Basic Notions of Probability and Information Theory

R. Basili

Course on *Deep Learning* a.a. 2024-25

March 13, 2025

<ロ> (四) (四) (三) (三) (三) (三)

Overview •	Information Theory	Entropy 0000	Joint-Entropy 0000	Mutual Information	Norms 0000000	References O
Outli	ine					

▲□▶ ▲□▶ ▲□▶ ▲□▶ ▲□ ● のへで

Outline

- Information Theory
- Entropy
- Joint-Entropy and Conditional entropy
- Mutual Information
- Cross-Entropy and Norms



How much information is there in knowing the outcome of ξ ?

< □ > < 同 > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > <



How much information is there in knowing the outcome of ξ ? Or equivalently:

How much uncertainty arises if the outcome ξ is unknown?



How much information is there in knowing the outcome of ξ ? Or equivalently:

How much uncertainty arises if the outcome ξ is unknown?

This is the information needed to specify which of the x_i has occurred. The problem is writing ξ .



How much information is there in knowing the outcome of ξ ? Or equivalently:

How much uncertainty arises if the outcome ξ is unknown?

This is the information needed to specify which of the x_i has occurred. The problem is writing ξ . Let us assume further that we only have a small set of symbols $A = \{a_k : k = 1, ...D\}$, that is a *coding alphabet*.



Thus each x_i will be represented by a string over A. Let us assume that ξ is *uniformly distributed*, i.e.

$$p_i = \frac{1}{M}$$
 $\forall i = 1, ..., M,$

▲□▶▲□▶▲□▶▲□▶ □ のQで

and that the coding alphabet is exactly $A = \{0, 1\}$.



Thus each x_i will be represented by a string over A. Let us assume that ξ is *uniformly distributed*, i.e.

$$p_i = \frac{1}{M}$$
 $\forall i = 1, ..., M,$

and that the coding alphabet is exactly $A = \{0, 1\}$. Thus, each x_i will be represented by a binary number. To use N binary digits to specify which x_i actually occurred means:

$$N: 2^{N-1} < M \le 2^N$$

Thus we need $N = \lceil \log_2 M \rceil$ digits. So what if the distribution is *nonuniform*, i.e., if the p_i s are not all equal?



How much uncertainty does a possible outcome with probability introduce?





How much uncertainty does a possible outcome with probability introduce? The basic assumption is that p_i will introduce equally much uncertainty regardless of the rest of the probabilities p_j with $j \neq i$.

< □ > < 同 > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > <



How much uncertainty does a possible outcome with probability introduce?

The basic assumption is that p_i will introduce equally much uncertainty regardless of the rest of the probabilities p_j with $j \neq i$.

We can thus reduce the problem to the case where all outcomes have probability p_i . In this case, there are $\frac{1}{p_i} = M_{p_i}$ possible outcomes.



How much uncertainty does a possible outcome with probability introduce?

The basic assumption is that p_i will introduce equally much uncertainty regardless of the rest of the probabilities p_j with $j \neq i$.

We can thus reduce the problem to the case where all outcomes have probability p_i . In this case, there are $\frac{1}{p_i} = M_{p_i}$ possible outcomes.

Example: if $p_i \approx 1$ then $M_{p_i} \approx 1$.



How much uncertainty does a possible outcome with probability introduce? We can thus reduce the problem to the case where all outcomes have probability p_i . In this case, there are $\frac{1}{p_i} = M_{p_i}$ possible outcomes.

< □ > < 同 > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > <



How much uncertainty does a possible outcome with probability introduce?

We can thus reduce the problem to the case where all outcomes have probability p_i . In this case, there are $\frac{1}{p_i} = M_{p_i}$ possible outcomes.

For a binary coding alphabet, we thus need

$$\log_2 M_{p_i} = \log_2 \frac{1}{p_i} = -\log_2 p_i$$

binary digits to specify that the outcome was x_i . Thus, the uncertainty introduced by p_i is in the general case

$$-\log_2 p_i$$

Overview O	Information Theory	Entropy ●000	Joint-Entropy 0000	Mutual Information	Norms 0000000	References O
Entre	opv					

Uncertainty of ξ

I -

The uncertainty introduced by the random variable ξ will be taken to be the *expectation value of the number of digits* required to specify its outcome.

