

Intro to Information Theory

Statistics and data analysis

Ben Galili

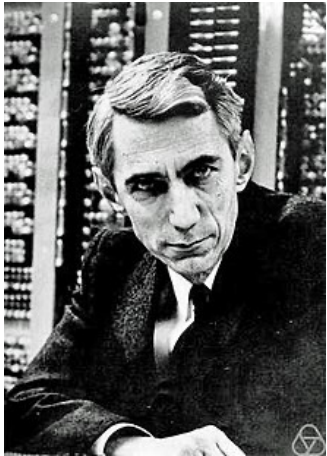
Zohar Yakhini

RUNI, Herzeliya

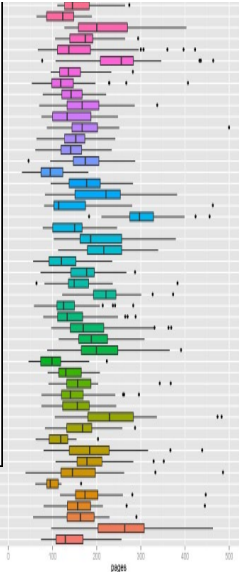
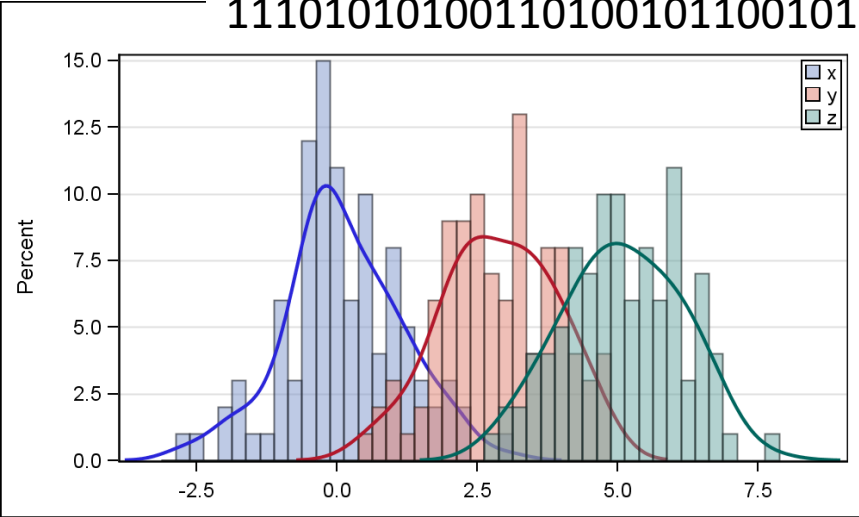
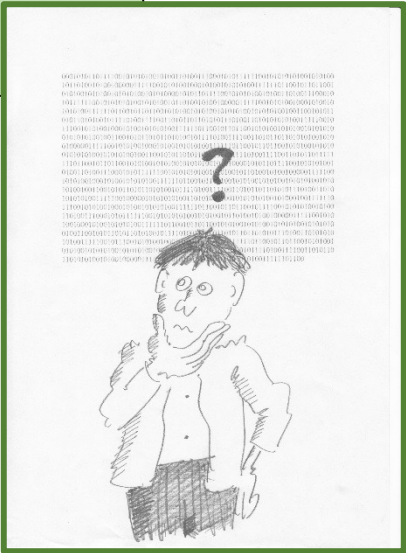
Shannon's entropy equation

$$H(X) = - \sum_{i=0}^{N-1} p_i \log_2 p_i$$

Mathematics for Life



0010011101010100101010100100100010
1010100010101111101011010011001001
11101010100111010010110010110110010



Background

The paper "A Mathematical Theory of Communication" was published by **Claude Shannon** in 1948.

This was the beginning of the information theory.

He defined the bit as an information measurement and the Entropy of information source as the minimal number of bits needed to code any message from the source

The Idea

The amount of information in a result of an experiment (that has more than one possible outcome) increases when the result is more surprising.

Example:

