

# Lecture 3: Information Inequalities

Tianpei Xie

Jan. 6th., 2023

## Contents

<b>1</b>	<b>Information Theory Basics</b>	<b>2</b>
1.1	Entropy, Relative Entropy, and Mutual Information . . . . .	2
1.2	Chain Rules for Entropy, Relative Entropy, and Mutual Information . . . . .	4
1.3	Log-Sum Inequalities and Convexity . . . . .	5
1.4	Data Processing Inequality . . . . .	5
<b>2</b>	<b>Information Inequalities</b>	<b>6</b>
2.1	Han's Inequality . . . . .	6
2.2	Applications of Han's Inequality . . . . .	8
2.2.1	Combinatorial Entropies . . . . .	8
2.2.2	Edge Isoperimetric Inequality on the Binary Hypercube . . . . .	8
2.3	$\Phi$ -Entropy . . . . .	8
2.4	Sub-Additivity of $\Phi$ -Entropy . . . . .	8
2.5	Duality and Variational Formulas . . . . .	9
2.6	Optimal Transport . . . . .	9
2.7	Pinsker's Inequality . . . . .	9
2.8	Birgé's Inequality . . . . .	9
2.9	The Brunn-Minkowski Inequality . . . . .	9

# 1 Information Theory Basics

## 1.1 Entropy, Relative Entropy, and Mutual Information

- **Definition (*Shannon Entropy*)** [Cover and Thomas, 2006]

Let  $(\Omega, \mathcal{F}, \mathbb{P})$  be a probability space and  $X : \mathbb{R} \rightarrow \mathcal{X}$  be a random variable. Define  $p(x)$  as the probability density function of  $X$  with respect to a base measure  $\mu$  on  $\mathcal{X}$ . **The Shannon Entropy** is defined as

$$\begin{aligned} H(X) &:= \mathbb{E}_p [-\log p(X)] \\ &= \int_{\Omega} -\log p(X(\omega)) d\mathbb{P}(\omega) \\ &= - \int_{\mathcal{X}} p(x) \log p(x) d\mu(x) \end{aligned}$$

- **Definition (*Conditional Entropy*)** [Cover and Thomas, 2006]

If a pair of random variables  $(X, Y)$  follows the joint probability density function  $p(x, y)$  with respect to a base product measure  $\mu$  on  $\mathcal{X} \times \mathcal{Y}$ . Then **the joint entropy** of  $(X, Y)$ , denoted as  $H(X, Y)$ , is defined as

$$H(X, Y) := \mathbb{E}_{X, Y} [-\log p(X, Y)] = - \int_{\mathcal{X} \times \mathcal{Y}} p(x, y) \log p(x, y) d\mu(x, y)$$

Then **the conditional entropy**  $H(Y|X)$  is defined as

$$\begin{aligned} H(Y|X) &:= \mathbb{E}_{X, Y} [-\log p(Y|X)] = - \int_{\mathcal{X} \times \mathcal{Y}} p(x, y) \log p(y|x) d\mu(x, y) \\ &= \mathbb{E}_X [\mathbb{E}_Y [-\log p(Y|X)]] = \int_{\mathcal{X}} p(x) \left( - \int_{\mathcal{Y}} p(y|x) \log p(y|x) d\mu(y) \right) d\mu(x) \end{aligned}$$

- **Proposition 1.1 (*Properties of Shannon Entropy*)** [Cover and Thomas, 2006]

Let  $X, Y, Z$  be random variables.

1. (**Non-negativity**)  $H(X) \geq 0$ ;
2. (**Chain Rule**)

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

Furthermore,

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

3. (**Sub-Additivity**)

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

4. (**Concavity**)  $H(p) := \mathbb{E}_p [-\log p(X)]$  is a concave function in terms of p.d.f.  $p$ , i.e.

$$H(\lambda p_1 + (1 - \lambda)p_2) \geq \lambda H(p_1) + (1 - \lambda)H(p_2)$$

for any two p.d.fs  $p_1, p_2$  on  $\mathcal{X}$  and any  $\lambda \in [0, 1]$ .