▲□▶▲□▶▲□▶▲□▶ ▲□ ● のへで

Overview O	Information Theory	Entropy ●000	Joint-Entropy 0000	Mutual Information	Norms 0000000	References O
Entre	onv					

Uncertainty of ξ

r J

The uncertainty introduced by the random variable ξ will be taken to be the *expectation value of the number of digits* required to specify its outcome. This is the expectation value of $-\log_2 P(\xi)$, i.e.

$$E[-\log_2 P(\xi)] = \sum_i -p_i \log_2 p_i$$

▲□▶▲□▶▲□▶▲□▶ □ のQで

Overview O	Information Theory	Entropy 0000	Joint-Entropy 0000	Mutual Information	Norms 0000000	References O
Entre	onv					

Entropy

P J

The entropy $H[\xi]$ of ξ is precisely the amount of uncertainty introduced by the random variable ξ and it is more often referred to a natural logarithm ln(.), so that

$$H[\xi] = E[-\ln p(\xi)] = \sum_{x_i \in \Omega_{\xi}} -p(x_i) \ln p(x_i) = \sum_{i}^{M} -p_i \ln p_i$$

▲□▶▲□▶▲□▶▲□▶ ▲□ ● のへで

Overview O	Information Theory	Entropy 00●0	Joint-Entropy 0000	Mutual Information	Norms 0000000	References O
Entre	onv					

Example 1: Dice

P J

In the Dado example, $\forall i = 1, ..., 6$, it follows that $p_i = \frac{1}{6}$.

$$H[\xi] = E[-\ln p(\xi)] = \sum_{x_i \in \Omega_{\xi}} -p(x_i)\ln p(x_i) = 6 \cdot \frac{1}{6}\ln 6 = 1,792$$

▲□▶ ▲□▶ ▲ □▶ ▲ □▶ ▲ □ ● ● ● ●

Overview O	Information Theory	Entropy 0000	Joint-Entropy 0000	Mutual Information	Norms 0000000	References O
Entr	onv					

Example 1: Dice

PJ

In the Dado example, $\forall i = 1, ..., 6$, it follows that $p_i = \frac{1}{6}$.

$$H[\xi] = E[-\ln p(\xi)] = \sum_{x_i \in \Omega_{\xi}} -p(x_i)\ln p(x_i) = 6 \cdot \frac{1}{6}\ln 6 = 1,792$$

Example 2: the Loosing Dice

A loosing Die: $p_1 = 1.00$, and $\forall i = 2, ..., 6, p_i = 0$.

$$H[\xi] = E[-\ln p(\xi)] = \sum_{x_i \in \Omega_{\xi}} -p(x_i)\ln p(x_i) = 1\ln 1 = 0$$

Overview O	Information Theory	Entropy 000●	Joint-Entropy 0000	Mutual Information	Norms 0000000	References O
Futu						

Consequence

Given a distribution p_i (i = 1, ..., M) for a discrete random variable ξ then for any other distribution q_i (i = 1, ..., M) over the same sample space Ω_{ξ} it follows that:

$$H[\xi] = -\sum_{i}^{M} p_i \ln p_i \le -\sum_{i}^{M} p_i \ln q_i$$

where equality holds **iff** the two distribution are the same, i.e. $\forall i = 1, ..., M$ $p_i = q_i$



Given two random variable ξ and η :

Joint-Entropy

the *joint entropy* of ξ and η is defined as:

$$H[\xi, \eta] = -\sum_{i=1}^{M} \sum_{j=1}^{L} p(x_i, y_j) \ln p(x_i, y_j)$$

▲□▶▲□▶▲□▶▲□▶ ▲□ ● のへで



Given two random variable ξ and η :

Joint-Entropy

the *joint entropy* of ξ and η is defined as:

$$H[\xi,\eta] = -\sum_{i=1}^{M} \sum_{j=1}^{L} p(x_i, y_j) \ln p(x_i, y_j) = H[\eta, \xi]$$