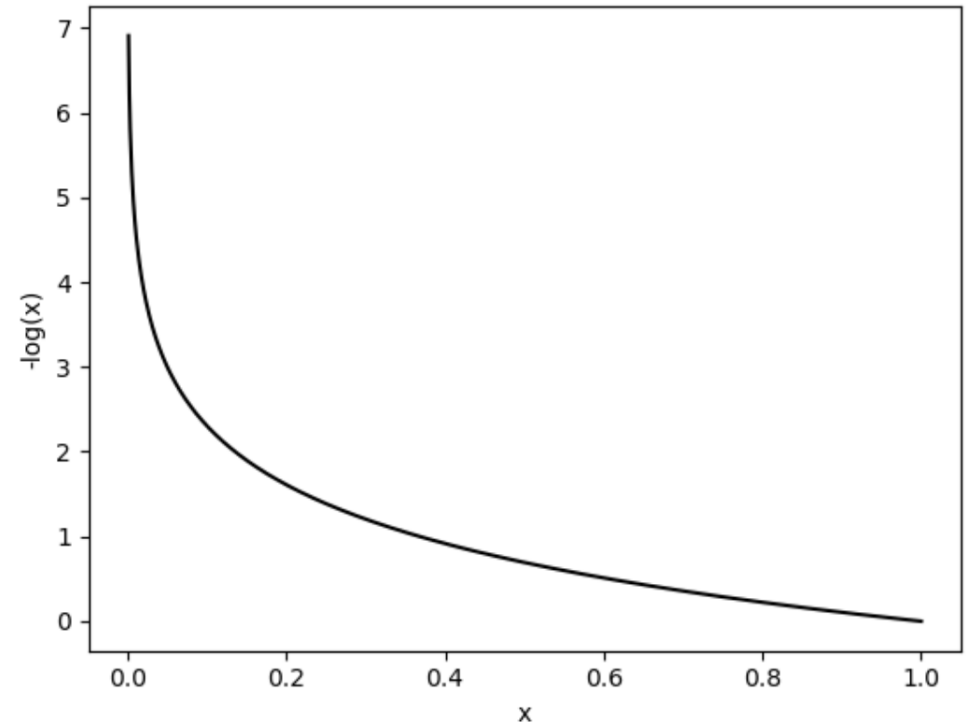
- We roll a fair die
- Where do we have more information? “the result is not 6” or “the result is 6”.

We want to measure the amount of information in a message/result

Information

$$I(p) = \log_2 \left(\frac{1}{p} \right)$$

- Large $p \rightarrow$ small $I(p)$
- $I(p) \geq 0$
- $I(1) = 0$
- If the events are independent
 $I(p_1, p_2) = I(p_1) + I(p_2)$
* $\log_2(p_1 p_2) = \log_2(p_1) + \log_2(p_2)$



Entropy – Definition

Let X be a discrete random variable with some PMF.

Let $H(X)$ be the Entropy of X (=the amount of information in the experiment defined by X):

$$H(X) = \sum_{x \in X} P(x) I(P(x)) = \sum_{x \in X} P(x) \log_2 \left(\frac{1}{P(x)} \right) = E \left[\log_2 \left(\frac{1}{P(x)} \right) \right]$$

Entropy = uncertainty measurement

Entropy – Comments

- Convention – $0 \log_2 0 = 0$
- The entropy depends only on the probabilities and NOT on the possible values of the random variable
- $H(X) \geq 0$

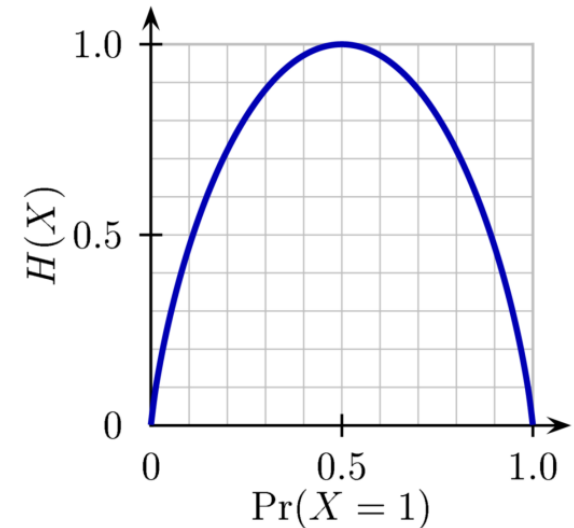
Entropy – Example 1

Let

$$X = \begin{cases} 0, & 1 - p \\ 1, & p \end{cases}$$

- $H(X) = p \log \left(\frac{1}{p} \right) + (1 - p) \log \left(\frac{1}{1-p} \right) = -p \log p - (1 - p) \log(1 - p)$
- $H(X) = \sum_x p \log \left(\frac{1}{p} \right) = -\sum_x p \log p$
- If $p = \frac{1}{2}$:

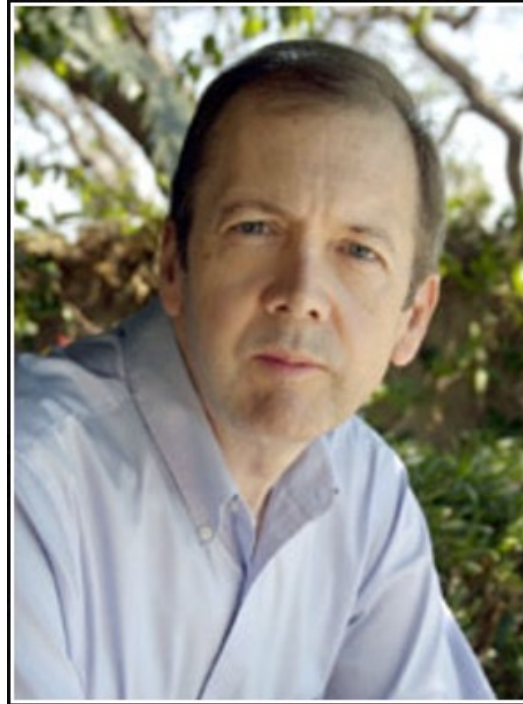
$$H(X) = -\frac{1}{2} \log \left(\frac{1}{2} \right) - \frac{1}{2} \log \left(\frac{1}{2} \right) = \log 2 = 1$$



Entropy – Example 2

Let

$$X = \begin{cases} a, & \frac{1}{2} \\ b, & \frac{1}{4} \\ c, & \frac{1}{8} \\ d, & \frac{1}{8} \end{cases}$$



Use "entropy" and you can never lose a debate, von Neumann told Shannon - because no one really knows what "entropy" is.

