = \mathbf{x}]$ (also $f_\theta(\mathbf{x})$) depends on \mathbf{x} only through $T(\mathbf{x})$, we could have concluded that the statistics $T(\mathbf{X})$ was sufficient.

The easiest way to identify and verify sufficient statistics is to show that the density $f_\theta(\mathbf{x})$ factorizes into a part that involves only the parameter θ and $T(\mathbf{x})$, and a part that involves only \mathbf{x} . This can be used as a working definition of sufficiency.

THEOREM: Let $f_\theta(\mathbf{x})$ (or $f(\mathbf{x}|\theta)$) denote the joint distribution of a data set \mathbf{X} , given parameter θ . A statistic $T(\mathbf{X})$ is a sufficient statistic for θ if and only if there exist functions $g(T(\mathbf{x}), \theta)$ and $h(\mathbf{x})$ such that, for all sample points \mathbf{x} and all parameter points θ , $f_\theta(\mathbf{x})$ factorizes into:

$$f_\theta(\mathbf{x}) = g(T(\mathbf{x}), \theta) \cdot h(\mathbf{x})$$

Notice that the unknown parameter θ interacts with the data \mathbf{x} only via the statistic $T(\mathbf{x})$, and $h(\mathbf{x})$ is independent of θ .

This was $\mathbb{P}_p[\mathbf{X} = \mathbf{x}] = p^k(1-p)^{n-k}$ in the previous example.

Example Sufficient statistics via factorization

EXAMPLE: Given a sample of n iid Bernoulli RVs X_1, \dots, X_n with unknown $\mathbb{P}[X_i = 1] = p$. Then $k = T(\mathbf{X}) = \sum_i X_i$ is a sufficient statistic for θ (assuming n is fixed).

Can you find the factorization $f_p(\mathbf{x}) = g(T(\mathbf{x}), p) \cdot h(\mathbf{x})$ in our earlier proof ?

We prove that by showing that the conditional distribution of \mathbf{X} given $T(\mathbf{X}) = k$ is independent of θ .

$$\mathbb{P}_p[\mathbf{X} = \mathbf{x}] = \prod_i^n p^{x_i} (1-p)^{\bar{x}_i} = p^k (1-p)^{n-k}$$

$$\mathbb{P}_p[\mathbf{X} = \mathbf{x} \wedge T(\mathbf{X}) = k] = \begin{cases} \mathbb{P}_p[\mathbf{X} = \mathbf{x}] & \text{if } \sum_i^n x_i = k \\ 0 & \text{otherwise} \end{cases}$$

$$\mathbb{P}_p[T(\mathbf{X}) = k] = \binom{n}{k} \cdot p^k (1-p)^{n-k}$$

$$\mathbb{P}_p[\mathbf{X} = \mathbf{x} | T(\mathbf{X}) = k] = \frac{\mathbb{P}_p[\mathbf{X}=\mathbf{x}, T(\mathbf{X})=k]}{\mathbb{P}_p[k]} = \frac{\cancel{p^k (1-p)^{n-k}}}{\binom{n}{k} \cdot \cancel{p^k (1-p)^{n-k}}} = \begin{cases} \binom{n}{k}^{-1} & \text{if } \sum_i^n x_i = k \\ 0 & \text{otherwise} \end{cases}$$

Thus, we have shown that $\mathbb{P}_p[\mathbf{X}|k] = \mathbb{P}[\mathbf{X}|k]$ is independent of p .

Example Sufficient statistics via factorization

EXAMPLE: Given a sample of n iid Bernoulli RVs X_1, \dots, X_n with unknown $\mathbb{P}[X_i = 1] = p$. Then $k = T(\mathbf{X}) = \sum_i X_i$ is a sufficient statistic for θ (assuming n is fixed).

Can you find the factorization $f_p(\mathbf{x}) = g(T(\mathbf{x}), p) \cdot h(\mathbf{x})$ in our earlier proof ?

We prove that by showing that the conditional distribution of \mathbf{X} given $T(\mathbf{X}) = k$ is independent of θ .

$$\mathbb{P}_p[\mathbf{X} = \mathbf{x}] = \prod_i^n p^{x_i} (1-p)^{\bar{x}_i} = \underbrace{p^k (1-p)^{n-k}}_{g(k, p)} \cdot \underbrace{1}_{h(\mathbf{x})}$$

$$\mathbb{P}_p[\mathbf{X} = \mathbf{x} \wedge T(\mathbf{X}) = k] = \dots$$

$$\mathbb{P}_p[T(\mathbf{X}) = k] = \binom{n}{k} \cdot p^k (1-p)^{n-k}$$

$$\mathbb{P}_p[\mathbf{X} = \mathbf{x} | T(\mathbf{X}) = k] = \frac{\mathbb{P}_p[\mathbf{X} = \mathbf{x}, T(\mathbf{X}) = k]}{\mathbb{P}_p[k]} = \frac{\cancel{p^k (1-p)^{n-k}}}{\binom{n}{k} \cdot \cancel{p^k (1-p)^{n-k}}} = \begin{cases} \binom{n}{k}^{-1} & \text{if } \sum_i^n x_i = k \\ 0 & \text{otherwise} \end{cases}$$

Thus, we have shown that $\mathbb{P}_p[\mathbf{X}|k] = \mathbb{P}[\mathbf{X}|k]$ is independent of p .