- **Definition (*Relative Entropy / Kullback-Leibler Divergence*)** [Cover and Thomas, 2006]

Suppose that  $P$  and  $Q$  are *probability measures* on a measurable space  $\mathcal{X}$ , and  $P$  is *absolutely continuous* with respect to  $Q$ , then the relative entropy or the Kullback-Leibler divergence is defined as

$$\mathbb{KL}(P \parallel Q) := \mathbb{E}_P \left[ \log \left( \frac{dP}{dQ} \right) \right] = \int_{\mathcal{X}} \log \left( \frac{dP(x)}{dQ(x)} \right) dP(x)$$

where  $\frac{dP}{dQ}$  is the *Radon-Nikodym derivative* of  $P$  with respect to  $Q$ . Equivalently, the KL-divergence can be written as

$$\mathbb{KL}(P \parallel Q) = \int_{\mathcal{X}} \left( \frac{dP(x)}{dQ(x)} \right) \log \left( \frac{dP(x)}{dQ(x)} \right) dQ(x)$$

which is *the entropy of  $P$  relative to  $Q$* . Furthermore, if  $\mu$  is a base measure on  $\mathcal{X}$  for which densities  $p$  and  $q$  with  $dP = p(x)d\mu$  and  $dQ = q(x)d\mu$  exist, then

$$\mathbb{KL}(P \parallel Q) = \int_{\mathcal{X}} p(x) \log \left( \frac{p(x)}{q(x)} \right) d\mu(x)$$

- **Definition (*Mutual Information*)** [Cover and Thomas, 2006]

Consider two random variables  $X, Y$  on  $\mathcal{X} \times \mathcal{Y}$  with joint probability distribution  $P_{(X,Y)}$  and marginal distribution  $P_X$  and  $P_Y$ . The mutual information  $I(X; Y)$  is *the relative entropy* between *the joint distribution  $P_{(X,Y)}$*  and *the product distribution  $P_X \otimes P_Y$* :

$$I(X; Y) = \mathbb{KL}(P_{(X,Y)} \parallel P_X \otimes P_Y) = \mathbb{E}_{P_{(X,Y)}} \left[ \log \frac{dP_{(X,Y)}}{dP_X \otimes dP_Y} \right]$$

If  $P_{(X,Y)}$  has a probability density function  $p(x, y)$  with respect to a base measure  $\mu$  on  $\mathcal{X} \times \mathcal{Y}$ , then

$$I(X; Y) = \int_{\mathcal{X} \times \mathcal{Y}} p(x, y) \log \left( \frac{p(x, y)}{p_X(x)p_Y(y)} \right) d\mu(x, y)$$

- **Proposition 1.2 (*Properties of Relative Entropy and Mutual Information*)** [Cover and Thomas, 2006]

Let  $X, Y$  be random variables.

1. (**Non-negativity**) Let  $p(x), q(x)$  be probability density function of  $P, Q$ .

$$\mathbb{KL}(P \parallel Q) \geq 0$$

with equality if and only if  $p(x) = q(x)$  almost surely. Therefore, the mutual information is non-negative as well:

$$I(X; Y) \geq 0$$

with equality if and only if  $X$  and  $Y$  are independent.

2. (**Finite Cardinality Domain**) Let  $|\mathcal{X}|$  be the number of elements in domain  $\mathcal{X}$  and  $X$  is a discrete random variables in  $\mathcal{X}$ . Then the relative entropy of probability distribution  $p$  with respect to uniform distribution  $u$  on  $\mathcal{X}$  is

$$\begin{aligned}\text{KL}(p \parallel u) &= \log |\mathcal{X}| - H(X) \geq 0 \\ \Rightarrow H(X) &\leq \log |\mathcal{X}|\end{aligned}$$

