

Z. Alirezaei

$$X \Rightarrow H(X) = - \sum_{i=1}^m p_i \cdot \log p_i$$

$$\text{Prob}(X=x_i) = p_i, \quad i=1 \dots m$$

Theorem 2.1.B.

$$a) \quad 0 \stackrel{(i)}{\leq} H(X) \stackrel{(ii)}{\leq} \log m$$

"=" in (i) \Leftrightarrow X has a singleton dist., i.e.,
 $\exists x_i: \text{Prob}(X=x_i) = 1$

"=" in (ii) \Leftrightarrow X is uniformly dist., i.e.,

$$\text{Prob}(X=x_i) = \frac{1}{m} \quad \forall i=1, \dots, m$$

$$b) \quad 0 \stackrel{(i)}{\leq} H(X|Y) \stackrel{(ii)}{\leq} H(X)$$

"=" in (i) \Leftrightarrow $\text{Prob}(X=x_i | Y=y_j) = 1 \quad \forall i, j$
 with $\text{Prob}(X=x_i, Y=y_j) > 0$, i.e.,
 X is totally dependent on Y .

"=" in (ii) \Leftrightarrow X, Y are stoch. independent.

$$c) \quad H(X) \stackrel{(i)}{\leq} H(X, Y) \stackrel{(ii)}{\leq} H(X) + H(Y)$$

"=" in (i) \Leftrightarrow Y is totally dependent on X

"=" in (ii) \Leftrightarrow X, Y are stoch. independent.

$$d) \quad H(X|Y, Z) \leq \min \{ H(X|Y), H(X|Z) \}$$

Proof:

$$a) \quad (ii) \quad H(X) = -\sum_{i=1}^m p_i \cdot \log p_i = \sum_{i=1}^m p_i \log \frac{1}{p_i}$$

$$\stackrel{\text{Lemma 2.1.6}}{\leq} \log \left(\sum_{i=1}^m p_i \cdot \frac{1}{p_i} \right)$$

$$= \log \left(\sum_{i=1}^m 1 \right) = \log m$$

(i) obvious

b) Exercise

c) (i) By the chain rule, Th. 2.1.5.:

$$H(X, Y) = H(X) + \underbrace{H(Y|X)}_{\geq 0} \geq H(X)$$

equality follows from ≥ 0 b) (ii).

$$(ii) \quad 0 \stackrel{(b)}{\leq} H(X) - H(X|Y) \stackrel{\text{Th. 2.1.5}}{=} H(X) - [H(X, Y) - H(Y)]$$

$$\Rightarrow H(X) + H(Y) \geq H(X, Y)$$

with equality from b) (i).

d) Analogous to b) (ii').

ref. Cover & Thomas, Mathem.

Definition 2.1.9.

Let X, Y, Z discrete r.v.

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

is called mutual information or synergetropy of X and Y .

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

is called conditional mutual information of X and Y given Z .