Proof Factorization Theorem (1/2)

PROOF (DISCRETE CASE): sufficient statistics \Leftrightarrow factorization $f_{\theta}(\mathbf{x}) = g(T(\mathbf{x}), \theta) \cdot h(\mathbf{x})$

FIRST DIRECTION sufficient statistics \Rightarrow factorization:

Assume $T(\mathbf{X})$ to be a sufficient statistics, i.e. $\theta \perp \mathbf{X} | T(\mathbf{X})$.

Let $f_{\theta}(\mathbf{x}, T(\mathbf{x}) = t)$ be the joint pdf of $\mathbb{P}_{\theta}[\mathbf{X} = \mathbf{x}, T(\mathbf{X}) = t]$.

$$f_{\theta}(\mathbf{x}) = ?$$

since T is a function of \mathbf{X} , and as long as $t = T(\mathbf{X})$

chain rule

by the definition of sufficient statistics $\theta \perp \mathbf{x} | t$

because t is a function of \mathbf{x} : $t = T(\mathbf{x})$

Proof Factorization Theorem (1/2)

PROOF (DISCRETE CASE): sufficient statistics \Leftrightarrow factorization $f_{\theta}(\mathbf{x}) = g(T(\mathbf{x}), \theta) \cdot h(\mathbf{x})$

FIRST DIRECTION sufficient statistics \Rightarrow factorization:

Assume $T(\mathbf{X})$ to be a sufficient statistics, i.e. $\theta \perp \mathbf{X} | T(\mathbf{X})$.

Let $f_{\theta}(\mathbf{x}, T(\mathbf{x}) = t)$ be the joint pdf of $\mathbb{P}_{\theta}[\mathbf{X} = \mathbf{x}, T(\mathbf{X}) = t]$.

$$\begin{aligned} f_{\theta}(\mathbf{x}) &= f_{\theta}(\mathbf{x}, t) && \text{since } T \text{ is a function of } \mathbf{X}, \text{ and as long as } t = T(\mathbf{X}) \\ &= g_{\theta}(t) \cdot h_{\theta}(\mathbf{x}|t) && \text{chain rule} \\ &= g_{\theta}(t) \cdot h(\mathbf{x}|t) && \text{by the definition of sufficient statistics } \theta \perp \mathbf{x} | t \\ &\quad \swarrow \quad \nwarrow && \\ &g(T(\mathbf{x}), \theta) \quad h(\mathbf{x}) && \text{because } t \text{ is a function of } \mathbf{x}: t = T(\mathbf{x}) \end{aligned}$$

This was $\mathbb{P}_p[\mathbf{X} = \mathbf{x}] = \prod_i^n p^{x_i}(1-p)^{\bar{x}_i} = p^k(1-p)^{n-k}$ in the previous example.

Proof Factorization Theorem (2/2)



OTHER DIRECTION: factorization \Rightarrow sufficient statistics:

Assume $f_{\theta}(\mathbf{x}) = g(t, \theta) \cdot h(\mathbf{x})$.

We need to show that the conditional probability distribution $f_{\theta}(\mathbf{x}|t)$ of \mathbf{X} given $T(\mathbf{X})$ is independent of θ , i.e. $f_{\theta}(\mathbf{x}|t) = f(\mathbf{x}|t)$.

$$f_{\theta}(t) = ?$$

definition of marginal probability distribution

since t is a function of \mathbf{x}

using our assumption

factoring out a common factor

$$f_{\theta}(\mathbf{x}|t) = ?$$

definition of conditional probability distribution

does not depend on θ , hence T is a sufficient statistic

Proof Factorization Theorem (2/2)

OTHER DIRECTION: factorization \Rightarrow sufficient statistics:

Assume $f_{\theta}(\mathbf{x}) = g(t, \theta) \cdot h(\mathbf{x})$.

We need to show that the conditional probability distribution $f_{\theta}(\mathbf{x}|t)$ of \mathbf{X} given $T(\mathbf{X})$ is independent of θ , i.e. $f_{\theta}(\mathbf{x}|t) = f(\mathbf{x}|t)$.

$$f_{\theta}(t) = \sum_{\mathbf{x}:T(\mathbf{x})=t} f_{\theta}(\mathbf{x}, t)$$

definition of marginal probability distribution

$$= \sum_{\mathbf{x}:T(\mathbf{x})=t} f_{\theta}(\mathbf{x})$$

since t is a function of \mathbf{x}

$$= \sum_{\mathbf{x}:T(\mathbf{x})=t} g(t, \theta) \cdot h(\mathbf{x})$$

using our assumption

$$= g(t, \theta) \cdot \sum_{\mathbf{x}:T(\mathbf{x})=t} h(\mathbf{x})$$

factoring out a common factor

$$f_{\theta}(\mathbf{x}|t) = \frac{f_{\theta}(\mathbf{x}, t)}{f_{\theta}(t)} = \frac{f_{\theta}(\mathbf{x})}{f_{\theta}(t)}$$

definition of conditional probability distribution

$$= \frac{\cancel{g(t, \theta)} \cdot h(\mathbf{x})}{\cancel{g(t, \theta)} \cdot \sum_{\mathbf{x}:T(\mathbf{x})=t} h(\mathbf{x})}$$

does not depend on θ , hence T is a sufficient statistic

Sufficient Statistics & Factorization Theorem

The concept of sufficient statistics is due to Sir Ronald Fisher around 1920, thus before the advent of information theory.