3. (**Symmetry**)  $I(X; Y) = I(Y; X)$
4. (**Information Gain via Conditioning**) The mutual information  $I(X; Y)$  is the reduction in the uncertainty of  $X$  due to the knowledge of  $Y$  (and vice versa)

$$\begin{aligned}I(X; Y) &= H(X) - H(X|Y) \\ &= H(Y) - H(Y|X) \\ &= H(X) + H(Y) - H(X, Y)\end{aligned}\tag{1}$$

5. (**Shannon Entropy as Self-Information**)  $I(X; X) = H(X)$

## 1.2 Chain Rules for Entropy, Relative Entropy, and Mutual Information

- **Proposition 1.3 (Conditioning Reduces Entropy)** [Cover and Thomas, 2006]  
From non-negativity of mutual information, we see that the entropy of  $X$  is non-increasing when conditioning on  $Y$

$$H(X|Y) \leq H(X)\tag{2}$$

where equality holds if and only if  $X$  and  $Y$  are independent.

- **Proposition 1.4 (Chain Rule for Entropy)** [Cover and Thomas, 2006]  
Let  $X_1, X_2, \dots, X_n$  be drawn according to  $p(x_1, x_2, \dots, x_n)$ . Then

$$H(X_1, X_2, \dots, X_n) = \sum_{i=1}^n H(X_i | X_{i-1}, \dots, X_1)\tag{3}$$

- **Proposition 1.5 (Sub-Additivity of Entropy)** [Cover and Thomas, 2006]  
Let  $X_1, X_2, \dots, X_n$  be drawn according to  $p(x_1, x_2, \dots, x_n)$ . Then

$$H(X_1, X_2, \dots, X_n) \leq \sum_{i=1}^n H(X_i)\tag{4}$$

with equality if and only if the  $X_i$  are independent.

- **Proposition 1.6 (Chain Rule for Mutual Information)** [Cover and Thomas, 2006]  
Let  $X_1, X_2, \dots, X_n, Y$  be drawn according to  $p(x_1, x_2, \dots, x_n, y)$ . Then

$$I(X_1, X_2, \dots, X_n; Y) = \sum_{i=1}^n H(X_i; Y | X_{i-1}, \dots, X_1)\tag{5}$$

where **the conditional mutual information** is defined as

$$I(X; Y|Z) := H(X|Z) - H(X|Y, Z) = \text{KL}(P_{(X,Y|Z)} \parallel P_{X|Z} \otimes P_{Y|Z})$$

- **Proposition 1.7 (Chain Rule for Relative Entropy)** [Cover and Thomas, 2006]  
Let  $P_{(X,Y)}$  and  $Q_{(X,Y)}$  be two probability measures on product space  $\mathcal{X} \times \mathcal{Y}$  and  $P \ll Q$ . Denote the marginal distributions  $P_X, Q_X$  and  $P_Y, Q_Y$  on  $\mathcal{X}$  and  $\mathcal{Y}$ , respectively.  $P_{Y|X}$  and  $Q_{Y|X}$  are conditional distributions (Note that  $P_{Y|X} \ll Q_{Y|X}$ ). Define **the conditional relative entropy** as

$$\mathbb{E}_X [\text{KL}(P_{Y|X} \parallel Q_{Y|X})] := \mathbb{E}_X \left[ \mathbb{E}_{P_{Y|X}} \left[ \log \left( \frac{dP_{Y|X}}{dQ_{Y|X}} \right) \right] \right].$$

Then the relative entropy of joint distribution  $P_{(X,Y)}$  with respect to  $Q_{(X,Y)}$  is

$$\text{KL}(P_{(X,Y)} \parallel Q_{(X,Y)}) = \text{KL}(P_X \parallel Q_X) + \mathbb{E}_X [\text{KL}(P_{Y|X} \parallel Q_{Y|X})] \quad (6)$$

In addition, let  $P$  and  $Q$  denote two joint distributions for  $X_1, X_2, \dots, X_n$ , let  $P_{1:i}$  and  $Q_{1:i}$  denote the marginal distributions of  $X_1, X_2, \dots, X_i$  under  $P$  and  $Q$ , respectively. Let  $P_{X_i|1\dots i-1}$  and  $Q_{X_i|1\dots i-1}$  denote the conditional distribution of  $X_i$  with respect to  $X_1, X_2, \dots, X_{i-1}$  under  $P$  and under  $Q$ .

