

Information Theory

Information theory: provides a theoretical ground for ML in general

Problems we aim at solving:

- how to complete a sequence?
1 2 3 5 7 11 13 17 ?
↳ why is "15" more "probable"? how to justify it?
- how to decide between 2 models for given data?
- how to set a ML problem?
↳ which criterion to optimize?
↳ how to measure a solution's performance?

This will lead us to

the following problems:

- how to quantify information?
- is there a "natural" distribution over numbers?
models?
- how to compress (losslessly) data? are there bounds?

Information Theory is also (not covered in this course):

- communication in a noisy channel (Shannon) \Rightarrow error-correcting codes
- hash functions

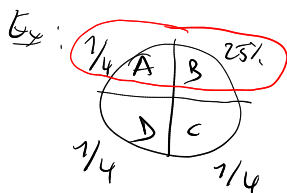
I - Information & Entropy

- quantify information

\rightarrow game of 20 binary questions

↳ level of information = number of questions to ask to identify the object

- pick an object according to a random law \Rightarrow can do better than dichotomy if I know the distribution
↳ fewer questions ↳ takes into account probabilities



- A, B, C, D: 4 possible objects

- I select one of them randomly

according to the law:

$$\begin{cases} p(A) = 1/4 \\ p(B) = 1/4 \\ p(C) = 1/4 \\ p(D) = 1/4 \end{cases}$$

Question 1: is it (A \cup B)?

yes

no D or C

Question 2: is it A?

is it C?

yes no
A B

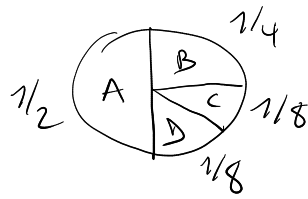
yes no
C D

\Rightarrow in all cases: you need to ask 2 questions to identify the object
↳ on average: 2 questions

⇒ The amount of information I need to provide you to identify the object, knowing that I picked it according to that law is 2 bits of information.

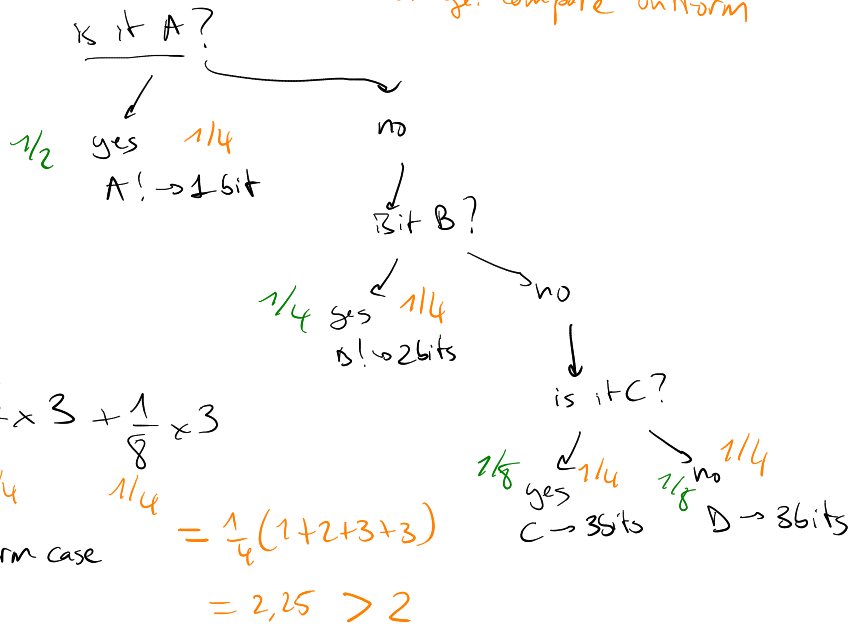
by answering binary questions

Non-uniform example:



green: non-uniform case

orange: compare uniform



On average:

$$= \frac{1}{2} \times 1 + \frac{1}{4} \times 2 + \frac{1}{8} \times 3 + \frac{1}{8} \times 3$$

$$= 1.75 < 2 \text{ in the uniform case}$$

$$= \frac{1}{4}(1+2+3+3) = 2.25 > 2$$

Desired properties:

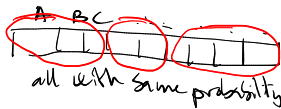
- continuous function H of the distrib^o of a random variable X :

$$H(\mathcal{Q}_X) = H((p_1, p_2, p_3, \dots))$$

$$\hookrightarrow H(p_1 + \epsilon, p_2 - \epsilon, p_3, \dots) \simeq H(p_1, p_2, p_3, \dots)$$

- symmetry: $H(p_1, p_2, p_3) = H(p_2, p_1, p_3)$

- additivity / independence to internal representation:



$$H(\text{uniform over } n \text{ bins}) =$$

$$H(\text{distribution over clusters})$$

average # of questions asked if proceeding in 2 steps

average # of questions with the direct process (1 step)

$$+ \sum_i p(\text{cluster } i) \times H(\text{uniform over bins inside cluster } i)$$

or

$$H(p_1, p_2, p_3, \dots) = H(p_1 + p_2, p_3, \dots) + (p_1 + p_2) H\left(\frac{p_1}{p_1 + p_2}, \frac{p_2}{p_1 + p_2}\right)$$

- uniform distribution has maximum entropy & its entropy increases with the number of bins

Entropy definition

Proposition [Khindin, 1957]: Any function F satisfying the properties above

can be written as: $F(p) = -C \sum_i p_i \log p_i$

$$\vec{p} = (p_1, p_2, p_3, \dots)$$

for some constant $C \in \mathbb{R}^+$

We choose $C=1, \log_2$: Shannon entropy

↳ unit: bit: number of binary questions needed

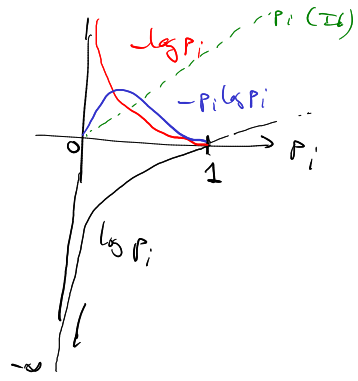
Examples:

- uniform law over n bins: $\forall i, p_i = 1/n$

$$H(p) = -\sum_i p_i \log p_i = \sum_i p_i \times \underbrace{(-\log p_i)}_{> 0}$$

$$= n \times \frac{1}{n} \times (-\log \frac{1}{n})$$

$$= \log_2 n \geq 0$$



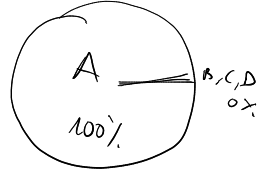
- Dirac peak:



$$H = 0$$

$$= \sum 0 \log 0 + 1 \times 0 = 0$$

ABCD previous case

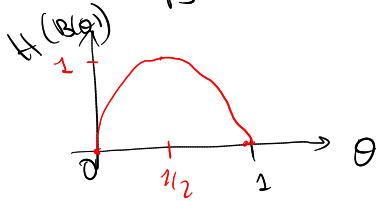


- Bernoulli $B(\theta)$

· toss a coin \rightarrow head θ
 \rightarrow tail $1-\theta$

- if $\theta = 0$ or 1 : Dirac peak $\rightarrow H = 0$

- if $\theta = 1/2$: maximum entropy: uniform over 2 bits $\rightarrow H = \log_2 2 = 1$



Other properties

- $H(X) \geq 0$
- $H(X) = 0$ iff $\exists i: p_i = 1$ (Dirac peak) no information
- $H(\text{any distribution over } n \text{ bins}) \leq \log_2 n$
- $H(F(X)) \leq H(X)$ For any (deterministic) function F (cf. later)

II - Why "log p"?