The factorization theorem is also varyingly called:

- Fisher's factorization theorem
- Fisher-Neyman factorization theorem
- Neyman-Fisher factorization theorem
- Halmos-Savage factorization theorem



Sir Ronald Fisher (1890–1962)

Normal (Gaussian) distribution: (μ, σ^2) are sufficient statistics

Example 6.2.9 (Normal sufficient statistic, both parameters unknown) Again assume that X_1, \dots, X_n are iid $n(\mu, \sigma^2)$ but, unlike Example 6.2.4, assume that both μ and σ^2 are unknown so the parameter vector is $\theta = (\mu, \sigma^2)$. Now when using the Factorization Theorem, any part of the joint pdf that depends on either μ or σ^2 must be included in the g function. From (6.2.1) it is clear that the pdf depends on the sample \mathbf{x} only through the two values $T_1(\mathbf{x}) = \bar{x}$ and $T_2(\mathbf{x}) = s^2 = \sum_{i=1}^n (x_i - \bar{x})^2 / (n - 1)$. Thus, we can define $h(\mathbf{x}) = 1$ and

$$\begin{aligned} g(\mathbf{t}|\theta) &= g(t_1, t_2 | \mu, \sigma^2) \\ &= (2\pi\sigma^2)^{-n/2} \exp\left(-\left(n(t_1 - \mu)^2 + (n-1)t_2\right) / (2\sigma^2)\right). \end{aligned}$$

Then it can be seen that

$$f(\mathbf{x}|\mu, \sigma^2) = g(T_1(\mathbf{x}), T_2(\mathbf{x}) | \mu, \sigma^2) h(\mathbf{x}). \quad (6.2.5)$$

Thus, by the Factorization Theorem, $T(\mathbf{X}) = (T_1(\mathbf{X}), T_2(\mathbf{X})) = (\bar{X}, S^2)$ is a sufficient statistic for (μ, σ^2) in this normal model. ||

Example 6.2.9 demonstrates that, for the normal model, the common practice of summarizing a data set by reporting only the sample mean and variance is justified. The sufficient statistic (\bar{X}, S^2) contains all the information about (μ, σ^2) that is available in the sample.

Exponential Family

The definition in terms of one *real-number* parameter can be extended to one *real-vector* parameter

$$\boldsymbol{\theta} \equiv [\theta_1, \theta_2, \dots, \theta_s]^\top .$$

A family of distributions is said to belong to a vector exponential family if the probability density function (or probability mass function, for discrete distributions) can be written as

$$f_X(x \mid \boldsymbol{\theta}) = h(x) g(\boldsymbol{\theta}) \exp \left(\boldsymbol{\eta}(\boldsymbol{\theta}) \cdot \mathbf{T}(x) \right)$$

- $T(x)$ is a *sufficient statistic* of the distribution. For exponential families, the sufficient statistic is a function of the data that holds all information the data x provides with regard to the unknown parameter values. This means that, for any data sets x and y , the likelihood ratio is the same, that is $\frac{f(x; \theta_1)}{f(x; \theta_2)} = \frac{f(y; \theta_1)}{f(y; \theta_2)}$ if $T(x) = T(y)$. This is true even if x and y are not equal to each other.

The dimension of $T(x)$ equals the number of parameters of θ and encompasses all of the information regarding the data related to the parameter θ . The sufficient statistic of a set of *independent identically distributed* data observations is simply the sum of individual sufficient statistics, and encapsulates all the information needed to describe the *posterior distribution* of the parameters, given the data (and hence to derive any desired estimate of the parameters). (This important property is discussed further *below*.)

