# Information Theory and Statistics
# Lecture 1: Entropy and information

Łukasz Dębowski
ldebowsk@ipipan.waw.pl

Ph. D. Programme 2013/2014

KAPITAŁ LUDZKI
NARODOWA STRATEGIA SPÓJNOŚCI

PhD
STUDIES

UNIA EUROPEJSKA
EUROPEJSKI
FUNDUSZ SPOŁECZNY

## Claude Shannon (1916–2001)

Entropy of a random variable on a probability space is the fundamental concept of information theory developed by Claude Shannon in 1948.

KAPITAŁ LUDZKI
NARODOWA STRATEGIA SPÓJNOŚCI

PhD
STUDIES

UNIA EUROPEJSKA
EUROPEJSKI
FUNDUSZ SPOŁECZNY

## Probability as a random variable

### Definition

Let $\mathbf{X}$ and $\mathbf{Y}$ be discrete variables and $\mathbf{A}$ be an event on a probability space $(\Omega, \mathcal{J}, \mathbf{P})$. We define $\mathbf{P(X)}$ as a discrete random variable such that

$$\mathbf{P(X)}(\omega) = \mathbf{P(X = x)} \iff \mathbf{X}(\omega) = \mathbf{x}.$$

Analogously we define $\mathbf{P(X|Y)}$ and $\mathbf{P(X|A)}$ as

$$\mathbf{P(X|Y)}(\omega) = \mathbf{P(X = x|Y = y)} \iff \mathbf{X}(\omega) = \mathbf{x} \text{ and } \mathbf{Y}(\omega) = \mathbf{y},$$
$$\mathbf{P(X|A)}(\omega) = \mathbf{P(X = x|A)} \iff \mathbf{X}(\omega) = \mathbf{x},$$

where the conditional probability is $\mathbf{P(B|A) = P(B \cap A)/P(A)}$ for $\mathbf{P(A) > 0}$.

KAPITAŁ LUDZKI
NARODOWA STRATEGIA SPÓJNOŚCI

PhD
STUDIES

UNIA EUROPEJSKA
EUROPEJSKI
FUNDUSZ SPOŁECZNY

**Entropy**
○○○○●○○○○○

KL divergence
○○○○○

Conditional entropy
○○○○○

Mutual information
○○○

Conditional MI
○○○○○○○

## Independence

We write $P(Y) = P(X_1, X_2, ..., X_n)$ for $Y = (X_1, X_2, ..., X_n)$.

### Definition (independence)

We say that random variables $X_1, X_2, ..., X_n$ are independent if

$$P(X_1, X_2, ..., X_n) = \prod_{i=1}^{n} P(X_i).$$

Analogously, we say that random variables $X_1, X_2, X_3, ...$ are independent if $X_1, X_2, ..., X_n$ are independent for any $n$.

### Example

Let $\Omega = [0, 1]$ be the unit section and let $P$ be the Lebesgue measure. Define real random variable $Y(\omega) = \omega$. If we consider its binary expansion $Y = \sum_{i=1}^{\infty} 2^{-i} Z_i$, where $Z_i : \Omega \rightarrow \{0, 1\}$, then $P(Z_1, Z_2, ..., Z_n) = 2^{-n} = \prod_{i=1}^{n} P(Z_i)$.

**Entropy**
○○○○●○○○○○

KL divergence
○○○○○

Conditional entropy
○○○○○

Mutual information
○○○

Conditional MI
○○○○○○○

# Expectation

---

**Definition (expectation)**

We define the expectation of a real random variable $\mathbf{X}$ as

$$\mathbb{E}\,\mathbf{X} := \int \mathbf{X}\,d\mathbf{P}.$$

---

For discrete random variables we obtain

$$\mathbb{E}\,\mathbf{X} = \sum_{x:\mathbf{P}(\mathbf{X}=x)>0} \mathbf{P}(\mathbf{X}=x) \cdot x.$$

Entropy
○○○○○●○○○○

KL divergence
○○○○○

Conditional entropy
○○○○○

Mutual information
○○○

Conditional MI
○○○○○○○

# Additivity of expectation

One of fundamental properties of the expectation is its additivity.

## Theorem

Let $\mathbf{X}, \mathbf{Y} \geq \mathbf{0}$. We have

$$\mathbb{E}(\mathbf{X} + \mathbf{Y}) = \mathbb{E}\,\mathbf{X} + \mathbb{E}\,\mathbf{Y}.$$

Remark: The restriction $\mathbf{X}, \mathbf{Y} \geq \mathbf{0}$ is made because, e.g., for $\mathbb{E}\,\mathbf{X} = \infty$ and $\mathbb{E}\,\mathbf{Y} = -\infty$ the sum is undefined.

