

Recap

Joint Entropy: $H(X, Y) = H(X) + \underbrace{H(Y|X)}_{\leq H(Y)} \rightarrow$ Conditional Entropy

Subadditivity : $H(X_1, \dots, X_m) = \sum_{i=1}^m H(X_i | X_1 \dots X_{i-1})$
 $\leq H(X_1) + \dots + H(X_m)$

Shearer's Lemma: Collection \mathcal{F} of subsets of $\{1, \dots, m\}$
Each $i \in [m]$ in at least t subsets

$$t \cdot H(X_1, \dots, X_m) \leq \sum_{S \in \mathcal{F}} H(X_S)$$

Mutual Information

$H(Y|X) \leq H(Y)$, but by how much?

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

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

$$= \underline{H(Y) + H(X) - H(X, Y)}$$

↳ Symmetric in X and Y

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

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

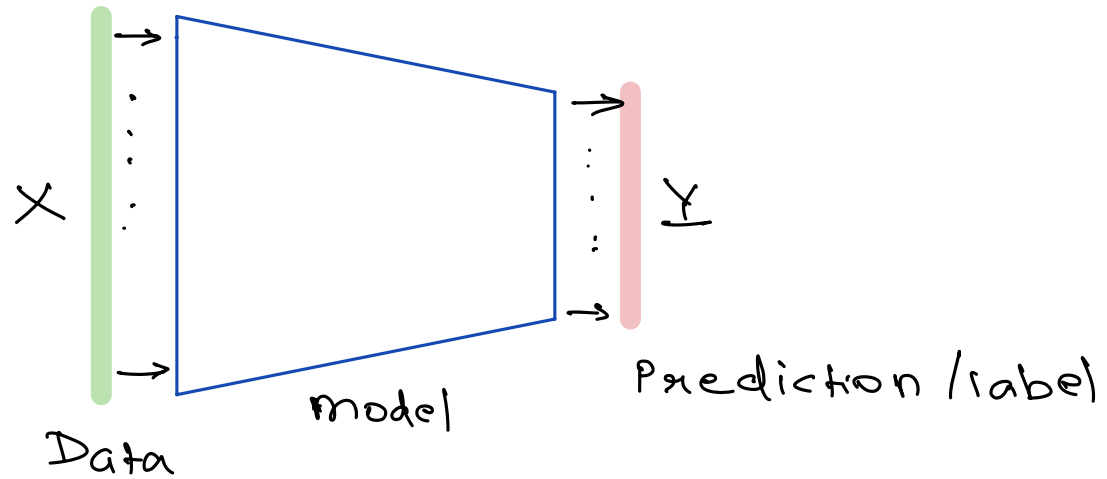
▶ $0 \leq I(X; Y) \leq H(X)$
When?

$(I(X; Y) \leq \min\{H(X), H(Y)\})$

Interesting when choosing Y (with other constraints)

- Error-correcting codes
- Statistical estimator
- Deep learning
- Communication Protocols
- Streaming algorithms

Eq. Learning representations from data



$$Z = \text{enc}(X)$$

encoding / embedding

maximize $I(Z; Y)$ (Mutual Information Coding)

maximize $I(Z; Y) - \beta I(Z; X)$ (Information Bottleneck)

E.g.

$$(X, Y, Z) = \begin{cases} 000 & \text{w.p. } 1/4 \\ 011 & \text{w.p. } 1/4 \\ 101 & \text{w.p. } 1/4 \\ 110 & \text{w.p. } 1/4 \end{cases}$$

$$X + Y + Z = 0$$

(mod 2)

$$I(X; Y) = 0$$

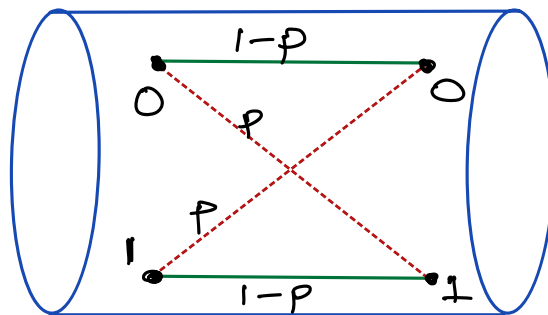
$$I(X; Y | Z) = \underbrace{H(X|Z)}_{=1} - \underbrace{H(X|Y, Z)}_{=0}$$

Conditioning does not necessarily reduce $I(X; Y)$

E.g.

