

$X_i$  i.i.d.  $\sim p(x)$  (PMF)

$$-\frac{1}{n} \log p^{(n)}(X_1, \dots, X_n) = -\frac{1}{n} \sum_{i=1}^n \log p(X_i)$$

$$\rightarrow E[-\log p(X_1)] = H(X)$$

(a.e., ~~and~~ hence in prob.)

$$A_\varepsilon^{(n)} = \left\{ (x_1, \dots, x_n) \in \mathcal{X}^n \mid 2^{-n(H(X)+\varepsilon)} \leq p^{(n)}(x_1, \dots, x_n) \leq 2^{-n(H(X)-\varepsilon)} \right\}$$

typical set.

Th. 2.4.4.

a) If  $(x_1, \dots, x_n) \in A_\varepsilon^{(n)}$  then

$$H(X) - \varepsilon \leq -\frac{1}{n} \log p^{(n)}(x_1, \dots, x_n) \leq H(X) + \varepsilon$$

b)  $P((X_1, \dots, X_n) \in A_\varepsilon^{(n)}) > 1 - \varepsilon$  for  $n$  sufficiently large

c)  $|A_\varepsilon^{(n)}| \leq 2^{n(H(X)+\varepsilon)}$  ( $\cdot$  / cardinality)

d)  $|A_\varepsilon^{(n)}| \geq (1-\varepsilon) 2^{n(H(X)-\varepsilon)}$  for  $n$  suff. large.  $\square$

Proof. a)  $\checkmark$  b)  $\checkmark$

$$\begin{aligned} c) \quad 1 &= \sum_{x \in \mathcal{X}^n} p^{(n)}(x) \geq \sum_{x \in A_\varepsilon^{(n)}} p^{(n)}(x) \geq \sum_{x \in A_\varepsilon^{(n)}} 2^{-n(H(X)+\varepsilon)} \\ &= 2^{-n(H(X)+\varepsilon)} |A_\varepsilon^{(n)}| \end{aligned}$$

d) For sufficiently large  $n$

$$\begin{aligned} 1 - \varepsilon &< P((X_1, \dots, X_n) \in A_\varepsilon^{(n)}) \leq \sum_{x \in A_\varepsilon^{(n)}} 2^{-n(H(X)-\varepsilon)} \\ &= |A_\varepsilon^{(n)}| 2^{-n(H(X)-\varepsilon)} \end{aligned}$$

$\square$

For given  $\varepsilon > 0$  and sufficiently large  $n$

$\mathcal{X}^n$  decomposes into a set  $T = A_\varepsilon^{(n)}$  (typical set)

such that

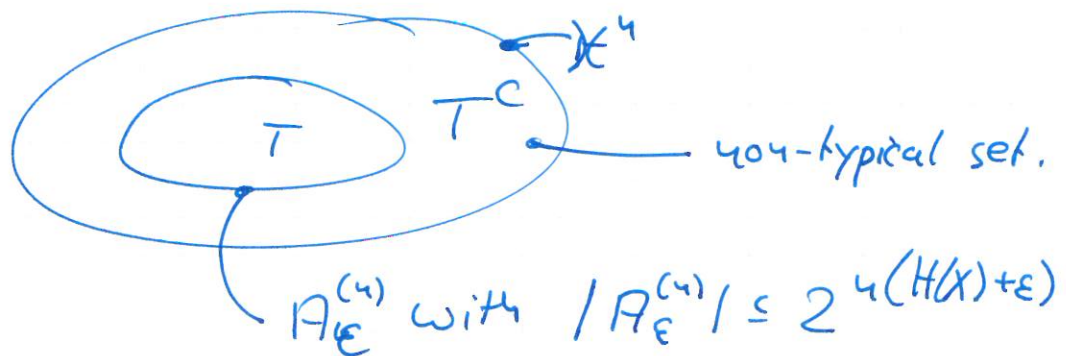
o  $P((X_1, \dots, X_n) \in T^c) \leq \varepsilon$

o For all  $x = (x_1, \dots, x_n) \in T$ :

$$\left| -\frac{1}{n} \log p^{(n)}(x_1, \dots, x_n) - H(X) \right| \leq \varepsilon$$

the normalized log-prob of all sequences in  $T$  is nearly equal and close to  $H(X)$ .

Graphically:



## The AEP and Data Compression

Let  $X_1, \dots, X_n$  i.i.d. with support  $\mathcal{X}$ ,  $X^{(n)} = (X_1, \dots, X_n)$ .

Aim: find a short description / encoding of  
all values  $x^{(n)} = (x_1, \dots, x_n) \in \mathcal{X}^n$ .

Key idea: index coding, allocate each of the  
 $|\mathcal{X}^n|$  values an index.

It holds •  $|A_\epsilon^{(n)}| \leq 2^{n(H(X) + \epsilon)}$  (Th. 2.4.4 c))

Indexing of all  $x^{(n)} \in A_\epsilon^{(n)}$  requires  
at most  $n(H(X) + \epsilon) + 1$

(1 bit extra since  $n(H(X) + \epsilon)$  may not be integer.)

•  $|\mathcal{X}^n|$  requires  $n \log |\mathcal{X}| + 1$  bits as indices.

Prefix each codeword for  $x^{(n)} \in A_\epsilon^{(n)}$  by 0

and each codeword for  $x^{(n)} \notin A_\epsilon^{(n)}$  by 1.

Let  $l(x^{(n)})$  denote the length of the codeword for  $x^{(n)}$ .

