

# Mutual Information

Definition (Mutual Information): Let  $(X, Y) \sim P(x \times y)$ .

The mutual information between  $X$  and  $Y$  is

$$I(X; Y) \triangleq D_{KL}(P_{XY} \parallel P_X \otimes P_Y)$$

where  $P_X$  and  $P_Y$  are the  $X$  and  $Y$  marginals of  $P_{XY}$  and  $P_X \otimes P_Y$  is the induced product measure.

Remark:

- (i) MI is a "fundamental" measure of dependence between random variables
- (ii) A typical interpretation of mutual information is a quantification of the amount of "info" that  $X$  and  $Y$  convey about each other.

Proposition (Basic Properties of MI):

$$\textcircled{1} \quad I(X; Y) \geq 0 \text{ w/ equality iff } X \perp\!\!\!\perp Y$$

$$\textcircled{2} \quad I(X; Y) = D_{KL}(P_{Y|X} \parallel P_Y | P_X)$$

$$\textcircled{3} \quad I(X; Y) = I(Y; X)$$

(4)  $I(X; Y) \geq I(X; f(Y))$  for any deterministic function  
w/ equality iff  $f$  is a bijection.

(5)  $\underbrace{I(X, Y; Z)}_{D_{KL}(P_{XYZ} || P_{XY} \otimes P_Z)} \geq I(X; Z) \leftarrow \text{More information}$

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

### Proof

① by definition

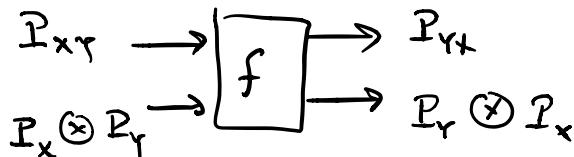
② chain rule of KL divergence

$$D_{KL}(P_{XY} || Q_{XY}) = D_{KL}(P_X || Q_X) + D_{KL}(P_{Y|X} || Q_{Y|X} | P_X)$$

• For us  $Q_{XY} = P_X \otimes P_Y$   
 $Q_X = P_X \quad Q_{Y|X} = P_Y$

$$\begin{aligned} &\Rightarrow D_{KL}(P_X || Q_X) + D_{KL}(P_{Y|X} || Q_{Y|X} | P_X) \\ &= D_{KL}(P_X || P_X) + D_{KL}(P_{Y|X} || P_Y | P_X) \\ &= 0 + D_{KL}(P_{Y|X} || P_Y | P_X) \end{aligned}$$

③ Let  $f(x,y) = (y,x)$  apply in BOTH directions



④ Establish later

⑤ Let  $f(x,y,z) = (x,z)$   $\nrightarrow$  use DPI

Ex

Proposition (I vs. H)

①

$$I(X;Y) = \begin{cases} H(X) & , X \text{ discrete} \\ \infty & , \text{o/w} \end{cases}$$

② Discrete X

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

③ Continuous X

$$\begin{aligned} I(X;Y) &= h(X) + h(Y) - h(X,Y) \\ &= h(X) - h(X|Y) \\ &= h(Y) - h(Y|X) \end{aligned}$$

Proof]

$$\textcircled{1} \quad D_{KL}(x; x) = D_{KL}(P_{x|x} \parallel P_x | P_x)$$

$$P_{x|x}(\cdot|x) = \delta_x(\cdot) \quad , \quad x \in \text{supp}(P_x)$$

(i) Discrete:  $\delta_x \ll P_x$

$$\begin{aligned} I(x; x) &= D_{KL}(P_{x|x} \parallel P_x | P_x) \\ &= \sum_{x \in X} p_x(x) \underbrace{D_{KL}(P_{x|x}(\cdot|x) \parallel P_x)}_{\delta_x(\cdot)} \\ &= \sum_{x \in X} p_x(x) \sum_{x' \in X} \delta_x(x') \frac{\log(\delta_x(x'))}{p_x(x')} \end{aligned}$$

$$= \sum_{x \in X} p_x(x) \log\left(\frac{1}{p_x(x)}\right) = H(X)$$

(ii) Continuous:  $I(x; x) = D_{KL}(P_{x|x} \parallel P_x \otimes P_x)$

and we will show  $P_{x|x} \not\ll P_x \otimes P_x$

Let

$$\Delta = \{(x, x) \mid x \in X\}$$

$$P_{x|x}(\Delta) = \int_{\Delta} dP_{x|x} = \int_X \int_X dP_{x|x}(x, x') \mathbf{1}_{\{x=x'\}} = \int_{x \in X} dP_x(x) \int_{x' \in X} dP_{x|x}(x'|x) \mathbf{1}_{\{x=x'\}}$$

$$= \int_{x \in X} dP_x(x) \int_{x' \in X} d\delta_x(x') = 1 > 0$$

$$\int_{\{x\}} d\delta_x(x') = \delta_x(\{x\}) = 1$$

However

$$\begin{aligned} P_x \otimes P_x (\Delta) &= \int_{\Delta} dP_x \otimes P_x (x, x') \\ &= \int_{\mathcal{X}} dP_x(x) \int_{\mathcal{X}} dP_x(x') \mathbb{1}_{\{x=x'\}} = 0 \\ &\quad \text{---} \\ &\quad \int_{\{x\}} dP_x(x') = P_x(\{x\}) = 0 \end{aligned}$$

□

This concludes the proof for continuous  $X$  case.