$$\text{KL}(P \parallel Q) = \sum_{i=1}^n \mathbb{E}_{P_{1:i-1}} [\text{KL}(P_{X_i|1\dots i-1} \parallel Q_{X_i|1\dots i-1})] \quad (7)$$

### 1.3 Log-Sum Inequalities and Convexity

- **Proposition 1.8 (Log-Sum Inequalities)** [Cover and Thomas, 2006]  
For non-negative numbers  $a_1, \dots, a_n$  and  $b_1, \dots, b_n$ ,

$$\sum_{i=1}^n a_i \log \frac{a_i}{b_i} \geq \left( \sum_{i=1}^n a_i \right) \log \frac{\sum_{i=1}^n a_i}{\sum_{i=1}^n b_i} \quad (8)$$

with equality if and only if  $\frac{a_i}{b_i}$  is constant.

- **Proposition 1.9 (Joint Convexity of Relative Entropy)** [Cover and Thomas, 2006]  
 $\text{KL}(p \parallel q)$  is **convex** in the pair  $(p, q)$ ; that is, if  $(p_1, q_1)$  and  $(p_2, q_2)$  are two pairs of probability density functions, then for  $\lambda \in [0, 1]$ ,

$$\text{KL}(\lambda p_1 + (1 - \lambda)p_2 \parallel \lambda q_1 + (1 - \lambda)q_2) \leq \lambda \text{KL}(p_1 \parallel q_1) + (1 - \lambda) \text{KL}(p_2 \parallel q_2) \quad (9)$$

- **Proposition 1.10** [Cover and Thomas, 2006]  
Let  $(X, Y) \sim p(x, y) = p(x)p(y|x)$ . The mutual information  $I(X; Y)$  is a **concave** function of  $p(x)$  for fixed  $p(y|x)$  and a **convex** function of  $p(y|x)$  for fixed  $p(x)$ .

### 1.4 Data Processing Inequality

- **Definition (Data Processing Markov Chain)**  
Random variables  $X, Y, Z$  are said to **form a Markov chain** in that order (denoted by  $X \rightarrow Y \rightarrow Z$ ) if the conditional distribution of  $Z$  depends only on  $Y$  and is **conditionally independent** of  $X$ . Specifically,  $X, Y$ , and  $Z$  form a Markov chain  $X \rightarrow Y \rightarrow Z$  if the joint probability mass function can be written as

$$p(x, y, z) = p(x)p(y|x)p(z|y)$$

- **Proposition 1.11** (*Data Processing Inequality*) [Cover and Thomas, 2006]  
If  $X \rightarrow Y \rightarrow Z$ , then

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

- **Corollary 1.12** [Cover and Thomas, 2006]  
In particular, if  $Z = g(Y)$ , we have

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

- **Corollary 1.13** [Cover and Thomas, 2006]  
If  $X \rightarrow Y \rightarrow Z$ , then

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

Thus, the dependence of  $X$  and  $Y$  is **decreased** (or remains unchanged) by the observation of a “**downstream**” random variable  $Z$ .

## 2 Information Inequalities

### 2.1 Han’s Inequality

- **Proposition 2.1** (*Han’s Inequality*) [Cover and Thomas, 2006, Boucheron et al., 2013]  
Let  $X_1, X_2, \dots, X_n$  be random variables. Then

$$H(X_1, X_2, \dots, X_n) \leq \frac{1}{n-1} \sum_{i=1}^n H(X_1, \dots, X_{i-1}, X_{i+1}, \dots, X_n) \quad (10)$$

**Proof:** For any  $i = 1, \dots, n$ , by the definition of the conditional entropy and the fact that conditioning reduces entropy,

$$\begin{aligned} H(X_1, X_2, \dots, X_n) &= H(X_1, \dots, X_{i-1}, X_{i+1}, \dots, X_n) + H(X_i | X_1, \dots, X_{i-1}, X_{i+1}, \dots, X_n) \\ &\leq H(X_1, \dots, X_{i-1}, X_{i+1}, \dots, X_n) + H(X_i | X_1, \dots, X_{i-1}). \end{aligned}$$

