

### III. Information Measures

- Shannon Entropy
- Differential Entropy
- Mutual Information

#### Shannon Entropy

Let  $X$  be a countable set  $\underline{P} \in \mathcal{P}(X)$

w/ pmf  $p$ . The Shannon entropy of  $X \sim \underline{P}$  is:

$$H(X) = H(\underline{P}) = H(p) := \mathbb{E}_{\underline{P}} \left[ \log \frac{1}{p(x)} \right]$$

Remark:  $H(X) = H(\underline{P})$  is a measure for the uncertainty/unpredictability of  $X \sim \underline{P}$ .

Example: For  $X \sim \underline{P}$ :

$$H(\underline{P}) = - D_{KL}(\underline{P} \parallel \#)$$

counting measure

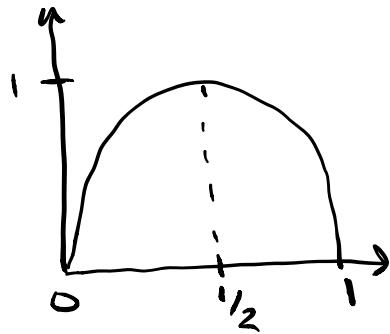
Furthermore, if  $|X| < \infty$ , then:

$$H(\underline{P}) = \log |X| - D_{KL}(\underline{P} \parallel \text{Unif}(x))$$

Example - Bernoulli:

$$X \sim \text{Ber}(\alpha), \quad \alpha \in [0, 1]$$

$$H(X) = H_b(\alpha) = \alpha \log \frac{1}{\alpha} + (1-\alpha) \log \frac{1}{1-\alpha}$$



Example -  $X \sim P = \text{Unif}(x)$

$$H(P) = \log |x|$$

Note:  $H(P) = +\infty$  is possible! Why?

Proposition: properties of  $H$ . Let  $X \sim P \in \mathcal{P}(x)$ ,  $X$  countable

1. Positivity:  $H(x) > 0$  w/ equality iff  $X$  is const.

2. Uniform Upper Bound: If  $|x| < \infty$ , then  $H(P) \leq \log |x|$

3. Invariance Under Relabeling:  $H(x) = H(f(x))$  for any bijective  $f$

4. Entropy of Functions: For any function  $f$ , we have  $H(x) \geq H(f(x))$  w/ equality iff  $f$  is a bijection

5. Concavity:  $P \mapsto H(P)$  is concave

## Definition (Joint Entropy)

For  $(X_1, \dots, X_n) \sim P_{x_1, \dots, x_n} \in \mathcal{P}(x_1, \dots, x_n)$  with joint pmf  $P_{x_1, \dots, x_n}$ , the joint shannon entropy is

$$H(X_1, \dots, X_n) := \mathbb{E}_{P_{x_1, \dots, x_n}} \left[ \log \frac{1}{P_{x_1, \dots, x_n}(X_1, \dots, X_n)} \right]$$

## Definition (Conditional Entropy)

Let  $(X, Y) \sim P_{XY} \in \mathcal{P}(x \times y)$  with pmf  $P_{XY}$ . The conditional entropy of  $Y$  given  $X$  is

$$H(Y|X) := \mathbb{E}_{P_{XY}} \left[ \log \frac{1}{P_{Y|X}(Y|X)} \right]$$

## Remark (Interpretation)

- Note that fixing  $x \in \mathcal{X}$ ,  $P_{Y|X}(\cdot|x) \in \mathcal{P}(y)$  so  $H(P_{Y|X}(\cdot|x))$  is a regular entropy as given in the first definition.
- Now:

$$H(Y|X) = \sum_{x \in \mathcal{X}} \sum_{y \in \mathcal{Y}} P_{XY}(x,y) \log \frac{1}{P_{Y|X=x}(y|x)}$$

$$= \sum_{x \in \mathcal{X}} \sum_{y \in \mathcal{Y}} P_X(x) P_{Y|X=x}(y|x) \log \frac{1}{P_{Y|X=x}(y|x)}$$

$$= \sum_{x \in \mathcal{X}} P_X(x) \sum_{y \in \mathcal{Y}} P_{Y|X=x}(y|x) \log \frac{1}{P_{Y|X=x}(y|x)}$$

$$= \sum_{x \in \mathcal{X}} P_X(x) H(Y|X=x) = \mathbb{E}_{x \sim P_X} [H(Y|X=x)]$$

## Proposition (properties continued)

5. Conditioning cannot increase entropy:  $H(X) \geq H(X|Y)$   
w/ equality iff  $X \perp\!\!\!\perp Y$

6. Entropy Chain Rule:

i) Small:

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

ii) Full:

$$\begin{aligned} H(X_1, \dots, X_n) &= H(X_1) + \sum_{i=2}^n H(X_i | X_1, \dots, X_{i-1}) \\ &\leq \sum_{i=1}^n H(X_i) \end{aligned}$$

w/ equality iff  $X_1, \dots, X_n$  mutually independent.