**Entropy**
○○○○○○○●○○○

KL divergence
○○○○○

Conditional entropy
○○○○○

Mutual information
○○○

Conditional MI
○○○○○○○

## Entropy

Some interpretation of entropy $\mathbf{H(X)}$ is the average uncertainty carried by a random variable $\mathbf{X}$. We expect that uncertainty adds for probabilistically independent sources. Formally, for $\mathbf{P(X, Y) = P(X)P(Y)}$, we postulate $\mathbf{H(X, Y) = H(X) + H(Y)}$. Because $\mathbf{log(xy) = log\, x + log\, y}$, the following definition comes as a very natural idea.

> **Definition (entropy)**
>
> The *entropy* of a discrete variable $\mathbf{X}$ is defined as
>
> $$\mathbf{H(X) := \mathbb{E}\left[-\log P(X)\right].} \qquad (1)$$
>
> Traditionally, it is assumed that $\mathbf{log}$ is the logarithm to the base $\mathbf{2}$.

Because $\mathbf{log\, P(X) \leq 0}$, we put the minus sign in the definition (1) so that entropy be positive.

KAPITAŁ LUDZKI
NARODOWA STRATEGIA SPÓJNOŚCI

PhD
STUDIES

UNIA EUROPEJSKA
EUROPEJSKI
FUNDUSZ SPOŁECZNY

# Entropy (continued)

Equivalently, we have

$$\mathbf{H(X)} = - \sum_{x:P(X=x)>0} \mathbf{P(X = x) \log P(X = x)},$$

We can verify that for $\mathbf{P(X, Y) = P(X)P(Y)}$,

$$\mathbf{H(X, Y) = \mathbb{E}\left[- \log P(X, Y)\right] = \mathbb{E}\left[- \log P(X) - \log P(X)\right]}$$
$$\mathbf{= \mathbb{E}\left[- \log P(X)\right] + \mathbb{E}\left[- \log P(X)\right] = H(X) + H(Y).}$$

# Entropy for a binary variable



Figure: Entropy $H(X) = -p \log p - (1 - p) \log(1 - p)$ for $P(X = 0) = p$ and $P(X = 1) = 1 - p$.

**Entropy**
○○○○○○○○○●

KL divergence
○○○○○

Conditional entropy
○○○○○

Mutual information
○○○

Conditional MI
○○○○○○○

# The range of entropy in general

Because function $f(p) = -p \log p$ is strictly positive for $p \in (0, 1)$ and equals $0$ for $p = 1$, it can be easily seen that:

### Theorem

$H(X) \geq 0$, *whereas* $H(X) = 0$ *if and only if* $X$ *assumes only a single value.*

This fact agrees with the idea that constants carry no uncertainty. On the other hand, assume that $X$ takes values $x \in \{1, 2, ..., n\}$ with equal probabilities $P(X = x) = 1/n$. Then

$$H(X) = -\sum_{x=1}^{n} \frac{1}{n} \log \frac{1}{n} = \sum_{x=1}^{n} \frac{1}{n} \log n = \log n.$$

As we will see, $\log n$ is the maximal value of $H(X)$ if $X$ assumes values in $\{1, 2, ..., n\}$. That fact agrees with the intuition that the highest uncertainty occurs for uniformly distributed variables.

Entropy
○○○○○○○○○○

KL divergence
●○○○○

Conditional entropy
○○○○○

Mutual information
○○○

Conditional MI
○○○○○○○

# Kullback-Leibler divergence

A discrete probability distribution is a function $\mathbf{p} : \mathbb{X} \to [0, 1]$ on a countable set $\mathbb{X}$ such that $\mathbf{p}(x) \geq 0$ and $\sum_x \mathbf{p}(x) = 1$.

### Definition (entropy revisited)

The entropy of a discrete probability distribution is denoted as

$$H(\mathbf{p}) := - \sum_{x : \mathbf{p}(x) > 0} \mathbf{p}(x) \log \mathbf{p}(x).$$

### Definition (KL divergence)

*Kullback-Leibler divergence*, or *relative entropy* of probability distributions $\mathbf{p}$ and $\mathbf{q}$ is defined as

$$D(\mathbf{p}||\mathbf{q}) := \sum_{x : \mathbf{p}(x) > 0} \mathbf{p}(x) \log \frac{\mathbf{p}(x)}{\mathbf{q}(x)}.$$