Exponential families have a large number of properties that make them extremely useful for statistical analysis. In many cases, it can be shown that *only* exponential families have these properties. Examples:

- Exponential families are the only families with *sufficient statistics* that can summarize arbitrary amounts of *independent identically distributed* data using a fixed number of values. (*Pitman–Koopman–Darmois* theorem)

Aggregates in Databases

Data Cube: A Relational Aggregation Operator Generalizing Group-By, Cross-Tab, and Sub-Totals*

JIM GRAY
SURAJIT CHAUDHURI
ADAM BOSWORTH
ANDREW LAYMAN
DON REICHART
MURALI VENKATRAO

Microsoft Research, Advanced Technology Division, Microsoft Corporation, One Microsoft Way, Redmond, WA 98052

FRANK PELLOW
HAMID PIRAHESH
IBM Research, 500 Harry Road, San Jose, CA 95120

Gray@Microsoft.com
SurajitC@Microsoft.com
AdamB@Microsoft.com
AndrewL@Microsoft.com
DonRei@Microsoft.com
MuraliV@Microsoft.com

Pellow@vnet.IBM.com
Pirahesh@Almaden.IBM.com

ICDE Influential Paper Awards

ICDE 2006

Jim Gray, Adam Bosworth, Andrew Layman, Hamid Pirahesh

[Data Cube: A Relational Aggregation Operator Generalizing Group-By, Cross-Tab, and Sub-Total](#), ICDE 1996

Citation: This seminal paper defined a simple SQL construct that enables one to efficiently compute aggregations over all combinations of group-by columns in a single query, where previous approaches required multiple queries. This feature has had significant impact on industry and is now incorporated in all major database systems.

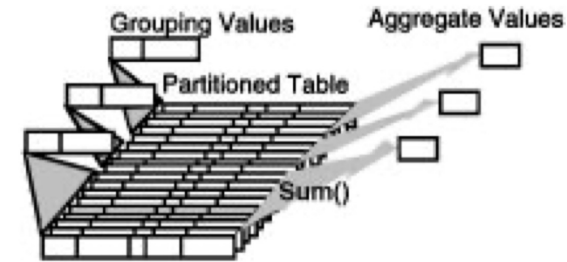


Figure 2. The GROUP BY relational operator partitions a table into groups. Each group is then aggregated by a function. The aggregation function summarizes some column of groups returning a value for each group.

Consider aggregating a two dimensional set of values $\{X_{ij} \mid i = 1, \dots, I; j = 1, \dots, J\}$. Aggregate functions can be classified into three categories:

Distributive: Aggregate function $F()$ is distributive if there is a function $G()$ such that $F(\{X_{i,j}\}) = G(\{F(\{X_{i,j} \mid i = 1, \dots, I\} \mid j = 1, \dots, J)\})$. COUNT(), MIN(), MAX(), SUM() are all distributive. In fact, $F = G$ for all but COUNT(). $G = \text{SUM}()$ for the COUNT() function. Once order is imposed, the cumulative aggregate functions also fit in the distributive class.

Algebraic: Aggregate function $F()$ is algebraic if there is an M -tuple valued function $G()$ and a function $H()$ such that $F(\{X_{i,j}\}) = H(\{G(\{X_{i,j} \mid i = 1, \dots, I\} \mid j = 1, \dots, J)\})$. Average(), standard deviation, MaxN(), MinN(), center_of_mass() are all algebraic. For Average, the function $G()$ records the sum and count of the subset. The $H()$ function adds these two components and then divides to produce the global average. Similar techniques apply to finding the N largest values, the center of mass of group of objects, and other algebraic functions. The key to algebraic functions is that a fixed size result (an M -tuple) can summarize the sub-aggregation.

Holistic: Aggregate function $F()$ is holistic if there is no constant bound on the size of the storage needed to describe a sub-aggregate. That is, there is no constant M , such that an M -tuple characterizes the computation $F(\{X_{i,j} \mid i = 1, \dots, I\})$. Median(), MostFrequent() (also called the Mode()), and Rank() are common examples of holistic functions.

We know of no more efficient way of computing super-aggregates of holistic functions than the 2^N -algorithm using the standard GROUP BY techniques. We will not say more about cubes of holistic functions.

Information Inequalities

Best reference:

[Yeung'08] Yeung, Information Theory and Network Coding, 2008. Ch 2.6, 2.7, 13, 14, 15

<http://iest2.ie.cuhk.edu.hk/~whyung/tempo/main2.pdf>

Basic inequalities

Shannon's information measures refer to entropy, conditional entropy, mutual information, and conditional mutual information (but not interaction information!).

They can be expressed as linear combinations of entropies:

$$H(X|Y) = ?$$

$$I(X; Y) = ?$$

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

Basic inequalities

Shannon's information measures refer to entropy, conditional entropy, mutual information, and conditional mutual information (but not interaction information!).

They can be expressed as linear combinations of entropies:

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

$$I(X; Y) = H(X) + H(Y) - H(X, Y)$$

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

by repeated expansion of conditional entropies; also holds if we replace variables with sets of variables

They are also special cases of conditional mutual information.

$$\begin{aligned} H(X) &= ? \\ H(X|Z) &= ? \\ I(X; Y) &= ? \end{aligned}$$

Assume φ to be degenerate RV that takes on a constant value with probability 1

Basic inequalities

Shannon's information measures refer to entropy, conditional entropy, mutual information, and conditional mutual information (but not interaction information!).