— William Poundstone —

AZ QUOTES

$$\begin{aligned} H(X) &= \frac{1}{2} \log 2 + \frac{1}{4} \log 4 + \frac{1}{8} \log 8 + \frac{1}{8} \log 8 \\ &= \frac{1}{2} \cdot 1 + \frac{1}{4} \cdot 2 + \frac{1}{8} \cdot 3 + \frac{1}{8} \cdot 3 = 1.75 \end{aligned}$$

Joint Entropy – Definition

Let X, Y be two discrete random variables.

Let $H(X, Y)$ be the Joint Entropy:

$$H(X, Y) = - \sum_{x \in X} \sum_{y \in Y} P(x, y) \log(P(x, y)) = -E[\log(P(x, y))]$$

Conditional Entropy – Definition

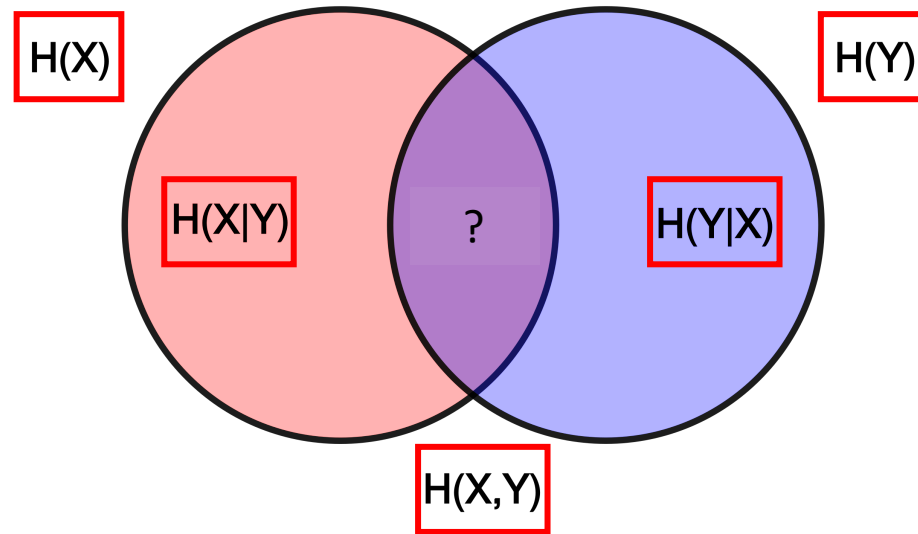
Let X, Y be two discrete random variables.

Let $H(Y|X)$ be the Conditional Entropy:

$$\begin{aligned} H(Y|X) &= \sum_{x \in X} P(x) H(Y|X = x) = - \sum_{x \in X} P(x) \sum_{y \in Y} P(y|x) \log(P(y|x)) \\ &= - \sum_{x \in X} \sum_{y \in Y} P(x, y) \log(P(y|x)) = -E_{P(x,y)} [\log(P(y|x))] \end{aligned}$$

Chain Rule

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



Chain Rule – Proof

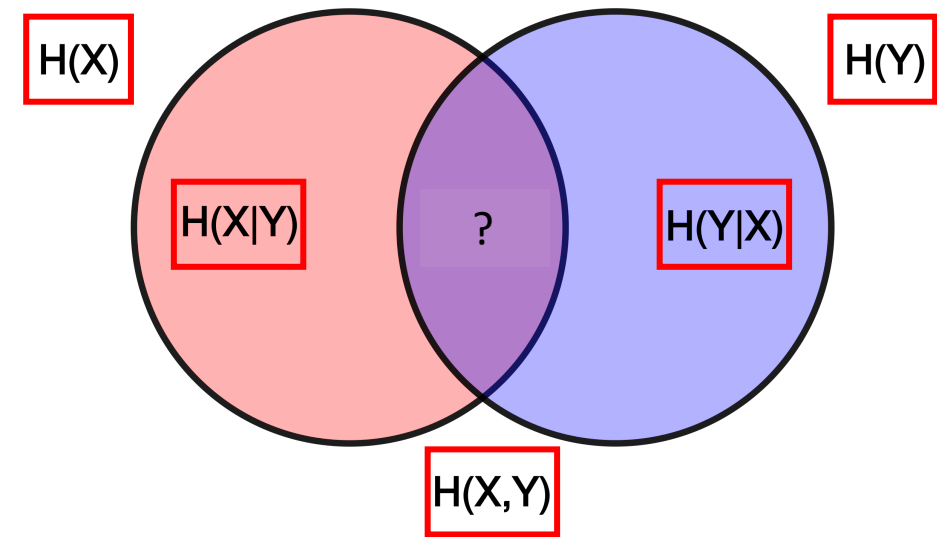
$$H(X, Y) = H(X) + H(Y|X) = H(Y) + H(X|Y)$$