▲□▶▲□▶▲□▶▲□▶ □ のQで

Overview O	Information Theory	Entropy 0000	Joint-Entropy 0●00	Mutual Information	Norms 0000000	References O

Conditional-entropy

Conditional Entropy

the *conditional entropy* $H[\xi|\eta]$ of ξ and η is defined as:

$$H[\xi|\eta] = -\sum_{j=1}^{L} p(y_j) \sum_{i=1}^{M} p(x_i|y_j) \ln p(x_i|y_j) = \\ = -\sum_{j=1}^{L} \sum_{i=1}^{M} p(x_i, y_j) \ln p(x_i|y_j)$$

▲□▶ ▲□▶ ▲ □▶ ▲ □▶ ▲ □ ● の Q ()

 Overview
 Information Theory
 Entropy
 Joint-Entropy

 0
 0000
 0000
 0000
 0000
 0000
 0000
 0000
 0000
 0000
 0000
 0000
 0000
 0000
 0000
 0000
 0000
 0000
 0000
 0000
 0000
 0000
 0000
 0000
 0000
 0000
 0000
 0000
 0000
 0000
 0000
 0000
 0000
 0000
 0000
 0000
 0000
 0000
 0000
 0000
 0000
 0000
 0000
 0000
 0000
 0000
 0000
 0000
 0000
 0000
 0000
 0000
 0000
 0000
 0000
 0000
 0000
 0000
 0000
 0000
 0000
 0000
 0000
 0000
 0000
 0000
 0000
 0000
 0000
 0000
 0000
 0000
 0000
 0000
 0000
 0000
 0000
 0000
 0000
 0000
 0000
 0000
 0000
 0000
 0000
 0000
 <

Mutual Information

Norms 0000000

▲□▶▲□▶▲□▶▲□▶ ▲□ ● のへで

References O

Conditional and joint entropy

Conditional and Joint Entropy

The conditional and joint entropies are related just like the conditional and joint probabilities:

 $H[\xi,\eta]=H[\eta]+H[\xi|\eta]$

erview Information Theory Entropy

Joint-Entropy 00●0 Mutual Information

Norms 0000000 References 0

Conditional and joint entropy

Conditional and Joint Entropy

The conditional and joint entropies are related just like the conditional and joint probabilities:

$$H[\xi,\eta] = H[\eta] + H[\xi|\eta]$$

Conveyed Information

The *information conveyed* by η , denoted $I[\xi|\eta]$, is the reduction in entropy of ξ by finding out the outcome of η . This is defined by:

 $I[\boldsymbol{\xi}|\boldsymbol{\eta}] = H[\boldsymbol{\xi}] - H[\boldsymbol{\xi}|\boldsymbol{\eta}]$

Joint-Entropy 000●

ヘロト 人間 とくほとくほとう

3

References

Conditional and joint entropy

Conditional and Joint Entropy

$$\begin{split} H[\xi,\eta] &= H[\eta] + H[\xi|\eta] \\ I[\xi|\eta] &= H[\eta] - H[\xi|\eta] \end{split}$$

Joint-Entropy 000● Mutual Information

lorms 0000000 References 0

Conditional and joint entropy

Conditional and Joint Entropy

$$\begin{split} H[\xi,\eta] &= H[\eta] + H[\xi|\eta] \\ I[\xi|\eta] &= H[\eta] - H[\xi|\eta] \end{split}$$

Consequences

Note that:

$$\begin{split} I[\xi|\eta] &= H[\xi] - H[\xi|\eta] = H[\xi] - (H[\xi,\eta] - H[\eta]) = \\ &= H[\xi] + H[\eta] - H[\xi,\eta] = H[\xi] + H[\eta] - H[\eta,\xi] = \\ &= H[\eta] + H[\xi] - H[\eta,\xi] = H[\eta] - H[\eta|\xi] = \\ &= I[\eta|\xi] \end{split}$$



Given two random variable ξ and η :

Mutual Information

the *mutual information* between ξ and η is defined as:

$$MI[\xi,\eta] = E[\ln\frac{P(\xi,\eta)}{P(\xi) \cdot P(\eta)}] =$$

=
$$\sum_{(x,y)\in\Omega_{(\xi,\eta)}} f_{(\xi,\eta)}(x,y) \ln\frac{f_{(\xi,\eta)}(x,y)}{f_{\xi}(x) \cdot f_{\eta}(y)}$$



Mutual Information measures the amount of information about a random variable ξ an observer receives when the outcome of a random variable η is available.