They can be expressed as linear combinations of entropies:

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

$$I(X; Y) = H(X) + H(Y) - H(X, Y)$$

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

also holds if we replace variables with sets of variables

They are also special cases of conditional mutual information.

$$H(X) = I(X; X|\varphi)$$

$$H(X|Z) = I(X; X|Z)$$

$$I(X; Y) = I(X; Y|\varphi)$$

Assume φ to be degenerate RV that takes on a constant value with probability 1

With the basic inequalities we refer to the fact that all Shannon's information measures are non-negative (because conditional mutual information is ≥ 0).

$$I(U; V|W) \geq 0$$

U, V, W can be arbitrary joint entropies

Shannon-type inequalities Γ_n (and constraints)

Shannon-type inequalities are inequalities on information measures implied by the basic inequalities and possibly additional constraints on the joint distribution of the RVs involved.

EXAMPLE: data-processing inequality for $X \rightarrow Y \rightarrow Z$:

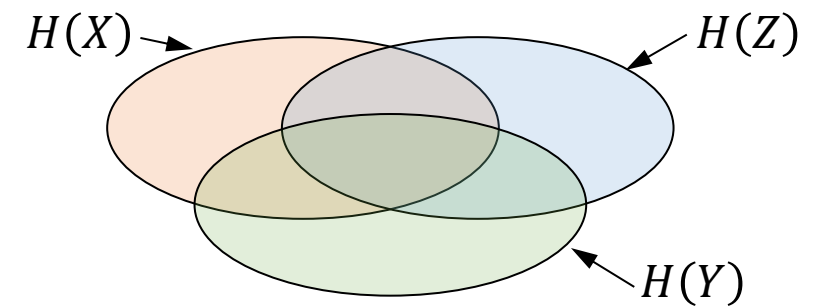
From $I(X; Z|Y) = 0$ and basic inequalities, we derived $I(X; Y) \geq I(X; Z)$

not a basic inequality

Information inequalities are the inequalities that govern the impossibilities in information theory. They imply that certain things cannot happen. For this reason, they are sometimes referred to as the **laws of information theory**.

EXAMPLE : $n = 3$ variables with given $k = 2^3 - 1 = 7$ joint entropies:

$$\begin{array}{lll} H(X) = 2 & H(X, Y) = 4 & H(X, Y, Z) = 5 \\ H(Y) = 3 & H(X, Z) = 4 & \\ H(Z) = 4 & H(Y, Z) = 4 & \end{array}$$



Find 3 RVs that fulfill those constraints ?

Shannon-type inequalities Γ_n (and constraints)

Shannon-type inequalities are inequalities on information measures implied by the basic inequalities and possibly additional constraints on the joint distribution of the RVs involved.

EXAMPLE: data-processing inequality for $X \rightarrow Y \rightarrow Z$:

From $I(X; Z|Y) = 0$ and basic inequalities, we derived $I(X; Y) \geq I(X; Z)$

not a basic inequality

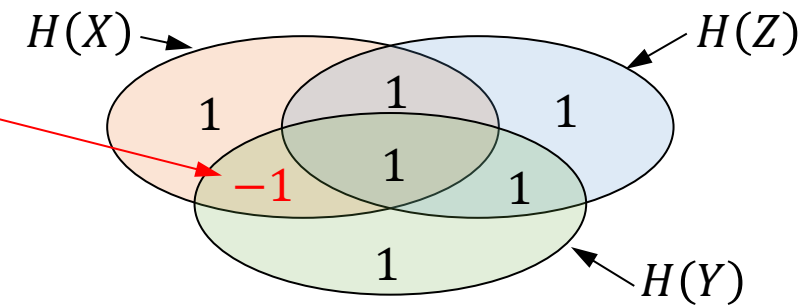
Information inequalities are the inequalities that govern the impossibilities in information theory. They imply that certain things cannot happen. For this reason, they are sometimes referred to as the **laws of information theory**.

EXAMPLE : $n = 3$ variables with given $k = 2^3 - 1 = 7$ joint entropies:

$$\begin{array}{lll} H(X) = 2 & H(X, Y) = 4 & H(X, Y, Z) = 5 \\ H(Y) = 3 & H(X, Z) = 4 & \\ H(Z) = 4 & H(Y, Z) = 4 & \end{array}$$

$$I(X; Y|Z) \neq 0$$

not possible ☹️



Almost all the information inequalities known to date are Shannon-type inequalities and thus implied by the basic inequalities.

Applications in Databases

Worst-case Optimal Join Algorithms

HUNG Q. NGO, University at Buffalo, SUNY
ELY PORAT, Bar-Ilan University
CHRISTOPHER RÉ, Stanford University
ATRI RUDRA, University at Buffalo, SUNY

PODS'12, JACM'18

Efficient join processing is one of the most fundamental and well-studied tasks in database research. In this work, we examine algorithms for natural join queries over many relations and describe a new algorithm to process these queries optimally in terms of worst-case data complexity. Our result builds on recent work by Atserias, Grohe, and Marx, who gave bounds on the size of a natural join query in terms of the sizes of the individual relations in the body of the query. These bounds, however, are not constructive: they rely on Shearer's entropy inequality, which is information-theoretic. Thus, the previous results leave open the question of whether there exist algorithms whose runtimes achieve these optimal bounds. An answer to this question may be interesting to database practice, as we show in this article that any project-join style plans, such as ones typically employed in a relational database management system, are asymptotically slower than the optimal for some queries. We present an algorithm whose runtime is worst-case optimal for all natural join queries. Our result may be of independent interest, as our algorithm also yields a constructive proof of the general fractional cover bound by Atserias, Grohe, and Marx without using Shearer's inequality. This bound implies two famous inequalities in geometry: the Loomis-Whitney inequality and its generalization, the Bollobás-Thomason inequality. Hence, our results algorithmically prove these inequalities as well. Finally, we discuss how our algorithm can be used to evaluate full conjunctive queries optimally, to compute a relaxed notion of joins and to optimally (in the worst-case) enumerate all induced copies of a fixed subgraph inside of a given large graph.