$$\begin{aligned} H(X, Y) &= - \sum_{x \in X} \sum_{y \in Y} P(x, y) \log(P(x, y)) \\ &= - \sum_{x \in X} \sum_{y \in Y} P(x, y) \log(P(x)P(y|x)) \\ &= - \sum_{x \in X} \sum_{y \in Y} P(x, y) \log(P(x)) - \sum_{x \in X} \sum_{y \in Y} P(x, y) \log(P(y|x)) \\ &= - \sum_{x \in X} P(x) \log(P(x)) - \sum_{x \in X} \sum_{y \in Y} P(x, y) \log(P(y|x)) \\ &= H(X) + H(Y|X) \end{aligned}$$

Chain Rule – Example

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

$\begin{matrix} X \\ Y \end{matrix}$	1	2	3	4	P(Y)
1	1/8	1/16	1/32	1/32	1/4
2	1/16	1/8	1/32	1/32	1/4
3	1/16	1/16	1/16	1/16	1/4
4	1/4	0	0	0	1/4
P(X)	1/2	1/4	1/8	1/8	



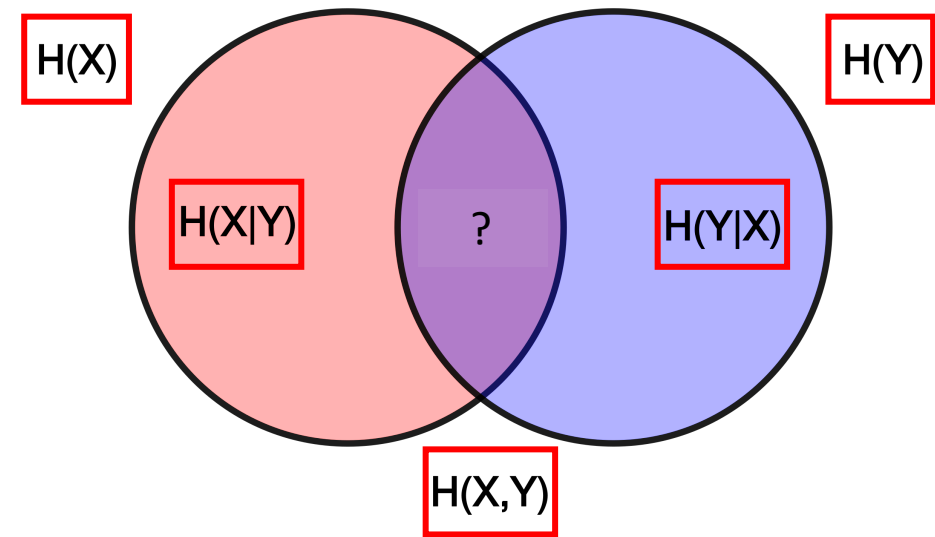
$H(X) = ?$

$H(Y) = ?$

Chain Rule – Example

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

$\begin{matrix} X \\ Y \end{matrix}$	1	2	3	4	P(Y)
1	1/8	1/16	1/32	1/32	1/4
2	1/16	1/8	1/32	1/32	1/4
3	1/16	1/16	1/16	1/16	1/4
4	1/4	0	0	0	1/4
P(X)	1/2	1/4	1/8	1/8	