- describe (encode) an integer $\in [1, 2^b]$: with b bits
 \rightarrow ask $b = \log n$ binary questions
 $= -\log \frac{1}{n}$ \rightarrow binary tree dict. binary

11111
 15 : 01111
 $[1, 5]$
 $[1, 32]$

- code any event with $-\log p$ bits

- For $n = 2^b$, uniform law: $\log n = -\log p$ $p = 1/n$
- if non uniform, $p = 2^{-s}$: see ABCD case before
- general probability: $\sim \lceil -\log p \rceil$

\hookrightarrow we'll see later that we can reach $-\log p$



$\frac{1}{8} \leq p \leq \frac{1}{4} \Rightarrow$ need 3 questions
 $p \geq 1/8$

$$\frac{1}{p} \leq 8 \quad -\log p \leq 3$$

- entropy $H(p) : = \mathbb{E} \left[\underbrace{\text{information needed to describe the choice of that event } x}_{-\log p(x)} \right]$
 $=$ average length to encode events from p

- Ex: english text

- ↳ random indep. characters
- ASCII encoding: 256 possible characters $\Rightarrow \log(256) = 8$ bits \leftrightarrow uniform law over 256 symbols
- 27 characters are much more probable } $\log(27) = 4.8 \leftrightarrow$ uniform - - - - - 27 ..
- ↳ 26 alphabet + 1 space
- use character frequencies in English $\Rightarrow H(\text{English}) = 4.1$

III Extensions : For pairs of variables / of laws

- Conditional entropy.

joint law over (x, y) :

$$H(x, y) = H((x, y)) = - \sum_{i, j} p(x_i, y_j) \log p(x_i, y_j)$$

define $H(Y|X)$ as average entropy of $p(y|x)$

$$\begin{aligned} \hookrightarrow & \sum_x p(x) H(Y|x) = \mathbb{E} [H(Y|X=x)] = \sum_x p(x) \underbrace{\sum_y p(y|x) \log p(y|x)}_{\text{distribute over } y \text{ (for fixed } x)} \\ & = - \sum_{x, y} p(x, y) \log p(y|x) \end{aligned}$$

$$H(Y|X) \geq 0$$

$$H(Y|X) \leq H(Y)$$

$$\begin{aligned} H(x, y) &= H(x) + H(y|x) \\ &= H(y) + H(x|y) \end{aligned}$$

Bayes: $p(x, y) = p(x) p(y|x)$

$$-\log p(x, y) = -\log p(x) - \log(p(y|x))$$

$$\mathbb{E}_{x, y} \left\{ \begin{aligned} H(x, y) &= H(x) + H(y|x) \end{aligned} \right.$$

Mutual information :

$$H(x, y) \leq H(x) + H(y)$$

entropy difference

$$\begin{aligned} \text{mutual information: } I(x, y) &= H(x) + H(y) - H(x, y) \\ &= H(x) - H(x|y) \\ &= H(y) - H(y|x) \end{aligned}$$

information gain when considering x & y together

↳ symmetric ($H(x|y) = H(y|x)$)

$$\hookrightarrow I(x, y) \geq 0$$

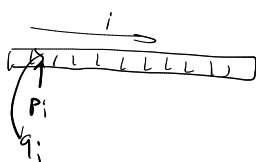
$$\hookrightarrow I(x, y) = 0 \text{ iff } x \perp y$$

$$\begin{aligned} &\downarrow \\ &\text{independent} \\ &p(x, y) = p(x) p(y) \end{aligned}$$

Relative entropy : Kullback-Leibler divergence

consider $x \sim p, y \sim q$

$$KL(p||q) = \sum_i p_i \log \frac{p_i}{q_i}$$



- $KL(P||Q) \geq 0$ $\forall p, q$ known as Gibbs inequality
 - $KL(P||P) = 0$
 - $KL(P||Q) = 0$ iff $P=Q$
- } KL as a sort of "distance" to compare 2 distributions P and Q

not symmetric: $KL(Q||P) \neq KL(P||Q)$ in general \Rightarrow not a "distance"

$$KL(P||Q) = \underbrace{E[-\log q]}_{\substack{\text{cross-entropy} \\ P, Q}} - \underbrace{H(P)}_{\substack{\text{entropy} \\ P}} = E[-\log p]$$