Decision Problems in Information Theory

Mahmoud Abo Khamis

relationalAI, Berkeley, CA, USA

Phokion G. Kolaitis

UC Santa Cruz and IBM Research - Almaden, Santa Cruz, CA, USA

Hung Q. Ngo

relationalAI, Berkeley, CA, USA

Dan Suciu

University of Washington, Seattle, WA, USA

ICALP'20

B.2 Application to Relational Query Evaluation

The problem of bounding the number of copies of a graph inside of another graph has a long and interesting history [17, 4, 14, 35]. The subgraph homomorphism problem is a special case of the relational query evaluation problem, in which we want to find an upper bound on the output size of a full conjunctive query. Using the entropy argument from [14], Shearer's lemma in particular, Atserias, Grohe, and Marx [5] established a tight upper bound on the answer to a full conjunctive query over a database. Note that Shearer's lemma is a Shannon-type inequality. Their result was extended to include functional dependencies and more generally degree constraints in a series of recent work in database theory [19, 2, 3]. All these results can be cast as applications of Shannon-type inequalities. For a simple example, let $R(X, Y), S(Y, Z), T(Z, U)$ be three binary relations (tables), each with N tuples, then their join $R(X, Y) \bowtie S(Y, Z) \bowtie T(Z, U)$ can be as large as N^2 tuples. However, if we further know that the functional dependencies $XZ \rightarrow U$ and $YU \rightarrow X$ hold in the output, then one can prove that the output size is $\leq N^{3/2}$, by using the following Shannon-type information inequality:

$$h(XY) + h(YZ) + h(ZU) + h(X|YU) + h(U|XZ) \geq 2h(XYZU) \quad (24)$$

Ngo, Porat, Re, Rudra. Worst-case Optimal Join Algorithms, JAC 2018 (PODS 2012). <https://doi.org/10.1145/3180143>,

Khamis, Kolaitis, Ngo, Suciu, "Decision Problems in Information Theory", ICALP 2020. <https://doi.org/10.4230/LIPIcs.ICALP.2020.106>

Gatterbauer, Aslam. Foundations and Applications of Information Theory: <https://northeastern-datalab.github.io/cs7840/>

Applications in Databases

1.1 The problems

In the history of query evaluation in database, logic and constraint satisfaction areas, there are three research threads which have yielded spectacular results recently.

Thread 1: size bound for conjunctive queries. From the seminal work of Grohe and Marx [33], Atserias, Grohe, and Marx [8], and Gottlob, Lee, Valiant and Valiant [30], we now know of a deep connection between the output size bound of a conjunctive query with (or without) functional dependencies (FD) and information theory. In particular, we can derive tight output size bounds by solving a convex optimization problem whose variables are marginal entropies. Briefly, the bound works as follows. Consider a conjunctive query Q represented by a multi-hypergraph $\mathcal{H} = (V, \mathcal{E})$, where $V = [n]$ is identified with the set of variables A_1, \dots, A_n . To each hyperedge $F \in \mathcal{E}$ there is an input relation R_F whose attributes are $(A_i)_{i \in F}$. A function $h : 2^V \rightarrow \mathbb{R}_+$ is called *entropic* if there exists a joint distribution on n variables such that $h(F)$ is the marginal entropy of the distribution on the variables in F , for every non-empty set $F \subseteq V$; by convention, $h(\emptyset) = 0$. Let Γ_n^* denote the set of all n -variable entropic functions.¹ Let CC denote the set of “cardinality constraints” of the form $N_F = |R_F|$, obtained from the input database instance. Let FD denote the set of “FD constraints” of the form $X \rightarrow Y$, where $\emptyset \subseteq X \subset Y \subseteq V$.² From the cardinality- and FD-constraints, we define two classes of set functions:

What do Shannon-type Inequalities, Submodular Width, and Disjunctive Datalog have to do with one another?

Mahmoud Abo Khamis
LogicBlox Inc.

Hung Q. Ngo
LogicBlox Inc.