Then

$$E[l(X^{(n)})] = \sum_{x^{(n)} \in \mathcal{X}^n} p(x^{(n)}) l(x^{(n)})$$

$$= \sum_{x^{(n)} \in A_\epsilon^{(n)}} p(x^{(n)}) l(x^{(n)}) + \sum_{x^{(n)} \notin A_\epsilon^{(n)}} p(x^{(n)}) l(x^{(n)})$$

$$\leq \sum_{x^{(n)} \in A_\epsilon^{(n)}} p(x^{(n)}) (n(H(X) + \epsilon) + 2)$$

$$+ \sum_{x^{(n)} \notin A_\epsilon^{(n)}} p(x^{(n)}) (n \log |\mathcal{X}| + 2)$$

$$\begin{aligned}
 &= P(X^{(n)} \in A_{\epsilon}^{(n)}) (n(H(X) + \epsilon) + 2) \\
 &\quad + P(X^{(n)} \notin A_{\epsilon}^{(n)}) (n \log |X| + 2) \\
 &\leq n(H(X) + \epsilon) + \epsilon n \log |X| + 2 \\
 &= n \left( \underbrace{H(X) + \epsilon \frac{2}{n}}_{\epsilon'} + \epsilon \log |X| + \frac{2}{n} \right) \\
 &= n(H(X) + \epsilon') \text{ for any } \epsilon' > 0 \text{ with } n \text{ suff. large.}
 \end{aligned}$$

It follows:

Th. 2.4.5.  $\{X_n\}$  i.i.d. For any  $\epsilon > 0$  there exists  $n \in \mathbb{N}$  and a binary code that maps each  $x^{(n)}$  one-to-one onto a binary string satisfying

$$E\left(\frac{1}{n} \ell(X^{(n)})\right) \leq H(X) + \epsilon \quad \square$$

Hence, for suff. large  $n$  there exists a code for  $x^{(n)}$  such that the expected average codeword length is arbitrary close to  $H(X)$ .

## 2.5. Differential Entropy

By now: Entropy for discrete r.v. with finite support.

Extension: discrete r.v. but countably many support points,  $\mathcal{X} = \{x_1, x_2, \dots\}$  and distr.  $\mathbb{K} p = (p_1, p_2, \dots)$

$$H(X) = - \sum_{i=1}^{\infty} p_i \log p_i.$$

Note: The sum may be infinite or may not even exist.

Important: Extension of entropy to r.v.  $X$  with a density  $f$ .

Def. 2.5.1. Let  $X$  be abs.-continuous with density  $f(x)$ .

$$h(X) = - \int_{-\infty}^{\infty} f(x) \log f(x) dx$$

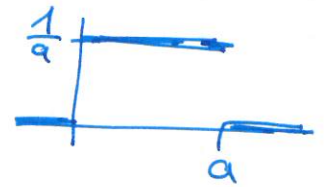
is called differential entropy of  $X$ .  $\perp$

Remarks.

- The integral in Def 2.5.1 may be infinite or may not even exist. ( $\rightarrow$  Exercises)
- As a general implicit assumption in defining  $h(X)$  we include: "provided the integral exists".
- $h(X) = E[-\log f(X)]$

Examples 2.5.2

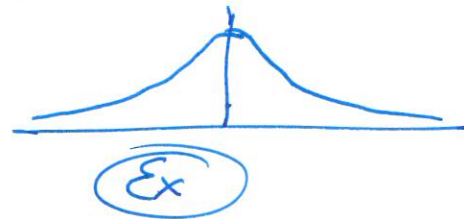
a)  $X \sim U(0, a)$ ,  $f(x) = \frac{1}{a} \mathbb{1}_{(0, a]}(x)$



$$h(X) = - \int_0^a \frac{1}{a} \log \frac{1}{a} dx$$

$$= \log a, \quad a > 0$$

b)  $X \sim N(\mu, \sigma^2)$ ,  $f(x) = \frac{1}{\sqrt{2\pi}\sigma} e^{-\frac{(x-\mu)^2}{2\sigma^2}}$ ,  $x \in \mathbb{R}$



$$h(X) = \frac{1}{2} \ln(2\pi e \sigma^2)$$

Def. 2.5.3.

a)  $X = (X_1, \dots, X_n)$  a  $n$ -vector with joint density  $f(x_1, \dots, x_n)$ .

$$h(X_1, \dots, X_n) = - \int \dots \int f(x_1, \dots, x_n) \log f(x_1, \dots, x_n) dx_1 \dots dx_n$$

is called joint differential entropy of  $X$ .

b)  $(X, Y)$  a  $n$ -vector with joint density  $f(x, y)$

and conditional density

$$f(x|y) = \frac{f(x, y)}{f(y)}, \quad \text{if } f(y) > 0,$$

and 0 otherwise.

$$h(X|Y) = - \iint f(x, y) \log f(x|y) dx dy$$

is called conditional diff. entropy of  $X$  given  $Y$ .

Example 2.5.4. (n-dim Gaussian distr.,  $X \sim N_n(\mu, C)$ )

$$\mu \in \mathbb{R}^n, C \in \mathbb{R}_{\text{sym}, \text{pos. def.}}^{n \times n}$$

$$f(x) = \frac{1}{(2\pi)^{n/2} |C|^{1/2}} \exp\left\{-\frac{1}{2} (x-\mu)^T C^{-1} (x-\mu)\right\},$$

$$x = (x_1, \dots, x_n)^T \in \mathbb{R}^n.$$

$$|\cdot| = |\det(\cdot)|$$

$$h(X) = \frac{1}{2} \ln((2\pi e)^n |C|) \quad \text{Ex.}$$

Def. 2.5.5. The mutual information between two r.v.  $X$  and  $Y$  with joint density  $f(x,y)$  is defined as

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

Interpretation: Amount of information about  $X$  from  $Y$  and vice versa.