

COMP2610 / COMP6261 - Information Theory

Lecture 8: Some Fundamental Inequalities

Assignment Project Exam Help

<https://eduassistpro.github.io>



Australian
National
University

Add WeChat edu_assist_pro

14 August 2018

Last time

Assignment Project Exam Help

- Decomposability of entropy

- Rel

<https://eduassistpro.github.io>

Add WeChat edu_assist_pr

- Mutual information

Review

Relative entropy (KL divergence):

$$D_{\text{KL}}(p||q) = \sum_{x \in \mathcal{X}} p(x) \log \frac{p(x)}{q(x)}$$

Assignment Project Exam Help

Mutual inf

<https://eduassistpro.github.io>

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

- Average reduction in uncertainty in X wh
- $I(X; Y) = 0$ when X, Y statistically indepe

Add WeChat edu_assist_pr

Conditional mutual information of X, Y given Z :

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

This time

Assignment Project Exam Help

Mutual information chain rule

Jensen'

"Informa

Data processing inequality

<https://eduassistpro.github.io>

Add WeChat edu_assist_pr

Outline

1 Chain Rule for Mutual Information

2 Convex Functions

3 Jens

4 Gibb

5 Information Cannot Hurt

6 Data Processing Inequality

7 Wrapping Up

Assignment Project Exam Help

<https://eduassistpro.github.io>

Add WeChat edu_assist_pr

1 Chain Rule for Mutual Information

2 Convex Functions

3 Jensen's Inequality

4 Gibbs

5 Information Cannot Hurt

6 Data Processing Inequality

7 Wrapping Up

Assignment Project Exam Help

<https://eduassistpro.github.io>

Add WeChat edu_assist_pr

Recall: Joint Mutual Information

Recall the mutual information between X and Y :

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

Assignment Project Exam Help

We can also

X_N and

$Y_1, \dots,$

<https://eduassistpro.github.io>

$I(X_1, \dots, X_N; Y_1, \dots, Y_M)$

$- H(X)$

Add WeChat: edu_assist_pro

Note that $I(X, Y; Z) \neq I(X; Y, Z)$ in general

- Reduction in uncertainty of X and Y given Z versus reduction in uncertainty of X given Y and Z

Chain Rule for Mutual Information

Let X, Y, Z be r.v. and recall that:

$$p(Z, Y) = p(Z|Y)p(Y)$$

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

$$I(X; Y, Z) = I(Y, Z; X) \quad \text{symmetry}$$

<https://eduassistpro.github.io>

$$I(X; Y, Z) = \underbrace{I(X; Y)}_{I(Y; X)} + \underbrace{I(X; Z|Y)}_{\text{defini}} \quad \text{edu_assist_pr}$$

Similarly, by symmetry:

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

Chain Rule for Mutual Information

General form

For any collection of random variables X_1, \dots, X_N and Y :

$$I(X_1, \dots, X_N; Y) = I(X_1; Y) + I(X_2, \dots, X_N; Y | X_1)$$

<https://eduassistpro.github.io>

Add WeChat edu_assist_pro

$$\begin{aligned} &= \sum_{i=1}^N I(X_i; Y | X_1, \dots, X_{i-1}) \\ &= \sum_{i=1}^N I(Y; X_i | X_1, \dots, X_{i-1}). \end{aligned}$$

1 Chain Rule for Mutual Information

2 Convex Functions

3 Jensen's Inequality

4 Gibbs

5 Information Cannot Hurt

6 Data Processing Inequality

7 Wrapping Up

Assignment Project Exam Help

<https://eduassistpro.github.io>

Add WeChat edu_assist_pr

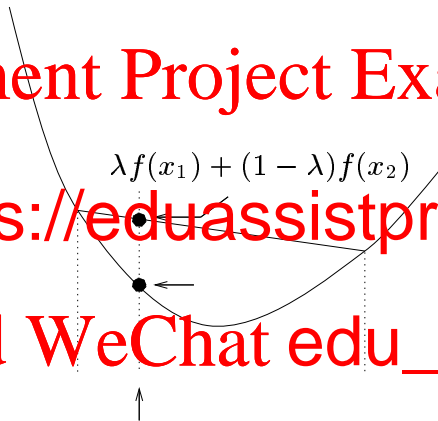
Convex Functions:

Introduction

Assignment Project Exam Help

<https://eduassistpro.github.io>

Add WeChat edu_assist_pr



$$0 \leq \lambda \leq 1$$

(Figure from Mackay, 2003)

A function is convex \smile if every chord of the function lies above the function

Convex and Concave Functions

Definitions

Definition

A function f is **convex** if for all x_1, x_2 and

$$0 \leq \lambda \leq 1$$

We say f is **strictly convex** \curvearrowright if for all x_1, x_2 and λ only
for $\lambda = 0$ and $\lambda = 1$.