Summing these  $n$  inequalities and using the chain rule for entropy, we get

$$\begin{aligned} nH(X_1, X_2, \dots, X_n) &\leq \sum_{i=1}^n H(X_1, \dots, X_{i-1}, X_{i+1}, \dots, X_n) + \sum_{i=1}^n H(X_i | X_1, \dots, X_{i-1}) \\ &= \sum_{i=1}^n H(X_1, \dots, X_{i-1}, X_{i+1}, \dots, X_n) + H(X_1, X_2, \dots, X_n) \end{aligned}$$

which is what we wanted to prove. ■

- **Proposition 2.2** (*Han’s Inequality for Relative Entropy*) [Boucheron et al., 2013]  
Let  $(\mathcal{X}, \mathcal{B})$  be a measurable space, and  $P$  and  $Q$  be probability measures on  $\mathcal{X}^n$  such that  $P = P_1 \otimes \dots \otimes P_n$  is a **product measure**. We denote the element of  $\mathcal{X}^n$  by  $x = (x_1, \dots, x_n)$  and write  $x_{(-i)} := (x_1, \dots, x_{i-1}, x_{i+1}, \dots, x_n)$  for the  $(n-1)$ -vector obtained by **leaving out the  $i$ -th component of  $x$**  (i.e. the  $i$ -th Jackknife sample of  $x$ ). Denote  $Q_{(-i)}$  and  $P_{(-i)}$  the

marginal distributions of  $Q$  and  $P$ . Let  $p_{(-i)}$  and  $q_{(-i)}$  denote the corresponding probability density function with respect to base measure  $\mu$  on  $\mathcal{X}$ .

$$\begin{aligned} q_{(-i)}(x_{(-i)}) &= \int_{y \in \mathcal{X}} q(x_1, \dots, x_{i-1}, y, x_{i+1}, \dots, x_n) d\mu(y) \\ p_{(-i)}(x_{(-i)}) &= \int_{y \in \mathcal{X}} p(x_1, \dots, x_{i-1}, y, x_{i+1}, \dots, x_n) d\mu(y) \\ &= \prod_{j \neq i} p_j(x_j). \end{aligned}$$

Then

$$\mathbb{KL}(Q \parallel P) \geq \frac{1}{n-1} \sum_{i=1}^n \mathbb{KL}(Q_{(-i)} \parallel P_{(-i)}) \quad (11)$$

or equivalently,

$$\mathbb{KL}(Q \parallel P) \leq \sum_{i=1}^n (\mathbb{KL}(Q \parallel P) - \mathbb{KL}(Q_{(-i)} \parallel P_{(-i)})) \quad (12)$$

**Proof:** From Han's inequality, we have

$$-H(Q) \geq -\frac{1}{n-1} \sum_{i=1}^n H(Q_{(-i)}).$$

Since

$$\mathbb{KL}(Q \parallel P) = -H(Q) + \mathbb{E}_Q[-\log P(X)]$$

and

$$\mathbb{KL}(Q_{(-i)} \parallel P_{(-i)}) = -H(Q_{(-i)}) + \mathbb{E}_{Q_{(-i)}}[-\log P_{(-i)}(X_{(-i)})],$$

it suffices to show that

$$\mathbb{E}_Q[-\log P(X)] = \frac{1}{n-1} \sum_{i=1}^n \mathbb{E}_{Q_{(-i)}}[-\log P_{(-i)}(X_{(-i)})].$$