## Differentiable Entropy

Let  $P \in \mathcal{P}(\mathbb{R}^d)$ . The differentiable entropy of  $X \sim P$  is:

$$h(X) = h(P) := -D_{KL}(P \parallel \text{Leb}(\mathbb{R}^d)).$$

In particular, if  $P \ll \text{Leb}(\mathbb{R}^d)$  with pdf  $p$ , then:

$$h(X) = h(P) = h(p) = \mathbb{E}_P \left[ \log \frac{1}{P(x)} \right] = \int_{\text{Supp}(P)} p(x) \log \frac{1}{p(x)} dx$$

## Examples

$$\textcircled{1} \quad X \sim P = \text{Unif}([0, a]), \quad a > 0$$

$$h(X) = \log(a)$$

$$\textcircled{2} \quad X \sim P = N(\mu, \sigma^2)$$

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

$$X \sim P = N(\vec{\mu}, \Sigma)$$

$$h(X) = \frac{1}{2} \log((2\pi e)^d \det(\Sigma))$$

The definition of  $h(X)$  extends to  $h(X_1, \dots, X_n)$  and  $h(X|Y)$  similarly to what we saw for Shannon entropy.

### Aside

$$H(X|Y) = \mathbb{E}_{y \sim P_Y} \left[ H(X|Y=y) \right] \quad \begin{matrix} \text{could have} \\ \text{discrete} \end{matrix} = \int p_Y(y) H(X|Y=y) dy$$

also have

$$h(X|Y) = \mathbb{E}_{y \sim P_Y} \left[ h(X|Y=y) \right] \quad \begin{matrix} \text{continuous} \\ \text{continuous} \end{matrix} = \sum_{y \in Y} p_Y(y) h(X|Y=y) dy$$

## Proposition (properties of $h$ )

1. Uniform distribution maximizes  $h$  over a bounded domain

If  $\text{supp}(P) \subseteq \mathbb{R}^d$  is bounded, then

$$h(P) \leq h(\text{Unif}(\text{supp}(P)))$$

## 2. Gaussian maximizes h under variance constraint

If  $\text{supp}(P) = \mathbb{R}^d$  and  $X \sim P$  with

$$\Sigma := \mathbb{E}[(X - \mathbb{E}[X])^T(X - \mathbb{E}[X])]$$

then

$$h(P) \leq h(N(0, \Sigma))$$

## 3. Scaling and Shifting $X \sim P \in \mathcal{P}(\mathbb{R}^d)$

i)  $h(X + \mu) = h(X)$

ii)  $h(aX) = h(X) + d \log |a|$

iii)  $h(AX) = h(X) + \log(\det(A))$  for invertible A

## 4. Conditioning Cannot Increase Entropy

$$h(X) \geq h(X|Y)$$

w/ equality iff  $X \perp\!\!\!\perp Y$

## 5. Chain Rule

$$h(X_1, \dots, X_n) = h(X_1) + \sum_{i=2}^n h(X_i | X_1, \dots, X_{i-1})$$

$$\leq \hat{\sum}_{i=1}^n h(X_i)$$

w/ equality iff  $(X_1, \dots, X_n)$  mutually independent.

**Remark (Warning):** Note that  $h$  and  $H$ , while seemingly quantifying similar ideas, sometimes have vastly different behavior.

(i) While  $H(P) \geq 0$ ,  $h(P)$  can be positive/negative/unbounded

(ii)  $H(X) \geq H(f(X)) \rightarrow$  not true for  $h$

(iii) For  $H$ , we have

$$\begin{aligned} H(X+Y) &\leq H(X+Y, Y) \\ &= H(X, Y) \leq H(X) + H(Y) \end{aligned}$$

BUT this doesn't hold for  $h$