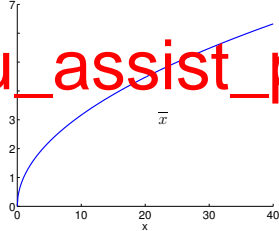
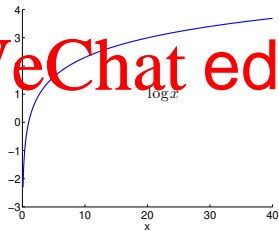
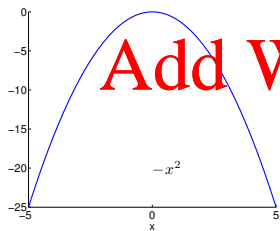
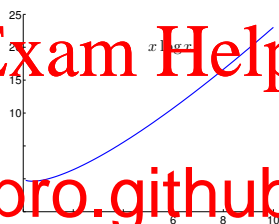
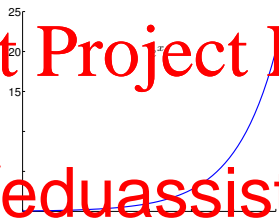
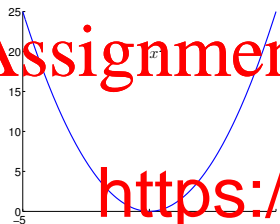
Similarly, a function f is **concave** \curvearrowleft if $-f$ is strictly convex
the function lies below the function.

Examples of Convex and Concave Functions

Assignment Project Exam Help

<https://eduassistpro.github.io>

Add WeChat edu_assist_pro



Verifying Convexity

Theorem (Cover & Thomas, Th 2.6.1)

If a function f has a second derivative that is non-negative (positive) over an interval, the function is convex (strictly convex) over that interval.

This allow

Example

- x^2 : $\frac{d}{dx} \left(\frac{d}{dx} (x^2) \right) = \frac{d}{dx} (2x) = 2$

- e^x : $\frac{d}{dx} \left(\frac{d}{dx} (e^x) \right) = \frac{d}{dx} (e^x) = e^x$

- $\sqrt{x}, x > 0$: $\frac{d}{dx} \left(\frac{d}{dx} (\sqrt{x}) \right) = \frac{1}{2} \frac{d}{dx} \left(\frac{1}{\sqrt{x}} \right) = -\frac{1}{4} \frac{1}{\sqrt{x}^3}$

Convexity, Concavity and Optimization

If $f(x)$ is concave \cap and there exists a point at which

Assignment Project Exam Help

$$\frac{df}{dx} = 0,$$

then $f(x)$

Note: the $f(x)$ is maximized at some x , it is

<https://eduassistpro.github.io>

- $f(x) = -|x|$: is maximized at $x = 0$ wh

Add WeChat edu_assist_pro

- $f(p) = \log p$ with $0 \leq p \leq 1$, is maximiz $\frac{df}{dp} = 1$

- Similarly for minimisation of convex functions

1 Chain Rule for Mutual Information

2 Convex Functions

3 Jensen's Inequality

4 Gibbs

5 Information Cannot Hurt

6 Data Processing Inequality

7 Wrapping Up


Assignment Project Exam Help

<https://eduassistpro.github.io>

Add WeChat edu_assist_pr

Jensen's Inequality for Convex Functions

Theorem: Jensen's Inequality

If f is a **convex**  function and X is a random variable then:

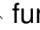
$$f(\mathbb{E}[X]) \leq \mathbb{E}[f(X)].$$

Moreover,
probabili

$\mathbb{E}[X]$ with

In other words, for a probability vector \mathbf{p} ,

$$f\left(\sum_{i=1}^N p_i x_i\right) \leq \sum_{i=1}^N p_i f(x_i).$$

Similarly for a concave  function: $\mathbb{E}[f(X)] \leq f(\mathbb{E}[X])$.