< □ > < 同 > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > <



Mutual Information measures the amount of information about a random variable ξ an observer receives when the outcome of a random variable η is available.



▲□▶ ▲□▶ ▲ □▶ ▲ □▶ ▲ □ ● ● ● ●



Mutual Information measures the amount of information about a random variable ξ an observer receives when the outcome of a random variable η is available.



How much information about the source output x_i does an observer gain by knowing the channel output y_i ?



Mutual Information measures the amount of information about a random variable ξ an observer receives when the outcome of a random variable η is known, in fact:

Mutual Information

$$MI[\xi,\eta] = H[\xi] - H[\xi|\eta] =$$

=
$$\sum_{(x,y)\in\Omega_{(\xi,\eta)}} f_{(\xi,\eta)}(x,y) \ln \frac{f_{(\xi,\eta)}(x,y)}{f_{\xi}(x) \cdot f_{\eta}(y)}$$

▲□▶▲□▶▲□▶▲□▶ ▲□ ● のへで



Mutual Information measures the amount of information about a random variable ξ an observer receives when the outcome of a random variable η is known, in fact:

Mutual Information

$$MI[\xi,\eta] = H[\xi] - H[\xi|\eta] =$$

=
$$\sum_{(x,y)\in\Omega_{(\xi,\eta)}} f_{(\xi,\eta)}(x,y) \ln \frac{f_{(\xi,\eta)}(x,y)}{f_{\xi}(x) \cdot f_{\eta}(y)}$$

▲□▶▲□▶▲□▶▲□▶ ▲□ ● のへで

Overview O	Information Theory	Entropy 0000	Joint-Entropy 0000	Mutual Information	Norms 00000
---------------	--------------------	-----------------	-----------------------	--------------------	----------------

▲□▶ ▲□▶ ▲□▶ ▲□▶ ▲□ ● のへで

Mutual Information

MI and H

$MI[\xi,\eta] = H[\xi] - H[\xi|\eta]$

Overview O Information Theory

Entro 000 Joint-Entropy 0000 Mutual Information

lorms 0000000

▲□▶ ▲□▶ ▲ □▶ ▲ □▶ ▲ □ ● ● ● ●

References 0

Mutual Information

MI and H

$$\begin{split} &MI[\xi,\eta]=H[\xi]-H[\xi|\eta]\\ &H[\xi,\eta]=H[\eta,\xi]\\ &H[\xi,\eta]=H[\eta]+H[\xi|\eta], \end{split}$$

Overview O Information Theory 0000 opy Jo

Mutual Information

lorms 0000000 References 0

Mutual Information

MI and H

 $egin{aligned} MI[\xi,\eta] &= H[\xi] - H[\xi|\eta] \ H[\xi,\eta] &= H[\eta,\xi] \ H[\xi,\eta] &= H[\eta] + H[\xi|\eta], \end{aligned}$

$$H[\boldsymbol{\xi}|\boldsymbol{\eta}] = H[\boldsymbol{\xi},\boldsymbol{\eta}] - H[\boldsymbol{\eta}]$$

◆□▶ ◆□▶ ◆目▶ ◆目▶ ●目 ● のへぐ

Entrop 0000 Joint-Entropy 0000 Mutual Information

lorms 0000000 References 0

Mutual Information

MI and H

$$egin{aligned} &MI[\xi,\eta] = H[\xi] - H[\xi|\eta] \ &H[\xi,\eta] = H[\eta,\xi] \ &H[\xi,\eta] = H[\eta] + H[\xi|\eta], \end{aligned}$$

$$H[\xi|\eta] = H[\xi,\eta] - H[\eta]$$

Symmetry

Note that mutual information is symmetric in ξ and η , that is $MI[\xi, \eta] = MI[\eta, \xi]$, as

$$H[\xi] - H[\xi|\eta] = H[\xi] + H[\eta] - H[\xi,\eta] = H[\eta] - H[\eta|\xi]$$

▲□▶▲□▶▲□▶▲□▶ ■ のへの



< □ > < 同 > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > <

Another way to look to mutual information is about the individual values (i.e. outcomes) $\xi = x_i$ and $\eta = y_j$.