Entropy
○○○○○○○○○○

KL divergence
○●○○○○

Conditional entropy
○○○○○

Mutual information
○○○

Conditional MI
○○○○○○○

## Convex and concave functions

### Definition (convex and concave functions)

A real function $f : \mathbb{R} \rightarrow \mathbb{R}$ is *convex* if

$$p_1 f(x_1) + p_2 f(x_2) \geq f(p_1 x_1 + p_2 x_2)$$

for $p_i \geq 0$ and $p_1 + p_2 = 1$. Moreover, $f$ is called *strictly convex* if

$$p_1 f(x_1) + p_2 f(x_2) > f(p_1 x_1 + p_2 x_2)$$

for $p_i > 0$ and $p_1 + p_2 = 1$. We say that function $f$ is *concave* if $-f$ is *convex*, whereas $f$ is *strictly concave* if $-f$ is *stricly convex*.

### Example

If function $f$ has a positive second derivative then it is strictly convex. Hence functions $h(x) = -\log x$ and $g(x) = x^2$ are strictly convex.

KAPITAŁ LUDZKI
NARODOWA STRATEGIA SPÓJNOŚCI

PhD
STUDIES

UNIA EUROPEJSKA
EUROPEJSKI
FUNDUSZ SPOŁECZNY

# Jensen inequality

**Theorem (Jensen inequality)**

If $\mathbf{f}$ is a convex function and $\mathbf{p}$ is a discrete probability distribution over real values then

$$\sum_{x:p(x)>0} p(x)f(x) \geq f\left(\sum_{x:p(x)>0} p(x) \cdot x\right).$$

Moreover, if $\mathbf{f}$ is strictly convex then

$$\sum_{x:p(x)>0} p(x)f(x) = f\left(\sum_{x:p(x)>0} p(x) \cdot x\right)$$

holds if and only if distribution $\mathbf{p}$ is concentrated on a single value.

The proof proceeds by induction on the number of values of $\mathbf{p}$.

# KL divergence is nonegative

### Theorem

*We have*

$$D(p||q) \geq 0,$$

*where the equality holds if and only if* $p = q$.

### Proof

By the Jensen inequality for $f(y) = -\log y$, we have

$$D(p||q) = -\sum_{x:p(x)>0} p(x) \log \frac{q(x)}{p(x)} \geq -\log \left( \sum_{x:p(x)>0} p(x) \frac{q(x)}{p(x)} \right)$$

$$= -\log \left( \sum_{x:p(x)>0} q(x) \right) \geq -\log 1 = 0.$$

Entropy
○○○○○○○○○○

KL divergence
○○○○●

Conditional entropy
○○○○○

Mutual information
○○○

Conditional MI
○○○○○○○

# The maximum of entropy

### Theorem

Let $\mathbf{X}$ assume values in $\{1, 2, ..., n\}$. We have $\mathbf{H(X)} \leq \log \mathbf{n}$, whereas $\mathbf{H(X)} = \log \mathbf{n}$ if and only if $\mathbf{P(X = x)} = 1/\mathbf{n}$.

Remark: If the range of variable $\mathbf{X}$ is infinite then entropy $\mathbf{H(X)}$ may be infinite.

### Proof

Let $\mathbf{p(x)} = \mathbf{P(X = x)}$ and $\mathbf{q(x)} = 1/\mathbf{n}$. Then

$$0 \leq \mathbf{D(p||q)} = \sum_{x:p(x)>0} \mathbf{p(x)} \log \frac{\mathbf{p(x)}}{1/\mathbf{n}} = \log \mathbf{n} - \mathbf{H(X)},$$

where the equality occurs if and only if $\mathbf{p} = \mathbf{q}$.

# Conditional entropy

The next important question is what is the behavior of entropy under conditioning. The intuition is that given additional information, the uncertainty should decrease. So should entropy. There are two distinct ways of defining conditional entropy.

---

**Definition (conditional entropy)**

*Conditional entropy* of a discrete variable $\mathbf{X}$ given event $\mathbf{A}$ is

$$\mathbf{H(X|A) := H(p)} \text{ for } \mathbf{p(x) = P(X = x|A)}.$$

*Conditional entropy* of $\mathbf{X}$ given a discrete variable $\mathbf{Y}$ is defined as

$$\mathbf{H(X|Y) := \sum_{y:P(Y=y)>0} P(Y = y)H(X|Y = y).}$$

---

Both $\mathbf{H(X|A)}$ and $\mathbf{H(X|Y)}$ are nonnegative.