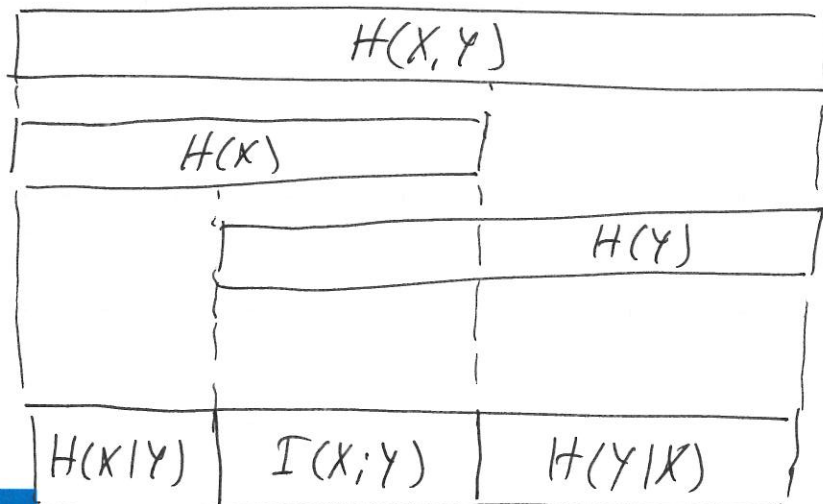
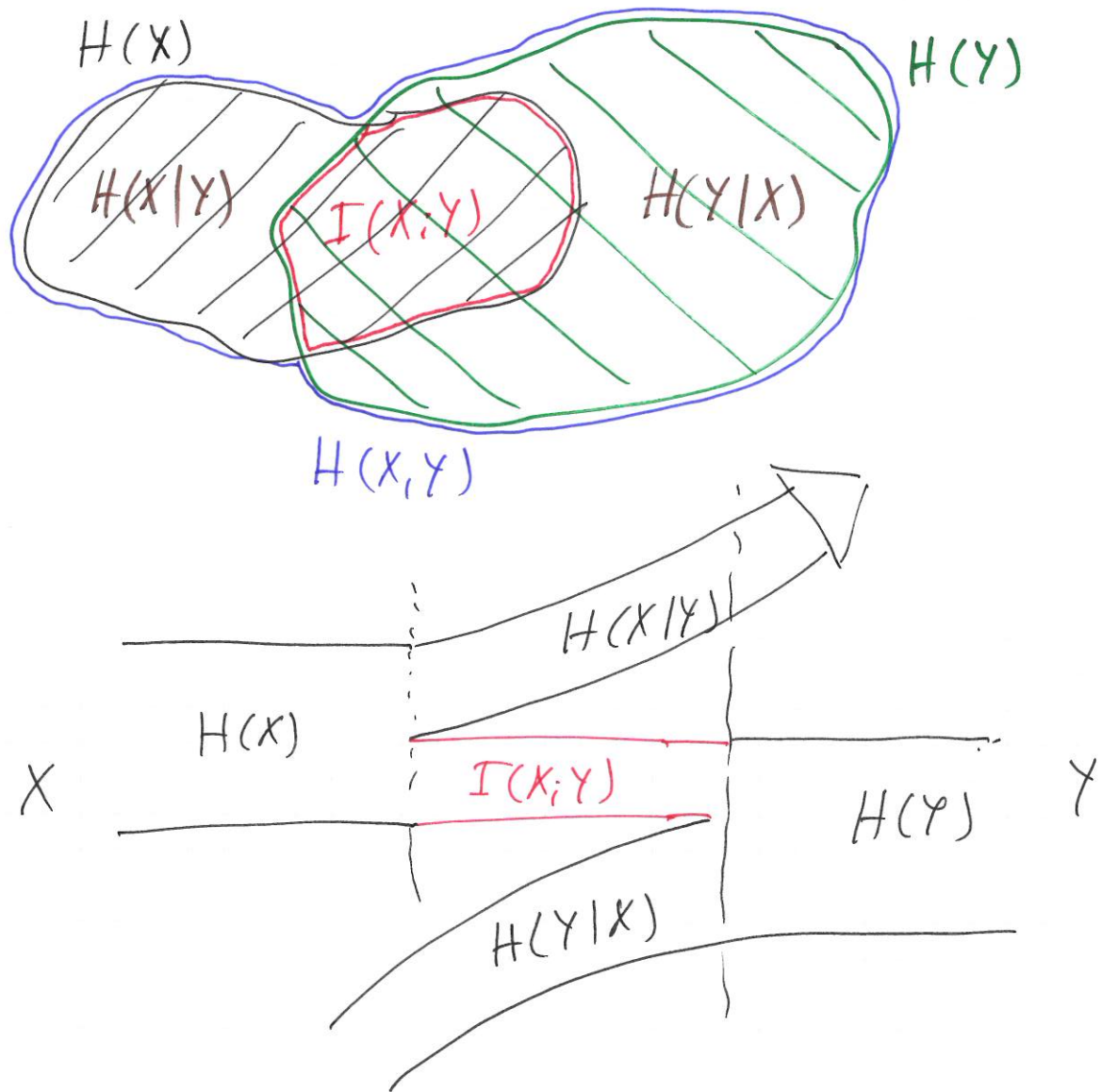
Interpretation: $I(X; Y)$ is the reduction in uncertainty about X when Y is given or the amount of information about X provided by Y .

see. Figures on the next page.

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

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

$$I(X; Y) = H(X) + H(Y) - H(X, Y)$$



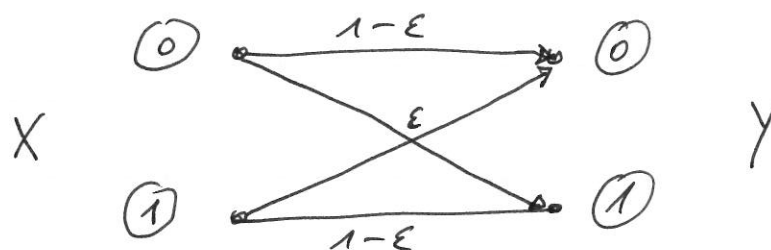
Note: by The 2.1.8 b) $I(X; Y) \geq 0$.