Dan Suciu
LogicBlox Inc. and
University of Washington

PODS'17

ABSTRACT

Recent works on bounding the output size of a conjunctive query with functional dependencies and degree bounds have shown a deep connection between fundamental questions in information theory and database theory. We prove analogous output bounds for disjunctive datalog rules, and answer several open questions regarding the tightness and looseness of these bounds along the way. The bounds are intimately related to Shannon-type information inequalities. We devise the notion of a “proof sequence” of a specific class of Shannon-type information inequalities called “Shannon flow inequalities”. We then show how a proof sequence can be used as symbolic instructions to guide an algorithm called PANDA, which answers disjunctive datalog rules within the size bound predicted. We show that PANDA can be used as a black-box to devise algorithms matching precisely the fractional hypertree width and the submodular width runtimes for aggregate and conjunctive queries *with* functional dependencies and degree bounds.

Our results improve upon known results in three ways. First, our bounds and algorithms are for the much more general class of disjunctive datalog rules, of which conjunctive queries are a special case. Second, the runtime of PANDA matches precisely the submodular width bound, while the previous algorithm by Marx has a runtime that is polynomial in this bound. Third, our bounds and algorithms work for queries with input cardinality bounds, functional dependencies, *and* degree bounds.

Overall, our results showed a deep connection between three seemingly unrelated lines of research; and, our results on proof sequences for Shannon flow inequalities might be of independent interest.

Applications in Databases

Applications of Information Inequalities to Database Theory Problems *

Dan Suciu
University of Washington
Seattle, USA
suciu@cs.washington.edu

June 6, 2024

LICS'23 keynote

Abstract

The paper describes several applications of information inequalities to problems in database theory. The problems discussed include: upper bounds of a query's output, worst-case optimal join algorithms, the query domination problem, and the implication problem for approximate integrity constraints. The paper is self-contained: all required concepts and results from information inequalities are introduced here, gradually, and motivated by database problems.

Information Inequalities v.s. Databases

Informally: $h(XY) \sim \log |\Pi_{XY}(R)|$. What do inequalities say about R ?

- $h(X) \leq h(XY) \leq h(XYZ)$
Says $|\Pi_X(R)| \leq |\Pi_{XY}(R)| \leq |R|$.
- $h(XY) + h(Z) \geq h(XYZ)$
Says $|\Pi_{XY}(R)| \cdot |\Pi_Z(R)| \geq |R|$.
- $h(XYZ|X) \geq h(XYZ|XY)$
Max frequency(X) is \geq max frequency(XY).
- **Careful!** $h(XZ) + h(YZ) \geq h(XYZ) + h(Z)$,
but $\underbrace{|\Pi_{XZ}(R)|}_3 \cdot \underbrace{|\Pi_{YZ}(R)|}_3 \not\geq \underbrace{|R|}_5 \cdot \underbrace{|\Pi_Z(R)|}_2$

X	Y	Z
a	x	m
a	y	m
b	x	m
b	y	m
a	x	n

Information inequalities Γ_n^*

Information inequalities are the inequalities that govern the impossibilities in information theory. They imply that certain things cannot happen. For this reason, they are sometimes referred to as the **laws of information theory**.

An information inequality or identity involves (linear combinations of) Shannon's information measures only (and possibly with constant terms) and is said to always hold if it holds for any joint distribution for the random variables involved.

There exist laws in information theory that are not implied by the basic inequalities (called non-Shannon-type inequalities). This celebrated result was published by [Zhang, Yeung'98]

PROPOSITION: The following information inequality always holds on any list of five random variables X, Y, Z, U, V , but is not implied by the basic inequalities:

$$H(X) + H(Y) + I(U; V|X) + I(U; V|Y) + 2I(U; V|Z) + I(U, V; Z) \geq H(X, Z) + 2I(U; V)$$

Key proof insight: $I(XY; Z|UV) = 0$ can be assumed for a different argument

A formalization of entropic vectors

- Given a set of n RVs $\Theta = \{X_1, \dots, X_n\}$, written as $\Theta = \{X_i\}, i \in [n]$. ↖ $\{1, 2, \dots, n\}$
- Associated with Θ are $2^n - 1$ joint entropies $H(X_1), \dots, H(X_1, \dots, X_n)$, written as $H_\Theta(\alpha) = H(X_\alpha)$ for any subset of $[n]$. Call the function $H_\Theta(\alpha), \alpha \in 2^{[n]}$ the **entropy function** of Θ .
- Example: $H(X_1, X_2, X_4)$ is $H_\Theta(\alpha)$ for $\alpha = \{1, 2, 4\}$.
- Together, the joint entropies form a point in the $2^n - 1$ dimensional **entropy space** $\mathbb{R}^{2^n - 1}$.
- In turn, a point in that space is called **entropic** if the point corresponds to the entropy function H_Θ of some set Θ of n RVs. Let $\Gamma_n^* \subset \mathbb{R}^{2^n - 1}$ be the **set of all entropic points**.
- How does that space $\Gamma_n^* \subset \mathbb{R}^{2^n - 1}$ look like?