This may be seen easily by noting that by the product property of  $P$ , we have  $p(x) = p_{(-i)}(x_{(-i)})p_i(x_i)$  for all  $i$ , and also  $p(x) = \prod_i p_i(x_i)$ , and therefore

$$\begin{aligned} \mathbb{E}_Q[-\log P(X)] &= \frac{1}{n} \sum_{i=1}^n \mathbb{E}_Q[-\log P_{(-i)}(X_{(-i)}) - \log P_i(X_i)] \\ &= \frac{1}{n} \sum_{i=1}^n \mathbb{E}_Q[-\log P_{(-i)}(X_{(-i)})] + \frac{1}{n} \sum_{i=1}^n \mathbb{E}_Q[-\log P_i(X_i)] \\ &= \frac{1}{n} \sum_{i=1}^n \mathbb{E}_Q[-\log P_{(-i)}(X_{(-i)})] + \frac{1}{n} \mathbb{E}_Q[-\log P(X)]. \end{aligned}$$

Rearranging, we obtain

$$\begin{aligned} \mathbb{E}_Q[-\log P(X)] &= \frac{1}{n-1} \sum_{i=1}^n \mathbb{E}_Q[-\log P_{(-i)}(X_{(-i)})] \\ &= \frac{1}{n-1} \sum_{i=1}^n \mathbb{E}_{Q_{(-i)}}[-\log P_{(-i)}(X_{(-i)})]. \quad \blacksquare \end{aligned}$$

## 2.2 Applications of Han's Inequality

### 2.2.1 Combinatorial Entropies

### 2.2.2 Edge Isoperimetric Inequality on the Binary Hypercube

## 2.3 $\Phi$ -Entropy

- **Definition ( $\Phi$ -Entropy)** [Boucheron et al., 2013]

Let  $\Phi : [0, \infty) \rightarrow \mathbb{R}$  be a **convex** function, and assign, to every **non-negative integrable random variable**  $X$ , the  $\Phi$ -entropy of  $X$  is defined as

$$H_\Phi(X) = \mathbb{E} [\Phi(X)] - \Phi(\mathbb{E} [X]). \quad (13)$$

- **Remark** By Jensen's inequality, the  $\Phi$ -entropy is *non-negative*

$$\begin{aligned} \Phi(\mathbb{E} [X]) &\leq \mathbb{E} [\Phi(X)] \\ \Rightarrow H_\Phi(X) &= \mathbb{E} [\Phi(X)] - \Phi(\mathbb{E} [X]) \geq 0. \end{aligned}$$

- **Example (*Special Examples for  $\Phi$ -Entropy*)**

1. For  $\Phi(x) = x^2$ , the  $\Phi$ -entropy of  $X$  is the **variance** of  $X$ :

$$H_\Phi(X) = \mathbb{E} [X^2] - (\mathbb{E} [X])^2 = \text{Var}(X).$$

2. For  $\Phi(x) = x \log x$ , the  $\Phi$ -entropy of  $X$  is defined as the entropy of  $X$

$$\text{Ent}(X) := \mathbb{E} [X \log X] - \mathbb{E} [X] \log (\mathbb{E} [X]). \quad (14)$$

Let  $(\Omega, \mathcal{B})$  be measurable space, and  $P$  and  $Q$  are probability measures on  $\Omega$  with  $P \ll Q$ . Define a random variable  $X$  by the *Radon-Nikodym derivative* of  $P$  with respect to  $Q$ ; that is,

$$X(\omega) := \begin{cases} \frac{dP}{dQ}(\omega) & Q(\omega) > 0 \\ 0 & \text{o.w.} \end{cases}.$$

We see that  $X$  is  $Q$ -measurable and  $dP = X dQ$  with  $\mathbb{E}_Q [X] = 1$ . Then the entropy of  $X$  is the relative entropy of  $P$  with respect to  $Q$ .

$$\text{Ent}(X) = \text{KL}(P \parallel Q) \quad (15)$$

## 2.4 Sub-Additivity of $\Phi$ -Entropy

- **Remark (*Sub-Additivity of Shannon Entropy*)**

Let  $X_1, X_2, \dots, X_n$  be drawn according to  $p(x_1, x_2, \dots, x_n)$ . Then

$$H(X_1, X_2, \dots, X_n) \leq \sum_{i=1}^n H(X_i)$$

with equality if and only if the  $X_i$  are independent.



