COMP2610 / COMP6261 - Information Theory

Lecture 8: Some Fundamental Inequalities

Australian
National
University

14 August 2018

Assignment Project Exam Help

- Decomposability of entropy

- Rel https://eduassistpro.github.i

Add WeChat edu_assist_pr

- Mutual information

# Review

Relative entropy (KL divergence):

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

Mutual inf

$$= H(X) - H($$

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

Conditional mutual information of $X, Y$ given $Z$:

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

Mutual information chain rule

Jensen'

"Informa

Data processing inequality

# Outline

Assignment Project Exam Help

https://eduassistpro.github.i

Add WeChat edu_assist_pr

Assignment Project Exam Help

https://eduassistpro.github.i

Add WeChat edu_assist_pr

# Recall: Joint Mutual Information

Recall the mutual information between $X$ and $Y$:

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

We can als ... $X_N$ and $Y_1, \ldots,$

$$I(\ _1 \ _N \ _1 \ _M \qquad _1 \ _N \ _1 \qquad Y_M)$$

$$- H(X$$

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 $Z$

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

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

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

$$\underbrace{\phantom{I(X; Y, Z)}}_{I(Y;X)}$$

$$I(X; Y, Z) = I(X; Y) + I(X; Z|Y) \quad \text{definition}$$

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 $X_1, \ldots, X_N$ and $Y$:

$$I(X_1, \ldots, X_N; Y) = I(X_1; Y) + I(X_2, \ldots, X_N; Y \mid X_1)$$

$$\cdots \mid X_1, X_2)$$

$$= \sum_{i=1}^{N} I(X_i; Y \mid X_1, \ldots, )$$

$$= \sum_{i=1}^{N} I(Y; X_i \mid X_1, \ldots, {}_{i-1}).$$

$$\lambda f(x_1) + (1 - \lambda) f(x_2)$$

$0 \leq \lambda \leq 1$   (Figure from Mackay, 2003)

A function is convex ⌣ if every chord of the function lies above the function

**Definition**
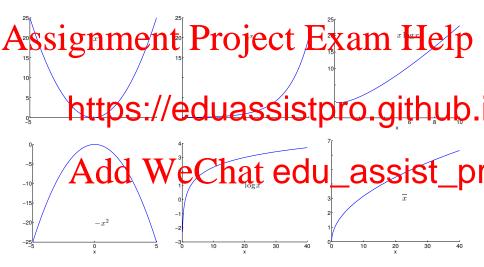
A function ... nd
$0 \leq \lambda \leq$ ...

We say $f$ is strictly convex $\smile$ if for all $x_1, x$ ... nly
for $\lambda = 0$ and $\lambda = 1$.

Similarly, a function $f$ is concave $\frown$ if $-f$ i ...
the function lies below the function.

# 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 $\smile$ (strictly convex $\smile$) over that interval.

This allow

Example

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

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

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

# Convexity, Concavity and Optimization

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

$$\frac{d f}{dx} = 0,$$

then $f(x$

Note: the q  zed at some $x$, it i

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

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

- Similarly for minimisation of convex functions

Assignment Project Exam Help

https://eduassistpro.github.i

Add WeChat edu_assist_pr

# Jensen's Inequality for Convex Functions

## Theorem: Jensen's Inequality

If $f$ is a convex $\smile$ 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 **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 $\frown$ function: $\mathbb{E}[f(X)] \leq f(\mathbb{E}[X])$ .

(1): $K = 2$

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

-

- $0 \le p \le 1$

we simply follow the definition of convexity:

$$\underbrace{p_1 f(x_1) + p_2 f(x_2)}_{\mathbb{E}[f(X)]} \ge \underbrace{\phantom{1 1 \quad 2 2}}_{\mathbb{E}[X]}$$

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

$$\sum_{i=1}^{K} p_i f(x_i) = p_K f(x_K) + (1 - p_K) \sum_{i=1}^{K-1} p'_i f(x_i)$$

induction hypothesis

$$\geq p_K f(x_K) + (1 - p_K) f\left(\sum_{i=1}^{K-1} p'_i x_i\right)$$

$$\underbrace{\phantom{p_K x_K + (1-p_K)\sum}}_{\sum_{i=1}^{K} p_i x_i}$$

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

# Jensen's Inequality Example: The AM-GM Inequality

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

Consider $X \in \{x_1, \ldots, x_N\}$, $X \geq 0$ with uniform probability distribution
$p = (\frac{1}{N}, \ldots, \frac{1}{N})$ and the strictly concave $\frown$ function $f(x) = \log x$

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

$$\log \left( \prod_{i=1}^{N} x_i \right)^{\frac{1}{N}} \leq \log \frac{1}{N} \sum_{i=1}^{N} x_i$$

$$\left( \prod_{i=1}^{N} x_i \right)^{\frac{1}{N}} \leq \frac{1}{N} \sum_{i=1}^{N} x_i$$

$$\sqrt[N]{x_1 x_2 \ldots x_N} \leq \frac{x_1 + x_2 \ldots + x_N}{N}$$

Assignment Project Exam Help

https://eduassistpro.github.i

Add WeChat edu_assist_pr

### Theorem

The relati $p(X)$
and $q(X$

$$\mathrm{KL}(\ \|\ ) \geq$$

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

# Gibbs' Inequality

## Proof (1 of 2)

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

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

$$= \log \sum_{x\in\mathcal{A}} q(x)$$

$$\leq \log \sum_{x\in\mathcal{X}} q(x)$$

$$= \log 1$$

$$= 0$$

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

Also, the la

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

Therefore $c = 1$ and $D_{\mathsf{KL}}(p \| q) = 0 \Leftrightarrow p(x) = q(x)$ 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 equal

**Proof**: We simply use the definition of mutual inform
inequality:

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

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

Information Cannot Hurt — Proof

## Theorem

For any two random variables $X, Y$:

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

with equal

Proof: We

$$I(X; Y) = H(X) - H(X \mid Y)$$
$$H(X \mid Y)$$

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:

$\mathbf{p}(X) = (1/8, 7/8)$

$\mathbf{p}(X|Y = 1) = (0, 1)$

$1/2, 1/2$

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

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

$H(X|Y = y_k)$ 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

Assignment Project Exam Help

https://eduassistpro.github.i

Add WeChat edu_assist_pr

# Markov Chain



X → Y → Z

## Definiti

Random
(denote
written as

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

Consequence

- $X \rightarrow Y \rightarrow Z$ 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 t                                                                s the
  pro

- No "clever" manipulation of the data can impro
  can be made from the data

- 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:

$$I(X; Y, Z) = I(X; Y) + I(X; Z \mid Y)$$

Therefo

$$I(X; Y) + I(X; Z \mid Y) = I(X; Z) + I$$

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

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

**Corollary**

*In particul

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

Functions of the data $Y$ cannot increa

# 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:

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.

Assignment Project Exam Help

https://eduassistpro.github.i

Add WeChat edu_assist_pr

# Summary & Conclusions

- Chain rule for mutual information.

- Convex Functions

- Jen

- Important inequalities regarding informati processing

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

# Next time

- Law of large numbers

- Mar

- Chebychev's inequality