Our earlier EXAMPLE: $n = 3$, thus $k = 2^3 - 1 = 7$ joint entropies, representing a point in \mathbb{R}^7

$$\begin{array}{lll} H(X) = 2 & H(X, Y) = 4 & H(X, Y, Z) = 5 \\ H(Y) = 3 & H(X, Z) = 4 & \\ H(Z) = 4 & H(Y, Z) = 4 & \end{array}$$



Entropic vectors

- Given a set of n RVs $\Theta = \{X_1, \dots, X_n\}$, written as $\Theta = \{X_i\}, i \in [n]$.
- Associated with Θ are $2^n - 1$ joint entropies $H(X_1), \dots, H(X_1, \dots, X_n)$, written as $H_\Theta(\alpha) = H(X_\alpha)$ for any subset of $[n]$. Call the function $H_\Theta(\alpha), \alpha \in 2^{[n]}$ the **entropy function** of Θ .
- Example: $H(X_1, X_2, X_4)$ is $H_\Theta(\alpha)$ for $\alpha = \{1, 2, 4\}$.
- Together, the joint entropies form a point in the $2^n - 1$ dimensional **entropy space** $\mathbb{R}^{2^n - 1}$.
- In turn, a point in that space is called **entropic** if the point corresponds to the entropy function H_Θ of some set Θ of n RVs. Let $\Gamma_n^* \subset \mathbb{R}^{2^n - 1}$ be the **set of all entropic points**.
- How does that space $\Gamma_n^* \subset \mathbb{R}^{2^n - 1}$ look like?

Our earlier EXAMPLE: $n = 3$, thus $k = 2^3 - 1 = 7$ joint entropies, representing a point in \mathbb{R}^7

$$H(X) = 2 \quad H(X, Y) = 4 \quad H(X, Y, Z) = 5$$

$$H(Y) = 3 \quad H(X, Z) = 4$$

$$H(Z) = 4 \quad H(Y, Z) = 4$$

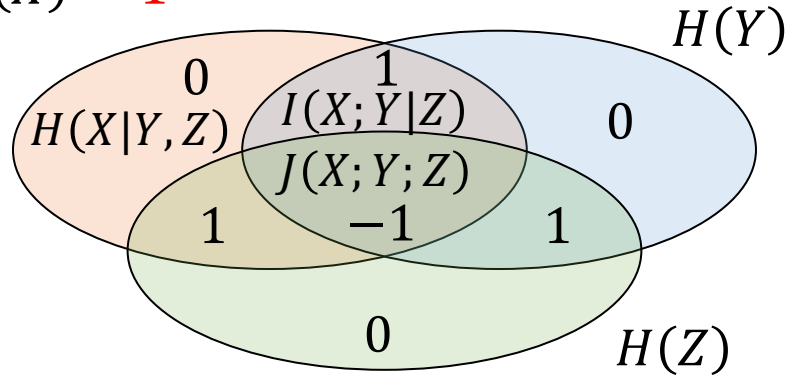
$$I(X; Y|Z) \neq 0$$

Thus this point $(2, 3, 4, 4, 4, 5) \notin \Gamma_3^*$?

A subtlety: entropic vectors Γ_n^* vs. almost entropic vectors $\bar{\Gamma}_n^*$

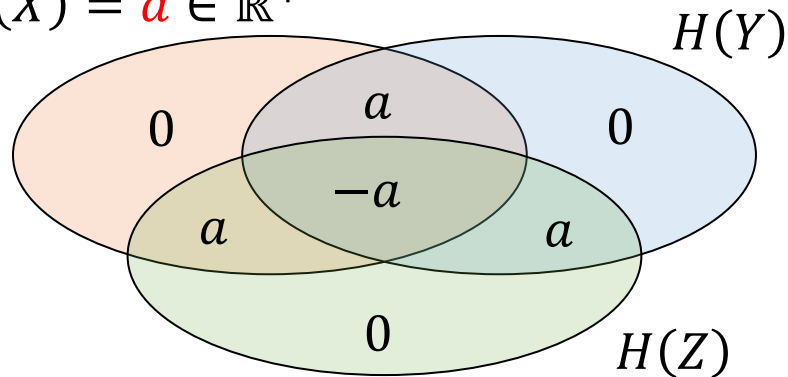
Our earlier "parity example":

$$H(X) = 1$$



More generally (from basic inequalities):

$$H(X) = a \in \mathbb{R}^+$$



However, a more careful analysis shows that all variables X, Y, Z need to be uniform for this example to work, which implies only discrete particular entropies as possible.

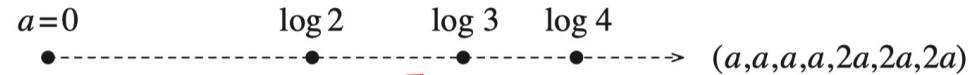


Fig. 15.2. The values of a for which $(a, a, a, 2a, 2a, 2a)$ is in Γ_3 .

Γ_n^* set of all entropic vectors

$\bar{\Gamma}_n^*$ set of all almost entropic vectors:
defined as **topological closure** of Γ_n^*

Γ_n subset of vectors that fulfill the
Shannon inequalities

The closure of a subset S of points in a topological space consists of all points in S **together with all limit points** of S .

Intuitively, it is possible to create a mixture model that models any rational number. The "closure" extends that to the real numbers.