0 w.p. $\frac{1}{2}$
1 w.p. $\frac{1}{2}$

X →



→ Y

Binary Symmetric Channel

$$I(X; Y) = ?$$

$$1 - H_2(p)$$

Chain rule

$$I(x_1, \dots, x_n; Y) = \sum_{i=1}^n I(x_i; Y \mid x_1, \dots, x_{i-1})$$

$$H(x_1, \dots, x_n) - H(x_1, \dots, x_n \mid Y)$$

$$\sum H(x_i \mid x_1, \dots, x_{i-1}) - H(x_i \mid Y, x_1, \dots, x_{i-1})$$

$$I(x_i; Y \mid x_1, \dots, x_{i-1})$$

Data Processing

▶ $I(X; Y) \geq I(X; g(Y))$ for any function g

Proof:

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

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

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

$$I(Y; Y) \geq I(Y; g(Y))$$

$$\begin{aligned} &= H(Y) \geq H(g(Y)) \\ &\quad - \cancel{H(g(Y)|Y)} \\ &\quad 0 \end{aligned}$$

Markov Chains

$$X \leftrightarrow Y \leftrightarrow Z$$

X, Z independent given Y

$$\triangleright I(X; Y) \geq I(X; Z)$$

Proof:

$$H(X) - H(X|Y)$$

$$H(X) - H(X|Y, Z)$$

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

$$I(X; Z|Y) = 0$$

by conditional independence

$$\geq H(X) - H(X|Z)$$

$$= I(X; Z)$$

Sufficient Statistics

For random variables X, Y :

$g(Y)$ is a **sufficient statistic** of Y for X if

$$I(X; Y) = I(X; g(Y))$$

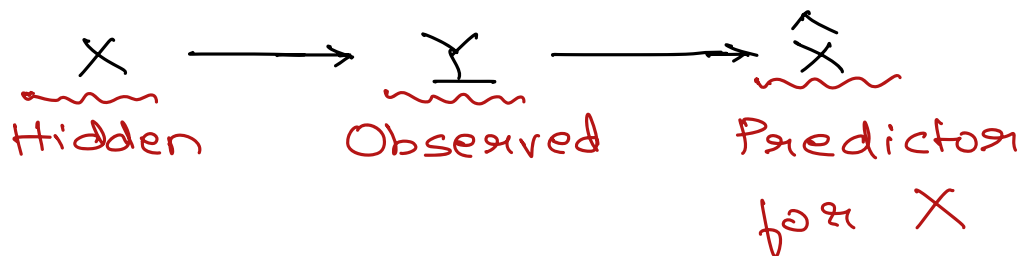
(No loss in
data processing)

Ex: $X = \begin{cases} p_1 & \text{w. p. } \frac{1}{2} \\ p_2 & \text{w. p. } \frac{1}{2} \end{cases} \rightarrow \begin{cases} Y = (Y_1, \dots, Y_n) \text{ iid} & Y_i = \begin{cases} 1 & \text{w.p. } p_1 \\ 0 & \text{w.p. } 1-p_1 \end{cases} \\ Y = (Y_1, \dots, Y_n) \text{ iid} & Y_i = \begin{cases} 1 & \text{w.p. } p_2 \\ 0 & \text{w.p. } 1-p_2 \end{cases} \end{cases}$

$$g(Y) = Y_1 + \dots + Y_n = \# \text{ of } 1\text{'s}$$

$$\text{Prove: } I(X; Y) = I(X; g(Y))$$

Fano's inequality



$$\text{Supp}(X) = \text{Supp}(\hat{X}) = \mathcal{X}$$

$$p_e = \mathbb{P}(\hat{X} \neq X)$$

\longrightarrow Data Processing

$$\triangleright H_2(p_e) + p_e \cdot \log(|\mathcal{X}| - 1) \geq H(X|\hat{X}) \geq H(X|Y)$$

Proof: $E = \mathbb{1}\{\hat{X} \neq X\}$ $\mathbb{P}(E=1) = p_e$

$$I(X, E | \hat{X}) = H(X|\hat{X}) - H(X|E, \hat{X})$$

$$I(X, E | \hat{X}) = H(E|\hat{X}) - \cancel{H(E|X, \hat{X})} \geq 0$$

$$\underbrace{H(E|\hat{X})}_{\leq H(E)} + \underbrace{H(X|E, \hat{X})}_{p_e H(X|\hat{X}, E=1) + (1-p_e) \cancel{H(X|\hat{X}, E=0)}} = H(X|\hat{X}) \quad \square$$