Overview Information Theory Entropy Joint-Entropy Mutual Information Norms References of the second second

Pointwise Mutual Information

Another way to look to mutual information is about the individual values (i.e. outcomes) $\xi = x_i$ and $\eta = y_j$.

Pointwise Mutual Information

Given the two random variable ξ and η : the *pointwise mutual information* between $\xi = x_i$ and $\eta = y_j$ is defined as:

▲□▶▲□▶▲□▶▲□▶ □ のQで

$$MI[x_i, y_j] = f_{(\xi, \eta)}(x_i, y_j) \ln \frac{f_{(\xi, \eta)}(x_i, y_j)}{f_{\xi}(x_i) \cdot f_{\eta}(y_j)}$$

 Overview
 Information Theory
 Entropy
 Joint-Entropy
 Mutual Information
 Norms
 References

 Operate
 0000
 0000
 0000
 0000
 0000000
 0000000

Pointwise Mutual Information

Another way to look to mutual information is about the individual values (i.e. outcomes) $\xi = x_i$ and $\eta = y_j$.

Pointwise Mutual Information

Given the two random variable ξ and η : the *pointwise mutual information* between $\xi = x_i$ and $\eta = y_j$ is defined as:

$$MI[x_i, y_j] = f_{(\xi, \eta)}(x_i, y_j) \ln \frac{f_{(\xi, \eta)}(x_i, y_j)}{f_{\xi}(x_i) \cdot f_{\eta}(y_j)} = P(x_i, y_j) \ln \frac{P(x_i, y_j)}{P(x_i) \cdot P(y_j)}$$

▲□▶▲□▶▲□▶▲□▶ □ のQで

O	erview	

Entropy 0000 Joint-Entrop 0000 Mutual Information

lorms 0000000

▲□▶▲□▶▲□▶▲□▶ ▲□ ● のへで

References 0

Pointwise Mutual Information

Pointwise Mutual Information (pmi)

$$MI[x_i, y_j] = P(x_i, y_j) \ln \frac{P(x_i, y_j)}{P(x_i) \cdot P(y_j)}$$

Joint-E 0000 Mutual Information

lorms 0000000

References

Pointwise Mutual Information

Pointwise Mutual Information (pmi)

$$MI[x_i, y_j] = P(x_i, y_j) \ln \frac{P(x_i, y_j)}{P(x_i) \cdot P(y_j)}$$

Use of the pmi

If $MI[x_i, y_j] >> 0$, there is a strong correlation between x_i and y_j If $MI[x_i, y_j] << 0$, there is a strong negative correlation. When $MI[x_i, y_j] \approx 0$ the two outcomes are almost independent.

Overview O	Information Theory	Entropy 0000	Joint-Entropy 0000	Mutual Information	Norms 0000000	References O
Perpl	exity					

Perplexity

The *perplexity* of a random variable ξ is the exponential of its entropy, i.e.

 $Perp[\xi] = e^{H[\xi]}$

▲□▶ ▲□▶ ▲ □▶ ▲ □▶ ▲ □ ● ● ● ●

Overview O	Information Theory	Entropy 0000	Joint-Entropy 0000	Mutual Information	Norms 0000000	References O
Perpl	exity					

Perplexity

The *perplexity* of a random variable ξ is the exponential of its entropy, i.e.

$$Perp[\xi] = e^{H[\xi]}$$

Example

Predicting the next *w* of a sequence of *n* words $w_k \in Dict$:

$$P(\xi_n = w | \xi_{n-1} = w_{n-1}, \xi_{n-2} = w_{n-2}, \dots, \xi_1 = w_1)$$

What is $Perp[(\xi_n, ..., \xi_1)]$? OSS: In case of a uniform distribution $P(\xi_n = w|...) = \frac{1}{|Dict|}$...