KAPITAŁ LUDZKI
NARODOWA STRATEGIA SPÓJNOŚCI

PhD
STUDIES

UNIA EUROPEJSKA
EUROPEJSKI
FUNDUSZ SPOŁECZNY

# Minimum of conditional entropy

### Theorem

$H(X|Y) = 0$ *holds if and only if* $X = f(Y)$ *for a certain function* $f$ *except for a set of probability* $0$.

### Proof

Observe that $H(X|Y) = 0$ if and only if $H(X|Y = y) = 0$ for all $y$ such that $P(Y = y) > 0$. This holds if and only if given $(Y = y)$ with $P(Y = y) > 0$, variable $X$ is concentrated on a single value. Denoting this value as $f(y)$, we obtain $X = f(Y)$, except for the union of those sets $(Y = y)$ which have probability $0$.

KAPITAŁ LUDZKI
NARODOWA STRATEGIA SPÓJNOŚCI

PhD
STUDIES

UNIA EUROPEJSKA
EUROPEJSKI
FUNDUSZ SPOŁECZNY

# Another formula for conditional entropy

We have

$$H(X|Y) = \mathbb{E}\left[-\log P(X|Y)\right].$$

### Proof

$$H(X|Y) = \sum_y P(Y=y)H(X|Y=y)$$

$$= -\sum_{x,y} P(Y=y)P(X=x|Y=y)\log P(X=x|Y=y)$$

$$= -\sum_{x,y} P(X=x, Y=y)\log P(X=x|Y=y)$$

$$= \mathbb{E}\left[-\log P(X|Y)\right].$$

## Conditional entropy and entropy

Because $\mathbf{P(Y)P(X|Y) = P(X,Y)}$, by the previous result,

$$\mathbf{H(Y) + H(X|Y) = H(X,Y).}$$

Hence

$$\mathbf{H(X,Y) \geq H(Y).}$$

# Inequality $H(X|A) \leq H(X)$ need not hold

### Example

Let $P(X = 0|A) = P(X = 1|A) = 1/2$, whereas $P(X = 0|A^c) = 1$ and $P(X = 1|A^c) = 0$. Assuming $P(A) = 1/2$, we have $P(X = 0) = (1/2) \cdot (1/2) + (1/2) = 3/4$ and $P(X = 0) = (1/2) \cdot (1/2) = 1/4$ so

$$H(X) = -\frac{3}{4} \log \frac{3}{4} - \frac{1}{4} \log \frac{1}{4} = \log 4 - \frac{3}{4} \log 3 = 0.811....$$

On the other hand, we have $H(X|A) = \log 2 = 1$.

Despite that fact, $H(X|Y) \leq H(X)$ holds in general. Thus entropy decreases given additional information on average.

## Mutual information

To show that $H(X)$ is greater than $H(X|Y)$, it is convenient to introduce another important concept.

---

**Definition (mutual information)**

*Mutual information* between discrete variables $X$ and $Y$ is

$$I(X; Y) := \mathbb{E}\left[\log \frac{P(X, Y)}{P(X)P(Y)}\right].$$

---

We have $I(X; X) = H(X)$. Thus entropy is sometimes called *self-information*.

Entropy
ooooooooooo

KL divergence
ooooo

Conditional entropy
ooooo

**Mutual information**
o●o

Conditional MI
ooooooo

# Mutual information is nonnegative

## Theorem

*We have*

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

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

## Proof

Let $p(x, y) = P(X = x, Y = y)$ and $q(x, y) = P(X = x)P(Y = y)$. Then we have

$$I(X; Y) = \sum_{(x,y):p(x,y)>0} p(x, y) \log \frac{p(x, y)}{q(x, y)} = D(p||q) \geq 0$$

with the equality exactly for $p = q$.

KAPITAŁ LUDZKI
NARODOWA STRATEGIA SPÓJNOŚCI

PhD
STUDIES

UNIA EUROPEJSKA
EUROPEJSKI
FUNDUSZ SPOŁECZNY

## Mutual information and entropy

By the definition of mutual information,

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

Hence

$$H(X) + H(Y) \geq H(X, Y),$$
$$H(X) \geq H(X|Y), \; I(X; Y).$$

Moreover, we have $H(X|Y) = H(Y)$ if $X$ and $Y$ are independent, which also agrees with intuition.