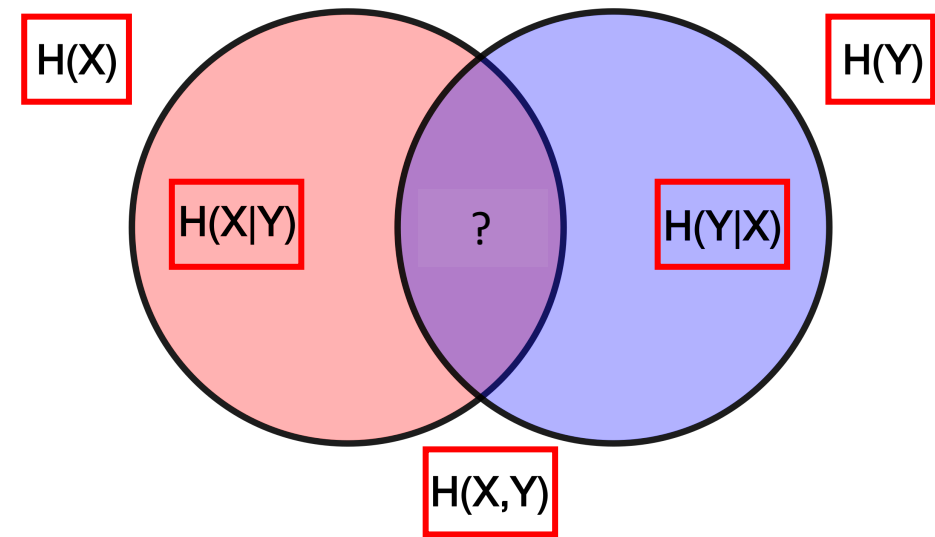
$$H(X|Y) = ?$$

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

Chain Rule – Example

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

$\begin{matrix} X \\ Y \end{matrix}$	1	2	3	4	P(Y)
1	1/8	1/16	1/32	1/32	1/4
2	1/16	1/8	1/32	1/32	1/4
3	1/16	1/16	1/16	1/16	1/4
4	1/4	0	0	0	1/4
P(X)	1/2	1/4	1/8	1/8	

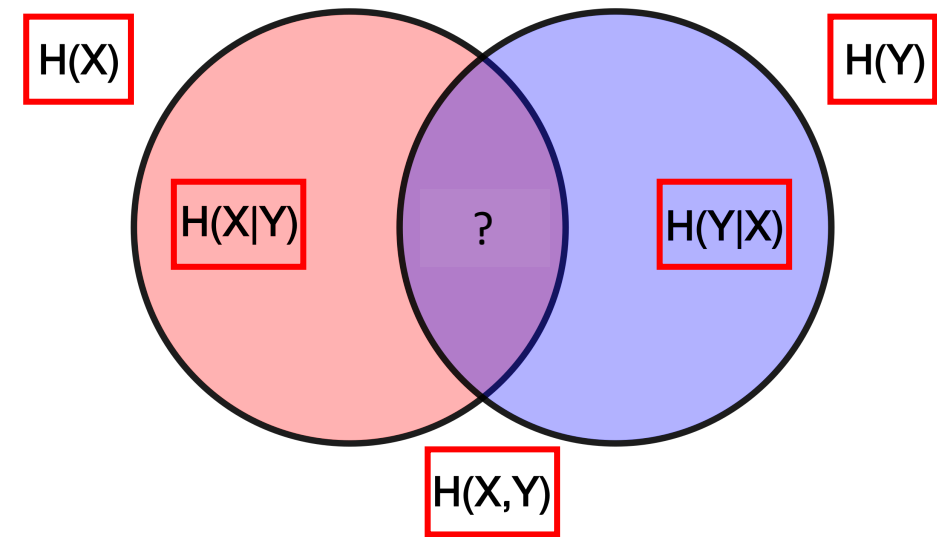


$$H(X, Y) = ?$$

Chain Rule – Example

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

$\begin{matrix} X \\ Y \end{matrix}$	1	2	3	4	P(Y)
1	1/8	1/16	1/32	1/32	1/4
2	1/16	1/8	1/32	1/32	1/4
3	1/16	1/16	1/16	1/16	1/4
4	1/4	0	0	0	1/4
P(X)	1/2	1/4	1/8	1/8	



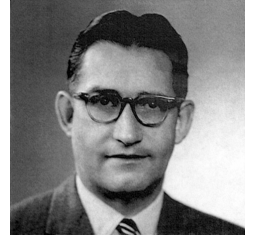
$$H(X) - H(X|Y) \stackrel{?}{=} H(Y) - H(Y|X)$$

We will define this term (the intersection) soon

Relative Entropy = Kullback-Leibler divergence



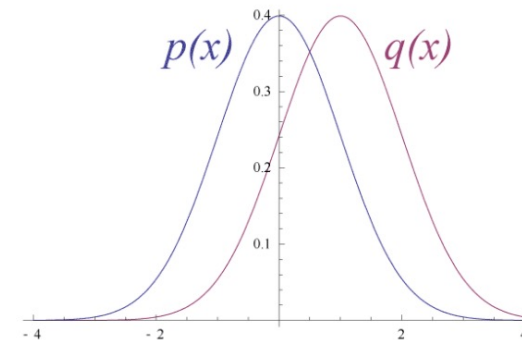
$$D_{\text{KL}}(P \parallel Q)$$



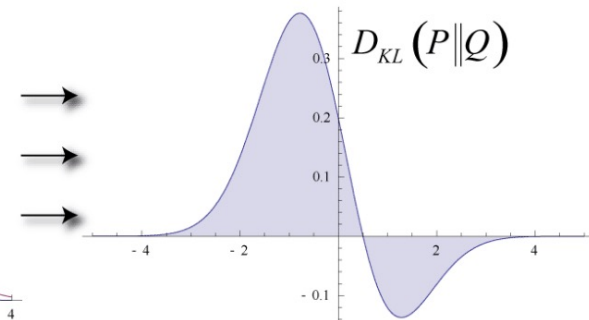
Measure the “distance” (difference) between the probability distribution P and the probability distribution Q .