Cross-entropy

If we have two distributions (collections of probabilities) p(x) and q(x) on Ω_{ξ} , then the *cross entropy* of *p* with respect to *q* is given by:

$$H_p[q] = -\sum_{x \in \Omega_{\xi}} p(x) \ln q(x)$$

▲□▶▲□▶▲□▶▲□▶ □ のQで



Cross-entropy

If we have two distributions (collections of probabilities) p(x) and q(x) on Ω_{ξ} , then the *cross entropy* of *p* with respect to *q* is given by:

$$H_p[q] = -\sum_{x \in \Omega_{\xi}} p(x) \ln q(x)$$

Minimality

$$H_p[q] = -\sum_{x \in \Omega_{\xi}} p(x) \ln q(x) \ge -\sum_{x \in \Omega_{\xi}} p(x) \ln p(x) \quad \forall q$$

implies that the cross entropy of a distribution q w.r.t. another distribution p is **minimal** when q is identical to p.

Entrop 0000 Joint-Entropy 0000 Mutual Information

Norms

References O

Cross-entropy as a Norm

Cross-entropy

$$H_p[q] = -\sum_{x \in \Omega_{\xi}} p(x) \ln q(x)$$

▲□▶ ▲圖▶ ▲臣▶ ★臣▶ = 臣 = のへで

Entropy 0000 Joint-Entropy 0000 Mutual Information

Norms

References O

Cross-entropy as a Norm

Cross-entropy

$$H_p[q] = -\sum_{x \in \Omega_{\xi}} p(x) \ln q(x)$$

Relative Entropy (or Kullback-Leibler distance)

$$D[p||q] = \sum_{x \in \Omega_{\xi}} p(x) \ln \frac{p(x)}{q(x)} = H_p[q] - H[p]$$

◆□ ▶ ◆昼 ▶ ◆臣 ▶ ◆臣 ● ● ● ●

Overview Information Theory Entropy Joint-Entrop 0 0000 0000 0000 Mutual Information

Norms 0000000 References 0

Cross-entropy and Norms

Relative Entropy (or Kullback-Leibler distance)

$$D[p||q] = \sum_{x \in \Omega_{\xi}} p(x) \ln \frac{p(x)}{q(x)} = H_p[q] - H[p]$$

KL distance: properties

 $D[p||q] \ge 0 \quad \forall q$

ロ > < 個 > < 目 > < 目 > < 目 > < 回 > < < の へ ()

Overview Information Theory Entropy Joint-Entro

Mutual Information

Norms 00●0000 References O

Cross-entropy and Norms

Relative Entropy (or Kullback-Leibler distance)

$$D[p||q] = \sum_{x \in \Omega_{\xi}} p(x) \ln \frac{p(x)}{q(x)} = H_p[q] - H[p]$$

KL distance: properties

$$D[p||q] \geq 0 \quad \forall q$$

$$D[p||q] = 0 \qquad \text{iff } q = p$$

(ロ)、(型)、(E)、(E)、(E)、(O)へ(C)

Overview Information Theory Entropy Joint-Entropy 0000 0000

Mutual Information

Norms 0000000

▲□▶▲□▶▲□▶▲□▶ ▲□ ● のへで

References O

Cross-entropy and Norms

Relative Entropy (or Kullback-Leibler distance)

$$D[p||q] = \sum_{x \in \Omega_{\xi}} p(x) \ln \frac{p(x)}{q(x)} = H_p[q] - H[p]$$

verview Information Theory Entropy Joint-En 0000 0000 0000 Mutual Information

Norms 0000000

▲□▶▲□▶▲□▶▲□▶ □ のQで

References O

Cross-entropy and Norms

Relative Entropy (or Kullback-Leibler distance)

$$D[p||q] = \sum_{x \in \Omega_{\xi}} p(x) \ln \frac{p(x)}{q(x)} = H_p[q] - H[p]$$

KL distance as a norm?