For an arbitrary non-discrete variable  $X$ :

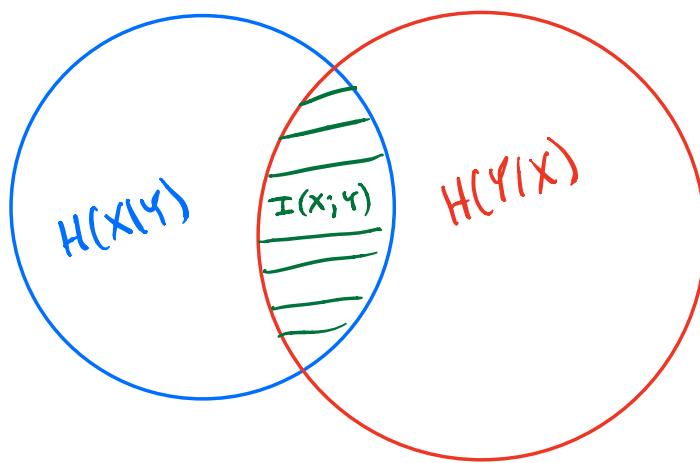
$$\begin{aligned} A &:= \{x \in \mathcal{X} \mid P_x(\{x\}) > 0\} \\ \Delta_A &:= \{(x, x) \mid x \in A^c\} \end{aligned}$$

and repeat continuous argument.

$$\begin{aligned} \textcircled{2} \quad I(X; Y) &= \sum_{x,y} P_{XY}(x, y) \log \frac{P_{XY}(x, y)}{P_X(x) P_Y(y)} \\ &= \sum_{x,y} P_{XY}(x, y) \log \frac{P_Y(y) P_{XY}(x|y)}{P_X(x) P_Y(y)} \\ &= \sum_{x,y} P_{XY}(x, y) \log \frac{1}{P_X(x)} - \sum_{x,y} P_{XY}(x, y) \log \frac{1}{P_{XY}(x|y)} \\ &= H(X) - H(X|Y) \end{aligned}$$

□

## Illustration

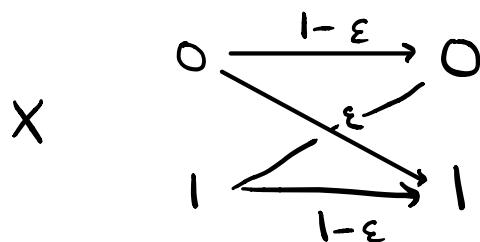


## Examples

① BSC  $P_x = x \sim \text{Ber}(\frac{1}{2}) \quad \text{II} \quad P_z = z \sim \text{Ber}(\varepsilon), \quad \varepsilon \in [0, \frac{1}{2}]$   
and

$$Y = X \oplus Z = \begin{cases} X \oplus 0 = X, & Z=0 \text{ w.p. } 1-\varepsilon \\ \pmod{2} & \\ X \oplus 1 = \bar{X}, & Z=1 \text{ w.p. } \varepsilon \end{cases}$$

$\Rightarrow$  So



is a binary symmetric channel w/ parameter  $\varepsilon$ .

$$I(X;Y) = \underbrace{H(Y)}_{(a)} - \underbrace{H(Y|X)}_{(b)}$$

② Start by computing pmf  $P_Y$

$$P_Y(0) = P_X(0) P_{Y|X}(0|0) + P_X(1) P_{Y|X}(0|1) = \frac{1}{2}(1-\varepsilon) + \frac{1}{2}\varepsilon = \frac{1}{2}$$

$$\Rightarrow Y \sim \text{Ber}(\frac{1}{2}) \Rightarrow H_b(\frac{1}{2}) = 1$$

⑥ Safe way

$$H(Y|X) = p_X(0) H(Y|X=0) + p_X(1) H(Y|X=1) = \dots$$

→ Instead, consider the following.

$$H(Y|X) = \sum_{x \in \{0,1\}} p_X(x) H(Y|X=x)$$

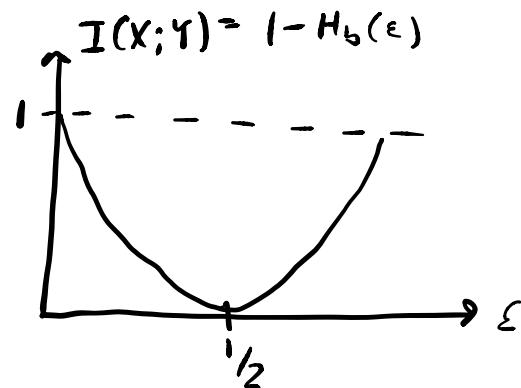
$$= \sum_{x \in \{0,1\}} p_X(x) H(X \oplus Z | X=x)$$

$$= \sum_{x \in \{0,1\}} p_X(x) \underbrace{H(X \oplus Z | X=x)}_{H(X \oplus Z | X=x) = H(Z|X=x)}$$

$$H(X \oplus Z | X=x) = H(Z|X=x)$$

$$= \sum_{x \in \{0,1\}} p_X(x) H(Z|X=x) = H(Z|X) = H(Z) = H_b(\varepsilon)$$

$$\Rightarrow I(X;Y) = 1 - H_b(\varepsilon)$$



## ② Bivariate Gaussian

$$(X, Y) \sim \mathcal{N} \left( \begin{bmatrix} 0 \\ 0 \end{bmatrix}, \begin{bmatrix} 1 & \rho \\ \rho & 1 \end{bmatrix} \right)$$

$$\begin{aligned} I(X; Y) &= h(X) + h(Y) - h(X, Y) \\ &= \frac{1}{2} \log(2\pi e) + \frac{1}{2} \log(2\pi e) - \frac{1}{2} \log((2\pi e)^2 (1-\rho)^2) \\ &= \frac{1}{2} \log \left( \frac{1}{(1-\rho)^2} \right) \end{aligned}$$