$$D_{\text{KL}}(P \parallel Q) = \sum_{x \in X} p(x) \log \left(\frac{p(x)}{q(x)} \right) = E_p \left[\log \left(\frac{p(x)}{q(x)} \right) \right]$$

- $D_{\text{KL}}(P \parallel Q) \neq D_{\text{KL}}(Q \parallel P)$
- $D_{\text{KL}}(P \parallel Q) \geq 0$
- $P = Q \Leftrightarrow D_{\text{KL}}(P \parallel Q) = 0$
- $0 \log \frac{0}{0} = 0, 0 \log \frac{0}{q} = 0, 0 \log \frac{p}{0} = \infty$



Original Gaussian PDF's



KL Area to be Integrated

Jensen's Inequality

$$\text{Var}(X) = E[X^2] - E^2[X] \geq 0$$

Thus

$$E[X^2] \geq E^2[X]$$

If we define $g(x) = x^2$, we can write the above inequality as

$$E[g(X)] \geq g(E[X])$$

The function $g(x) = x^2$ is an example of convex function.

Jensen's inequality states that, for any convex function g , we have

$$E[g(X)] \geq g(E[X])$$



Johan Ludvig
William Valdemar
Jensen (1859-1925)

Jensen's Inequality

$g(x)$ is convex if and only if $-g(x)$ is concave.

We can state the definition for convex and concave functions in the following way:

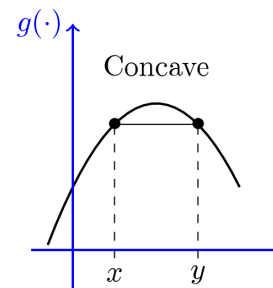
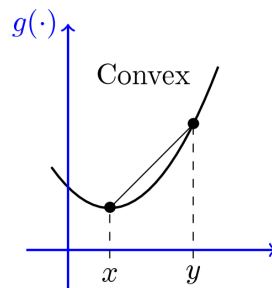
Consider a function $g: I \rightarrow \mathbb{R}$, where I is an interval in \mathbb{R} .

We say that g is a **convex** function if, for any two points x and y in I and any $\alpha \in [0,1]$, we have

$$g(\alpha x + (1 - \alpha)y) \leq \alpha g(x) + (1 - \alpha)g(y)$$

We say that g is **concave** if

$$g(\alpha x + (1 - \alpha)y) \geq \alpha g(x) + (1 - \alpha)g(y)$$



Jensen's Inequality

More generally, for a convex function $g: I \rightarrow \mathbb{R}$ and x_1, x_2, \dots, x_n in I and nonnegative real numbers α_i such that $\alpha_1 + \alpha_2 + \dots + \alpha_n = 1$, we have

$$g(\alpha_1 x_1 + \alpha_2 x_2 + \dots + \alpha_n x_n) \leq \alpha_1 g(x_1) + \alpha_2 g(x_1) + \dots + \alpha_n g(x_n)$$

If $n = 2$, the above statement is the definition of convex functions.
We can extend it to higher values of n by induction.

Jensen's Inequality

Now, consider a discrete random variable X with n possible values x_1, x_2, \dots, x_n .

In the previous equation,

$$g(\alpha_1 x_1 + \alpha_2 x_2 + \dots + \alpha_n x_n) \leq \alpha_1 g(x_1) + \alpha_2 g(x_1) + \dots + \alpha_n g(x_n)$$

we can choose $\alpha_i = P(X = x_i)$.

Then, the left-hand side becomes $g(E[X])$ and the right-hand side becomes $E[g(X)]$.

Jensen's Inequality:

If $g(x)$ is a convex function, and $E[g(X)]$ and $g(E[X])$ are finite, then

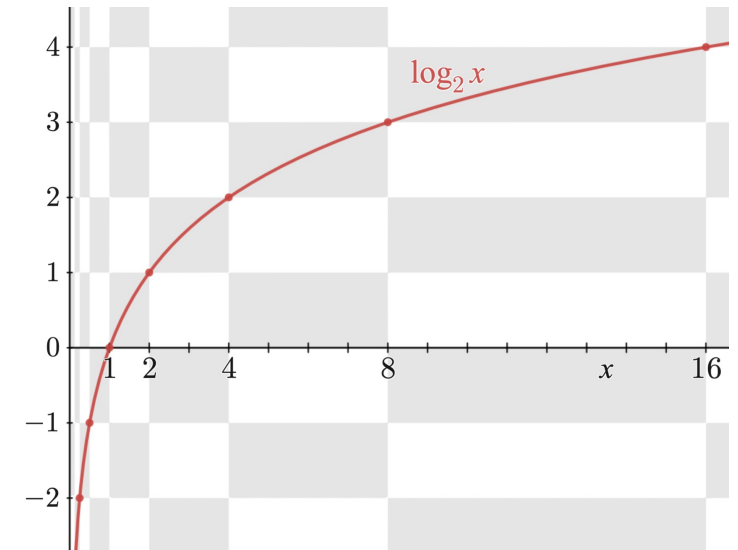
$$E[g(X)] \geq g(E[X])$$

Relative Entropy = Kullback-Leibler divergence

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

Proof:

$$\begin{aligned} D_{\text{KL}}(P \parallel Q) &= \sum_{x \in X} p(x) \log \left(\frac{p(x)}{q(x)} \right) \\ &= - \sum_{x \in X} p(x) \log \left(\frac{q(x)}{p(x)} \right) \\ &= -E_p \left[\log \left(\frac{q(x)}{p(x)} \right) \right] \\ &\geq -\log \left(E_p \left[\frac{q(x)}{p(x)} \right] \right) \text{ (by Jensen's Inequality for concave function log)} \\ &= -\log \left(\sum p(x) \frac{q(x)}{p(x)} \right) = -\log \left(\sum q(x) \right) = 0 \end{aligned}$$



Entropy Max Value

Consider a discrete random variable X with k possible values x_1, x_2, \dots, x_k .

$$H(X) = \sum_{x \in X} p(x) \log \frac{1}{p(x)} = E \left[\log \frac{1}{p(x)} \right]$$

(by Jensen's Inequality
for concave function \log)

$$\leq \log E \left[\frac{1}{p(x)} \right] = \log \sum_{x \in X} p(x) \frac{1}{p(x)} = \log k$$

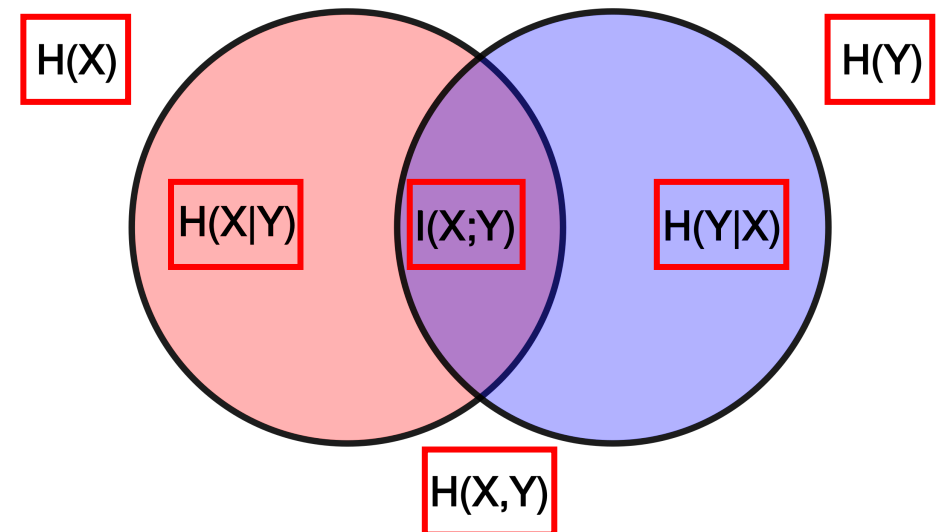
$$H(X) \leq \log k$$

Mutual Information

Let X and Y be two random variables with probability distributions $P(X)$ and $P(Y)$ respectively, and a joint distribution $P(X, Y)$.

The mutual information $I(X; Y)$ is the relative entropy between the joint distribution and the marginal distributions

$$\begin{aligned} I(X; Y) &= D_{KL}(P(x, y) \| P(x)P(y)) \\ &= \sum_{x \in X} \sum_{y \in Y} P(x, y) \log \left(\frac{P(x, y)}{P(x)P(y)} \right) \\ &= E_{p_{X,Y}} \left[\log \left(\frac{P(x, y)}{P(x)P(y)} \right) \right] \end{aligned}$$

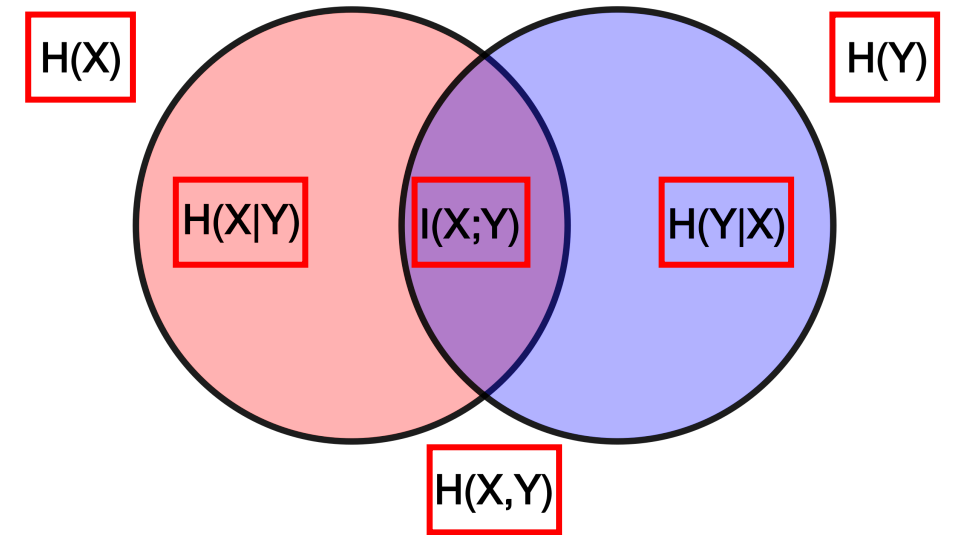


Mutual Information

$$I(X; Y) = \sum_{x \in X} \sum_{y \in Y} P(x, y) \log \left(\frac{P(x, y)}{P(x)P(y)} \right)$$