Jensen's Inequality for Convex Functions

Proof by Induction

Assignment Project Exam Help

- ▶ Two-state random variable $X \in \{x_1, x_2\}$

▶ <https://eduassistpro.github.io>

- ▶ $0 \leq p \leq 1$

we simply follow the definition of convexity:

$$\underbrace{p_1 f(x_1) + p_2 f(x_2)}_{\mathbb{E}[f(X)]} \geq \underbrace{p_1 p_1 + p_2 p_2}_{\mathbb{E}[X]}$$

Jensen's Inequality for Convex Functions

Proof by Induction — Cont'd

(2) $(K - 1) \rightarrow K$: Assuming the theorem is true for distributions with $K - 1$ states, and writing: $p'_i = p_i / (1 - p_K)$ for $i = 1, \dots, K - 1$:

$$p f(x) = p f(x) + (1 - p) p' f(x)$$

<https://eduassistpro.github.io> on hypothesis

Add WeChat $\sum_{i=1}^K p_i x_i$ edu_assist_pro

$$\sum_{i=1}^K p_i f(x_i) \geq f\left(\sum_{i=1}^K p_i x_i\right) \Rightarrow \mathbb{E}[f(X)] \geq f(\mathbb{E}[x]) \quad \text{equality case?}$$

Jensen's Inequality Example: The AM-GM Inequality

Recall that for a **concave** \cap function: $\mathbb{E}[f(X)] \leq f(\mathbb{E}[X])$.

Consider $X \in \{x_1, \dots, x_N\}$, $X \geq 0$ with uniform probability distribution

$\mathbf{p} = (\frac{1}{N}, \dots, \frac{1}{N})$ and the strictly concave \cap function $f(x) = \log x$

$$\frac{1}{N} \quad \frac{1}{N}$$

<https://eduassistpro.github.io>

$$\log \frac{1}{N} \leq \log \frac{1}{N}$$

Add WeChat edu_assist_pro

$$\left(\prod_{i=1}^N x_i \right)^{\frac{1}{N}} \leq \frac{1}{N} \sum_{i=1}^N x_i$$
$$\sqrt[N]{x_1 x_2 \dots x_N} \leq \frac{x_1 + x_2 \dots + x_N}{N}$$

1 Chain Rule for Mutual Information

2 Convex Functions

3 Jensen's Inequality

4 Gibb

5 Information Cannot Hurt

6 Data Processing Inequality

7 Wrapping Up

Assignment Project Exam Help

<https://eduassistpro.github.io>

Add WeChat edu_assist_pr

Assignment Project Exam Help

Theorem

The relative
entropy of $p(X)$
and $q(X)$

$p(X)$

<https://eduassistpro.github.io>

$$D_{\text{KL}}(p \parallel q) \geq 0$$

with equality if and only if $p(x) = q(x)$ for all x

Add WeChat edu_assist_pro

Gibbs' Inequality

Proof (1 of 2)

Recall that: $D_{\text{KL}}(p||q) = \sum_{x \in \mathcal{X}} p(x) \log \frac{p(x)}{q(x)} = \mathbb{E}_{p(X)} \left[\log \frac{p(X)}{q(X)} \right]$

Let $\mathcal{A} = \{x : p(x) > 0\}$. Then:

<https://eduassistpro.github.io>

Add WeChat edu_assist_pro

$$\begin{aligned} &= \log \sum_{x \in \mathcal{A}} q(x) \\ &\leq \log \sum_{x \in \mathcal{X}} q(x) \\ &= \log 1 \\ &= 0 \end{aligned}$$

Gibbs' Inequality

Proof (2 of 2)

Since $\log u$ is strictly convex we have equality if $\frac{p(x)}{q(x)} = c$ for all x . Then:

Also, the la

$\sum_{x \in \mathcal{X}} q(x) = \sum_{x \in \mathcal{X}} p(x)$

Therefore $c = 1$ and $D_{\text{KL}}(p||q) = 0 \Leftrightarrow p(x) = q(x)$ for all x .

Alternative proof: Use the fact that $\log x \leq x - 1$.

Non-Negativity of Mutual Information

Corollary

For any two random variables X, Y :

with equality

Proof: We simply use the definition of mutual information inequality:

$$I(X; Y) = D_{\text{KL}}(p(X, Y) \parallel$$

with equality if and only if $p(X, Y) = p(X)p(Y)$, i.e. X and Y are independent.