By definition it holds that:

$$\begin{aligned}
 I(X; Y) &= - \underbrace{\sum_i P(X_i) \cdot \log P(X_i)}_{= H(X)} + \sum_{i,j} P(X_i, Y_j) \cdot \log P(X_i | Y_j) \\
 &= - \sum_{i,j} P(X_i, Y_j) \cdot \log P(X_i) + \sum_{i,j} P(X_i, Y_j) \cdot \log P(X_i | Y_j) \\
 &= \sum_{i,j} P(X_i, Y_j) \cdot \log \frac{P(X_i | Y_j)}{P(X_i)} \\
 &= \sum_{i,j} P(X_i, Y_j) \cdot \log \frac{P(X_i, Y_j)}{P(X_i) \cdot P(Y_j)} \quad (*)
 \end{aligned}$$

which shows symmetry in X and Y .

Example 2.1.10. Binary symmetric channel (BSC)



symbol error with probability ϵ , $0 \leq \epsilon \leq 1$.

Hence:

$$P(Y=0 | X=0) = P(Y=1 | X=1) = 1-\epsilon$$

$$P(Y=0 | X=1) = P(Y=1 | X=0) = \epsilon$$

Assume that input symbols are uniformly distributed: $P(X=0) = P(X=1) = \frac{1}{2}$.

Then for the joint distributions:

$$P(X=0, Y=0) = P(Y=0 | X=0) \cdot P(X=0) = (1-\epsilon) \cdot \frac{1}{2}$$

⋮

X \ Y	0	1	
0	$\frac{1}{2}(1-\epsilon)$	$\frac{\epsilon}{2}$	$\frac{1}{2}$
1	$\frac{\epsilon}{2}$	$\frac{1}{2}(1-\epsilon)$	$\frac{1}{2}$
	$\frac{1}{2}$	$\frac{1}{2}$	

$$\text{Further: } P(X=0 | Y=0) = \frac{P(X=0, Y=0)}{P(Y=0)} = 1-\epsilon$$

$$P(X=1 | Y=1) = 1-\epsilon$$

$$P(X=0 | Y=1) = P(X=1 | Y=0) = \epsilon$$

For $\log = \log_2$

$$\Rightarrow H(X) = H(Y) = -\frac{1}{2} \cdot \log \frac{1}{2} - \frac{1}{2} \cdot \log \frac{1}{2} = 1 \text{ bit.}$$

$$H(X, Y) = 1 - (1-\epsilon) \cdot \log(1-\epsilon) - \epsilon \cdot \log \epsilon$$

$$H(X|Y) = -(1-\epsilon) \cdot \log(1-\epsilon) - \epsilon \cdot \log \epsilon$$

$$0 \leq I(X; Y) = 1 + (1-\epsilon) \cdot \log(1-\epsilon) + \epsilon \cdot \log \epsilon \leq 1$$

Def. 2.1.11 (Kullback-Leibler divergence)

Let $p = (p_1, \dots, p_m)$, $q = (q_1, \dots, q_m)$ be stoch. vectors.

$$D(p \parallel q) = \sum_{i=1}^m p_i \cdot \log \frac{p_i}{q_i}$$

is called KL divergence between p and q
(or relative entropy).

$D(p \parallel q)$ measures the divergence (distance, dissimilarity) between distributions p and q . However, it is not a metric, neither symmetric nor satisfies the triangle inequality.

It measures how difficult it is for p to pretend it were q .

Theorem 2.1.12.

- a) $D(p||q) \geq 0$ with "=" iff $p=q$.
- b) $D(p||q)$ is convex in the pair (p, q) .
- c) $I(X; Y) = D\left(\left(p(x_i, y_j)\right)_{i,j} \parallel \left(p(x_i) \cdot p(y_j)\right)_{i,j}\right)$

proof.

a) By definition and Cor. 2.1.8.

b) Use the log-sum inequality Lemma 2.1.7.
Let p, r and q, s be stoch. vectors.

For all $i = 1, \dots, m$ it holds that

$$\begin{aligned} & (\lambda p_i + (1-\lambda) \cdot r_i) \cdot \log \frac{\lambda p_i + (1-\lambda) \cdot r_i}{\lambda q_i + (1-\lambda) \cdot s_i} \\ & \leq \lambda p_i \log \frac{\lambda p_i}{\lambda q_i} + (1-\lambda) \cdot r_i \log \frac{(1-\lambda) \cdot r_i}{(1-\lambda) \cdot s_i} \end{aligned}$$

Summing over all $i = 1, \dots, m$, it follows $\forall \lambda \in [0, 1]$

$$D(\lambda p + (1-\lambda)r \parallel \lambda q + (1-\lambda)s) \leq \lambda D(p||q) + (1-\lambda) \cdot D(r||s)$$

c) By definition. ~~□~~

Note: $D(p||q) \neq D(q||p)$