$$= \sum_{x, y} P(x, y) \log \left(\frac{P(x|y)}{P(x)} \right)$$

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

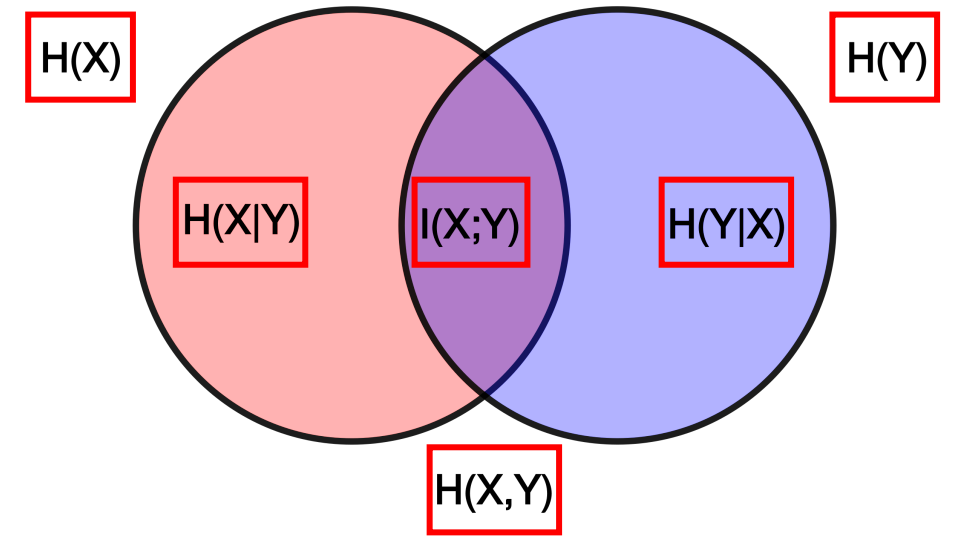


$$I(X; Y) = H(X) - H(X|Y) = H(Y) - H(Y|X)$$

$$H(X, Y) = H(Y) + H(X) - I(X; Y) \quad (\text{from the chain rule})$$

Mutual Information

- $I(X; X) = H(X) - H(X|X) = H(X)$
(Self information)
- $I(X; Y) \geq 0$
(Kullback-Leibler)
- $I(X; Y) = 0$ iff $\log \left(\frac{P(x,y)}{P(x)P(y)} \right) = 0$
iff X, Y are independent
- $H(X) \geq H(X|Y)$



Example

$\begin{matrix} X \\ Y \end{matrix}$	1	2	$P(Y)$
1	0	$\frac{3}{4}$	$\frac{3}{4}$
2	$\frac{1}{8}$	$\frac{1}{8}$	$\frac{1}{4}$
$P(X)$	$\frac{1}{8}$	$\frac{7}{8}$	

$$H(X) = H\left(\frac{1}{8}, \frac{7}{8}\right) = 0.544$$

$$H(Y) = H\left(\frac{1}{4}, \frac{3}{4}\right)$$

$$H(X|Y = 1) = 0$$

$$H(X|Y = 2) = 1$$

$$H(X|Y) = P(Y = 1)H(X|Y = 1) + P(Y = 2)H(X|Y = 2) = \frac{3}{4} \cdot 0 + \frac{1}{4} \cdot 1 = 0.25$$

Entropy of a Function of a Random Variable

- Let X be a random variable and let $g(X)$ be a function on X
 $X \in \mathbb{R} \quad g: \mathbb{R} \rightarrow \mathbb{R}$

$$X = \begin{cases} 1 & \frac{1}{3} \\ 0 & \frac{1}{3} \\ -1 & \frac{1}{3} \end{cases} \quad g(x) = x^2 = \begin{cases} 1 & \frac{2}{3} \\ 0 & \frac{1}{3} \end{cases}$$

- The entropy of a variable can only decrease when the latter is passed through a function

Entropy of a Function of a Random Variable

$$\begin{aligned} H(X, G(X)) &= H(X) + H(G(X)|X) \\ H(G(X)|X) &= 0 \rightarrow H(X, G(X)) = H(X) \end{aligned}$$

$$H(X) = H(X, G(X)) = H(G(X)) + H(X|G(X)) \geq H(G(X))$$

$$H(X) \geq H(G(X))$$

Summary

- Information
- Entropy
- Joint Entropy
- Conditional Entropy
- Relative Entropy
- Mutual Information
- Entropy of a function