1 Chain Rule for Mutual Information

2 Convex Functions

3 Jensen's Inequality

4 Gibbs

5 Information Cannot Hurt

6 Data Processing Inequality

7 Wrapping Up

Assignment Project Exam Help

<https://eduassistpro.github.io>

Add WeChat edu_assist_pr

Conditioning Reduces Entropy

Information Cannot Hurt — Proof

Theorem

For any two random variables X, Y ,

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

with equality

Proof: We

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

with equality if and only if $p(X, Y) = p(X)p(Y)$, i.e X and Y are independent.

Data are helpful, they don't increase uncertainty on average.

Conditioning Reduces Entropy

Information Cannot Hurt — Example (from Cover & Thomas, 2006)

Let X, Y have the following joint distribution:

$p(X, Y)$		$p(X)$ = (1/8, 7/8)	
		1	2
1	0	3/4	
2	1/2	1/2	

$H(X|Y=1) = 1 \text{ bit}$

We see that in this case $H(X|Y=1) < H(X)$.

However, $H(X|Y) = \sum_{y \in \{1,2\}} p(y) H(X|Y=y) = \frac{3}{4} H(X|Y=1) + \frac{1}{2} H(X|Y=2)$

$H(X|Y = y_k)$ may be greater than $H(X)$ but the average: $H(X|Y)$ is always less or equal to $H(X)$.

Information cannot hurt on average

1 Chain Rule for Mutual Information

2 Convex Functions

3 Jensen's Inequality

4 Gibb

5 Information Cannot Hurt

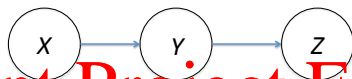
6 Data Processing Inequality

7 Wrapping Up

Assignment Project Exam Help

<https://eduassistpro.github.io>

Add WeChat edu_assist_pr



Assignment Project Exam Help

Definiti

Random
(denote
written as

$$p(X, Y, Z) = p(X)p(Y|X)p(Z|Y)$$

Consequences:

- $X \rightarrow Y \rightarrow Z$ if and only if X and Z are conditionally independent given Y .
- $X \rightarrow Y \rightarrow Z$ implies that $Z \rightarrow Y \rightarrow X$.
- If $Z = f(Y)$, then $X \rightarrow Y \rightarrow Z$

Data-Processing Inequality

Definition

Theorem

if $X \rightarrow Y \rightarrow Z$ then: $I(X; Y) \geq I(X; Z)$

- X is the source of the information, Y is the processed data, and Z is the destination. The data processing inequality states that the mutual information between X and Z is less than or equal to the mutual information between X and Y .
- No “clever” manipulation of the data can improve the information that Y contains about X .
- No processing of Y , deterministic or random, can increase the information that Y contains about X .

Data-Processing Inequality

Proof

Recall that the chain rule for mutual information states that:

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

Therefore

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

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

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

Data-Processing Inequality

Functions of the Data

Assignment Project Exam Help

Corollary

In particular

<https://eduassistpro.github.io>

Proof: $X \rightarrow Y \rightarrow g(Y)$ forms a Markov chain.

Add WeChat [edu_assist_pro](#)

Functions of the data Y cannot increase

Data-Processing Inequality

Observation of a “Downstream” Variable

Corollary

If $X \rightarrow Y \rightarrow Z$ then $I(X; Y|Z) \leq I(X; Y)$

Proof: We use again the chain rule for mutual information:

<https://eduassistpro.github.io>

Therefore:

$$I(X; Y) + \underbrace{I(X; Z|Y)}_0 = I(X; Z) + I(X; Y|Z)$$

$$I(X; Y|Z) = I(X; Y) - I(X; Z) \quad \text{but } I(X; Z) \geq 0$$

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

The dependence between X and Y cannot be increased by the observation of a “downstream” variable.

1 Chain Rule for Mutual Information

2 Convex Functions

3 Jensen's Inequality

4 Gibbs

5 Information Cannot Hurt

6 Data Processing Inequality

7 Wrapping Up

Assignment Project Exam Help

<https://eduassistpro.github.io>

Add WeChat edu_assist_pr

Summary & Conclusions

- Chain rule for mutual information

- Convex Functions

- Jensen's Inequality

- Important inequalities regarding information processing

- **Reading:** Mackay §2.6 to §2.10, Cover & Thomas §2.5 to §2.8

Next time

Assignment Project Exam Help

- Law of large numbers

- Mar

<https://eduassistpro.github.io>

Add WeChat edu_assist_pr

- Chebychev's inequality