• **Proposition 2.3 (Sub-Additivity of The Entropy)** [Boucheron et al., 2013]

Let  $\Phi(x) = x \log x$ , for  $x > 0$  and  $\Phi(0) = 0$ . Let  $Z_1, Z_2, \dots, Z_n$  be independent random variables taking values in  $\mathcal{X}$ , and let  $f : \mathcal{X}^n \rightarrow [0, \infty)$ . Letting  $X = f(Z_1, Z_2, \dots, Z_n)$ , we have

$$\mathbb{E} [\Phi(X)] - \Phi(\mathbb{E} [X]) \leq \sum_{i=1}^n \mathbb{E} [\mathbb{E}_{(-i)} [\Phi(X)] - \Phi(\mathbb{E}_{(-i)} [X])], \quad (16)$$

where  $\mathbb{E}_{(-i)} [\cdot]$  is the conditional expectation operator conditioning on  $Z_{(-i)}$ . Introducing the notation  $\text{Ent}_{(-i)}(X) = \mathbb{E}_{(-i)} [\Phi(X)] - \Phi(\mathbb{E}_{(-i)} [X])$ , this can be re-written as

$$\mathbb{E} [\Phi(X)] - \Phi(\mathbb{E} [X]) \leq \mathbb{E} \left[ \sum_{i=1}^n \text{Ent}_{(-i)}(X) \right]. \quad (17)$$

**Proof:** The proposition is a direct consequence of Han's inequality for relative entropies. First note that if the inequality is true for a random variable  $X$ , then it is also true for  $cX$  where  $c$  is a positive constant. Hence, we may assume that  $\mathbb{E} [X] = 1$ . Now define the probability measure  $P$  on  $\mathcal{X}^n$  by its probability density function  $p$  given by

$$p(z) = f(z)q(z), \quad \forall z \in \mathcal{X}^n$$

where  $q$  denote the probability density of  $Z := (Z_1, Z_2, \dots, Z_n)$  and  $Q$  the corresponding probability measure. Then

$$\text{Ent}(X) := \mathbb{E} [X \log X] - \mathbb{E} [X] \log (\mathbb{E} [X]) = \text{KL} (P \parallel Q)$$

which, by Han's inequality for relative entropy

$$\text{Ent}(X) = \text{KL} (P \parallel Q) \leq \sum_{i=1}^n (\text{KL} (P \parallel Q) - \text{KL} (P_{(-i)} \parallel Q_{(-i)}))$$

However, straightforward calculation shows that

$$\sum_{i=1}^n (\text{KL} (P \parallel Q) - \text{KL} (P_{(-i)} \parallel Q_{(-i)})) = \sum_{i=1}^n \mathbb{E} [\mathbb{E}_{(-i)} [\Phi(X)] - \Phi(\mathbb{E}_{(-i)} [X])]$$

and the statement follows. ■

- **Remark** The Efron-Stein inequality is the special case of the inequality when  $\Phi(x) = x^2$ ,

$$\begin{aligned} \mathbb{E} [\Phi(X)] - \Phi(\mathbb{E} [X]) &\leq \sum_{i=1}^n \mathbb{E} [\mathbb{E}_{(-i)} [\Phi(X)] - \Phi(\mathbb{E}_{(-i)} [X])] \\ &\Rightarrow \text{Var}(X) \leq \sum_{i=1}^n \mathbb{E} [\text{Var}_{(-i)}(X)] \end{aligned}$$

## 2.5 Duality and Variational Formulas

## 2.6 Optimal Transport

## 2.7 Pinsker's Inequality

## 2.8 Birgé's Inequality

## 2.9 The Brunn-Minkowski Inequality

## References

Stéphane Boucheron, Gábor Lugosi, and Pascal Massart. *Concentration inequalities: A nonasymptotic theory of independence*. Oxford university press, 2013.

Thomas M. Cover and Joy A. Thomas. *Elements of information theory (2. ed.)*. Wiley, 2006. ISBN 978-0-471-24195-9.