Unfortunately, as

$D[p||q] \neq D[q||p]$

the KL distance is *not* a valid metric in the classical terms. It is a *measure of the dissimilarity* between p and q.

Norms, Similarity and Learning

Why ranking probability distributions is necessary?

• During a learning process we need to figure out the circumstances (i.e. the state of affairs of the world) under which a certain concept/class/property manifest.

Why ranking probability distributions is necessary?

• During a learning process we need to figure out the circumstances (i.e. the state of affairs of the world) under which a certain concept/class/property manifest.

Mutual Information

Norms

0000000

References

• This make a direct reference to the probability of some (stochastic) event. Stochastic events are used to describe circumstances and properties.

Why ranking probability distributions is necessary?

• During a learning process we need to figure out the circumstances (i.e. the state of affairs of the world) under which a certain concept/class/property manifest.

Mutual Information

Norms

0000000

- This make a direct reference to the probability of some (stochastic) event. Stochastic events are used to describe circumstances and properties.
- Moreover, learning proceeds from experience, i.e. known facts or previous classified examples, to rules, i.e. probability joint distributions over *decisions* and *circumstances*

Why ranking probability distributions is necessary?

• During a learning process we need to figure out the circumstances (i.e. the state of affairs of the world) under which a certain concept/class/property manifest.

Mutual Information

Norms

0000000

- This make a direct reference to the probability of some (stochastic) event. Stochastic events are used to describe circumstances and properties.
- Moreover, learning proceeds from experience, i.e. known facts or previous classified examples, to rules, i.e. probability joint distributions over *decisions* and *circumstances*
- Learning in general means to induce the proper probability distributions from the known examples. There are several many ways to do it!!!

Why ranking probability distributions is necessary?

• During a learning process we need to figure out the circumstances (i.e. the state of affairs of the world) under which a certain concept/class/property manifest.

Mutual Information

Norms

0000000

- This make a direct reference to the probability of some (stochastic) event. Stochastic events are used to describe circumstances and properties.
- Moreover, learning proceeds from experience, i.e. known facts or previous classified examples, to rules, i.e. probability joint distributions over *decisions* and *circumstances*
- Learning in general means to induce the proper probability distributions from the known examples. There are several many ways to do it!!!

Why ranking probability distributions is necessary?

• **Consequences.** In general, we need to compare different inductive hypothesis (*IH*), that are different probability distributions *q_i* of the same decision,

Norms

0000000

References

Mutual Information

Why ranking probability distributions is necessary?

- **Consequences.** In general, we need to compare different inductive hypothesis (*IH*), that are different probability distributions *q_i* of the same decision,
- In order to do it, we measure the agreement of our hypothesis with the observations (i.e. a pool of annotated data kept aside, the *held out*, to validate the different q_i)

Mutual Information

Norms

0000000

Norms, Similarity and Learning

Information Theory

Why ranking probability distributions is necessary?

- **Consequences.** In general, we need to compare different inductive hypothesis (*IH*), that are different probability distributions *q_i* of the same decision,
- In order to do it, we measure the agreement of our hypothesis with the observations (i.e. a pool of annotated data kept aside, the *held out*, to validate the different q_i)

Mutual Information

Norms

イロト 不得 とうほう イヨン

-

0000000

References

• The result is an estimate of the similarity between the probability *q_i* induced at the *i*-th learning stage with the probability *p* characterizing the known examples.

Why ranking probability distributions is necessary?

- **Consequences.** In general, we need to compare different inductive hypothesis (*IH*), that are different probability distributions *q_i* of the same decision,
- In order to do it, we measure the agreement of our hypothesis with the observations (i.e. a pool of annotated data kept aside, the *held out*, to validate the different q_i)

Mutual Information

Norms

0000000

- The result is an estimate of the similarity between the probability *q_i* induced at the *i*-th learning stage with the probability *p* characterizing the known examples.
- The KL divergence $D[p||q] = H_p(q) H(p)$ can be the suitable dissimilarity function.