# Conditional mutual information

In a similar fashion to conditional entropy, we define:

### Definition (conditional mutual information)

*Conditional mutual information* between discrete variables $\mathbf{X}$ and $\mathbf{Y}$ given event $\mathbf{A}$ is

$$I(X; Y|A) := D(p||q) \text{ for } p(x, y) = P(X = x, Y = y|A)$$
$$\text{and } q(x, y) = P(X = x|A)P(Y = y|A).$$

*Conditional mutual information* between discrete variables $\mathbf{X}$ and $\mathbf{Y}$ given variable $\mathbf{Z}$ is defined as

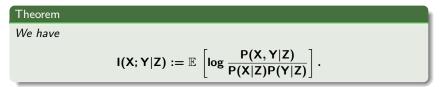$$I(X; Y|Z) := \sum_{z: P(Z=z) > 0} P(Z = z)I(X; Y|Z = z).$$

Both $\mathbf{I(X; Y|A)}$ and $\mathbf{I(X; Y|Z)}$ are nonnegative.

Entropy
0000000000

KL divergence
00000

Conditional entropy
00000

Mutual information
000

Conditional MI
0●00000

# Another formula for CMI

As in the case of conditional entropy, this proposition is true:

**Theorem**

*We have*

$$I(X; Y|Z) := \mathbb{E}\left[\log \frac{P(X, Y|Z)}{P(X|Z)P(Y|Z)}\right].$$

# Conditional independence

## Definition (conditional independence)

Variables $X_1, X_2, ..., X_n$ are *conditionally independent* given $Z$ if

$$P(X_1, X_2, ..., X_n | Z) = \prod_{i=1}^{n} P(X_i | Z).$$

Variables $X_1, X_2, X_3, ...$ are conditionally independent given $Z$ if $X_1, X_2, ..., X_n$ are conditionally independent given $Z$ for any $n$.

## Theorem

*We have*

$$I(X; Y | Z) \geq 0,$$

*with equality iff* $X$ *and* $Y$ *are conditionally independent given* $Z$.

KAPITAŁ LUDZKI
NARODOWA STRATEGIA SPÓJNOŚCI

PhD
STUDIES

UNIA EUROPEJSKA
EUROPEJSKI
FUNDUSZ SPOŁECZNY

## Conditional independence (examples)

### Example

Let $Y = f(Z)$ be a function of variable $Z$, whereas $X$ be an arbitrary variable. Variables $X$ and $Y$ are conditionally independent given $Z$. Indeed,

$$P(X = x, Y = y | Z = z) = P(X = x | Z = z)\mathbf{1}_{\{y=f(z)\}}$$
$$= P(X = x | Z = z)P(Y = y | Z = z).$$

### Example

Let variables $X$, $Y$, and $Z$ be independent. Variables $U = X + Z$ and $W = Y + Z$ are conditionally independent given $Z$. Indeed,

$$P(U = u, W = w | Z = z) = P(X = u - z, Y = w - z)$$
$$= P(X = u - z)P(Y = w - z)$$
$$= P(U = u | Z = z)P(W = w | Z = z).$$

## Markov chains

### Definition (Markov chain)

A stochastic process $(X_i)_{i=-\infty}^{\infty}$ is called a *Markov chain* if

$$P(X_i|X_{i-1}, X_{i-2}, ..., X_{i-n}) = P(X_i|X_{i-1})$$

holds for any $n$.

### Example

For a Markov chain $(X_i)_{i=-\infty}^{\infty}$, variables $X_i$ and $X_k$ are conditionally independent given $X_j$ if $i \leq j \leq k$. Indeed, after some algebra we obtain $P(X_k|X_i, X_j) = P(X_k|X_j)$, and hence

$$P(X_i, X_k|X_j) = P(X_i|X_j)P(X_k|X_i, X_j) = P(X_i|X_j)P(X_k|X_j).$$

KAPITAŁ LUDZKI
NARODOWA STRATEGIA SPÓJNOŚCI

PhD
STUDIES

UNIA EUROPEJSKA
EUROPEJSKI
FUNDUSZ SPOŁECZNY

## Conditional MI and MI

**Theorem**

*We have*

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

*Remark:* Hence, variables $X$ and $(Y, Z)$ are independent iff $X$ and $Z$ are independent and $X$ and $Y$ are independent given $Z$.

**Proof**

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

$$= \mathbb{E}\left[\log \frac{P(X, Y, Z)P(Z)}{P(X, Z)P(Y, Z)}\right] + \mathbb{E}\left[\log \frac{P(X, Z)}{P(X)P(Z)}\right]$$

$$= \mathbb{E}\left[\log \frac{P(X, Y, Z)}{P(X)P(Y, Z)}\right] = I(X; Y, Z).$$

## CMI and entropy

### Theorem

*If entropies* $\mathbf{H(X)}$, $\mathbf{H(Y)}$, *and* $\mathbf{H(Z)}$ *are finite, we have*

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

*where all terms are finite and nonnegative.*

The proof is left as an easy exercise.