$KL \geq 0$ means cross-entropy \geq entropy

Rewriting mutual information

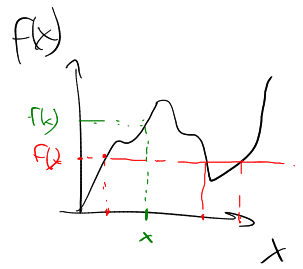
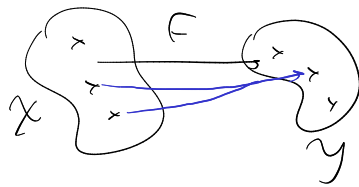
$$\begin{aligned}
 I(X;Y) &= H(X) + H(Y) - H(X,Y) \\
 &= -\sum_x p(x) \log p(x) - \sum_y q(y) \log q(y) + \sum_{x,y} p(x,y) \log p(x,y) \\
 &= -\sum_{x,y} p(x,y) \log p(x) - \sum_{x,y} p(x,y) \log p(y) + \dots \\
 &= -\sum_{x,y} p(x,y) \log \frac{p(x)p(y)}{p(x,y)} \\
 &= + KL(P(X,Y) || P(X)P(Y))
 \end{aligned}$$

p : joint law over $X \times Y$
 $P = X$
 $q = Y$
 $p(x,y) = p(x)p(y|x) = q(y)p(x|y)$

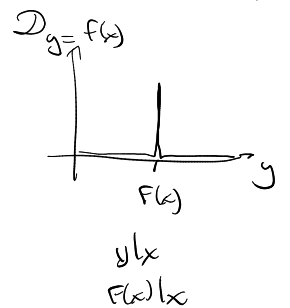
how far the joint law $P(X,Y)$ is from a product of independent laws
 i.e. how far X, Y are from being independent
 $X, Y \text{ indep} \Leftrightarrow p(x,y) = p(x)p(y)$

Information can only be lost by processing
 - consider a function F : deterministic

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



$$H(F(X)) \leq H(X)$$



III-B Sequences of non-independent variables

X_1, X_2, \dots, X_n sequence

$$H(X_1, X_2, \dots, X_n) = H(X_1) + H(X_2|X_1) + H(X_3|X_1, X_2) + \dots$$

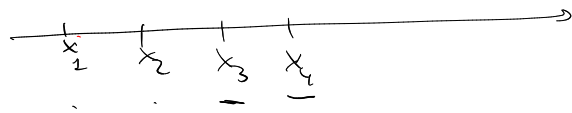
(Bayes)

$$P(X_1, \dots, X_n) = p(x_1)p(x_2|x_1)p(x_3|x_1, x_2) \dots$$

entropy rate:

$$H(X) = \lim_{n \rightarrow \infty} \frac{H(X_1, X_2, \dots, X_n)}{n}$$

Examples:
 - i.i.d. : $H(X) = \lim_{n \rightarrow \infty} \frac{H(X_1) + H(X_2) + H(X_3) + \dots + H(X_n)}{n} = \lim_{n \rightarrow \infty} \frac{n \cdot H(X_1)}{n} = H(X_1)$
 (identically distributed: $\text{Gauss}(X_i) = \text{Gauss}(X_j)$)
 - independent $X_i \perp X_j$
 - Markov chain: $p(X_n | X_1, \dots, X_{n-1}) = p(X_n | X_{n-1})$ order-1 M.C.



$$X_2 = 2X_1 \pm 1$$

$$X_3 = 2X_2 \pm 1 = 4X_1(\pm 1) \pm 1$$

$$p(X_3 | X_2) = p(X_3 | X_1, X_2)$$

$$H(X) = \lim_{n \rightarrow \infty} \frac{H(X_1) + H(X_2 | X_1) + H(X_3 | X_1, X_2) + \dots + H(X_n | X_1, \dots, X_{n-1})}{n}$$

$$= \lim_{n \rightarrow \infty} \frac{H(X_1) + (n-1)H(X_2 | X_1)}{n}$$

$$= H(X_2 | X_1)$$

$\frac{1}{n} H(X_1) \rightarrow 0$
 $\frac{n-1}{n} H(X_2 | X_1) \rightarrow 1$

III-C Continuous distributions

From discrete to continuous distributions:

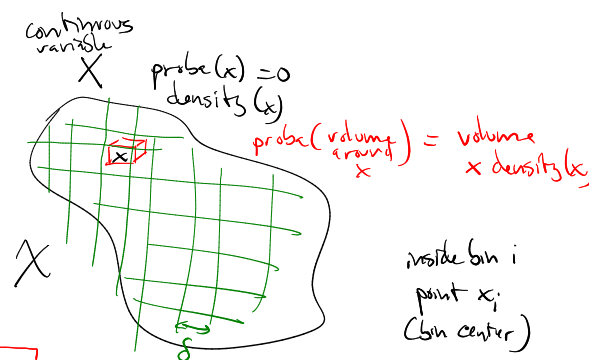
- continuous space X with probability density p

$$H(X_\delta) = - \sum_{\text{bins } i} p_\delta(i) \log p_\delta(i)$$

$$= - \sum_{\text{bins } i} \int_{\text{bin } i} p(x) \log p(x) dx$$

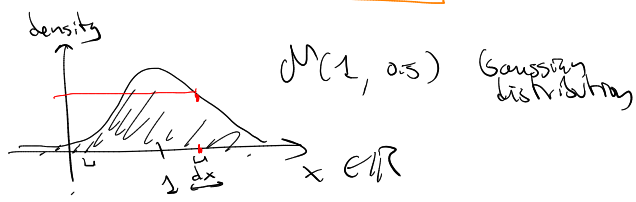
$$= - \int_X p(x) \log p(x) dx - \int_X p(x) D \log \delta$$

$\delta \rightarrow 0$
 $\int_X p \log p \rightarrow \int_X p \log p$
 $\int_X p D \log \delta \rightarrow 0$



discretization step \rightarrow back to discrete setting X_δ bins discrete random variable

Differential entropy: $H(X) := - \int_X p(x) \log p(x) dx$



$$p(x = 1.230470000) = 0$$

$$p(x \in [1.23 \dots \pm 0.001])$$

$$p_\delta(\text{bin } i) = \int_{\text{bin } i} p(x) dx$$

$$p(x)$$

$$p(x) dx$$

$$dp(x)$$

! $H(X)$ is not necessarily ≥ 0 !

$$\int_X p(x) dx = 1$$

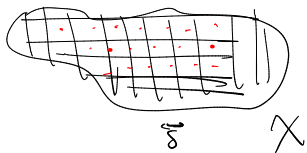
$$\int_X dp(x) = 1$$

$$\hookrightarrow KL(P||Q) = \int_X p(x) \log \frac{p(x)}{q(x)} dx$$

↳ same properties as before: $KL(P||Q) \geq 0$
and minimum (0) is reached when $P=Q$

because: $KL(P||Q) = \mathbb{E}_P[-\log Q] - H(P)$

Details about integrals



$$\sum_{\text{bins } i} \underbrace{\text{volume}(i)}_{\delta^D} f(\text{center of } i) \xrightarrow[\substack{\delta \rightarrow 0 \\ \text{discretized step} \\ \text{bin volume}}]{\text{}} \int_X f(x) dx$$

- D: dimension of the space X
- D=1: $x \in \mathbb{R}$
 - D=2: $x \in \mathbb{R}^2$
 - D=3:

III-D Example

EEG head set

- type letters just by focusing on them → no muscle involved
- control a mouse just by thinking left/right

↳ issue: very noisy signals

→ maximize information retrieved by optimizing protocol



→ optimal way to flash letters in order to identify the letter we focus on

Noisy-channel coding theorem [Shannon, 1948]

- transmitter: $X \rightarrow \text{noisy channel} \rightarrow Y$: receiver

- capacity: $C = \sup_{P_X} I(X; Y)$

- For any $\epsilon > 0$ and for any transmission rate $R \leq \text{channel capacity } C$,
 \exists encoding & decoding scheme that transmits data at rate R
with error probability $\leq \epsilon$ for a sufficiently large block length

- For any rate $R \geq C$, $p(\text{error}) \rightarrow 1$ - - - - -