Why ranking probability distributions is necessary?

- **Consequences.** In general, we need to compare different inductive hypothesis (*IH*), that are different probability distributions *q_i* of the same decision,
- In order to do it, we measure the agreement of our hypothesis with the observations (i.e. a pool of annotated data kept aside, the *held out*, to validate the different q_i)
- The result is an estimate of the similarity between the probability *q_i* induced at the *i*-th learning stage with the probability *p* characterizing the known examples.
- The KL divergence $D[p||q] = H_p(q) H(p)$ can be the suitable dissimilarity function.
- The probability \hat{q} (such that \hat{q} minimizes $\forall i D[p||q_i]$) is returned.

Norms

0000000

Mutual Information

Overview O	Information Theory	Entropy 0000	Joint-Entropy 0000	Mutual Information	Norms 0000000	References O	
Norm							

What makes a function a norm?

Overview O	Information Theory	Entropy 0000	Joint-Entropy 0000	Mutual Information	Norms 000000	References O		
N 7	NT							

Norm

What makes a function a norm? Any binary mapping m between a set of objects $D \times D$ and the real numbes is a norm **iff**:

▲□▶▲□▶▲□▶▲□▶ □ のQで

Axioms

• (*Positive*) $m(X, Y) \ge 0$ $\forall X, Y \in D$ whereas $m(X, Y) = 0 \rightarrow X = Y.$

Overview O	Information Theory	Entropy 0000	Joint-Entropy 0000	Mutual Information	Norms 0000000	References O	
Norm							

What makes a function a norm? Any binary mapping *m* between a set of objects $D \times D$ and the real numbes is a norm **iff**:

Axioms

- (*Positive*) $m(X, Y) \ge 0$ $\forall X, Y \in D$ whereas $m(X, Y) = 0 \rightarrow X = Y$.
- (Simmetry) m(X, Y) = m(Y, X) $\forall X, Y \in D$

▲□▶▲□▶▲□▶▲□▶ ▲□ ● のへで

Overview O	Information Theory	Entropy 0000	Joint-Entropy 0000	Mutual Information	Norms 0000000	References O

Norm

What makes a function a norm? Any binary mapping m between a set of objects $D \times D$ and the real numbes is a norm **iff**:

Axioms

- (*Positive*) $m(X, Y) \ge 0$ $\forall X, Y \in D$ whereas $m(X, Y) = 0 \rightarrow X = Y$.
- (Simmetry) m(X, Y) = m(Y, X) $\forall X, Y \in D$
- (Triangle inequality) $m(X,Y) \le m(X,Z) + m(Z,Y)$ $\forall X, Y, Z \in D$

▲□▶▲□▶▲□▶▲□▶ ▲□ ● のへで

Overview O	Information Theory	Entropy 0000	Joint-Entropy 0000	Mutual Information	Norms 000000	References O
3.7						

Norm

What makes a function a norm? Any binary mapping m between a set of objects $D \times D$ and the real numbes is a norm **iff**:

Axioms

- (*Positive*) $m(X, Y) \ge 0$ $\forall X, Y \in D$ whereas $m(X, Y) = 0 \rightarrow X = Y$.
- (Simmetry) m(X, Y) = m(Y, X) $\forall X, Y \in D$
- (Triangle inequality) $m(X,Y) \le m(X,Z) + m(Z,Y)$ $\forall X, Y, Z \in D$

Euclidean Norm

$$\sqrt[2]{\sum_{x\in\Omega(\xi)}(p(x)-q(x))^2}$$

Overview O	Information Theory	Entropy 0000	Joint-Entropy 0000	Mutual Information	Norms 0000000	References
Refere	ences					

Elementary Information Theory

 in (Krenn & Samuelsson, 1997), Brigitte Krenn, Christer Samuelsson, *The Linguist's Guide to Statistics Don't Panic*, Univ. of Saarlandes, 1997. URL:

http://nlp.stanford.edu/fsnlp/dontpanic.pdf

▲□▶▲□▶▲□▶▲□▶ ▲□ ● のへで