

Review of Week 1

Chapter 1 : Information measures

rv X with values in $\{1, 2, 5\}$



smaller probability
⇒ more surprised

1.1 Surprisal and Entropy

surprisal is a function of p_x , say $s(p)$

Properties we require:

1) Monotonicity: $s(1) = 0$,

$s(x)$ increases monotonically as p decreases

2) Additivity: $s(p_x p_y) = s(p_x) + s(p_y)$

probability of observing
 $x=x$ and $y=y$
for independent X and Y

3) Normalisation: $s(\frac{1}{2}) = 1$

$$\Rightarrow s(p) = \log \frac{1}{p}$$

Define $S(X) = s(P_X) = \log \frac{1}{P_X(x)}$ as the

surprisal of the random experiment X

governed by the pmf $P_X(x)$

The entropy of X is $H(X) = \mathbb{E}[S(X)]$.

We have

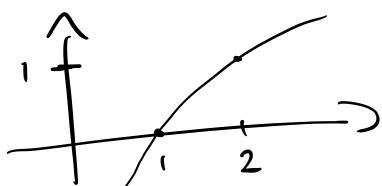
$$H(X) = \mathbb{E} \left[\log \frac{1}{P_X(x)} \right] = \sum_x P_X(x) \log \frac{1}{P_X(x)}$$

use convention that $0 \log 0 = 0$
 $0 \log \infty = 0$

$$\lim_{\epsilon \rightarrow 0} \epsilon \log \frac{1}{\epsilon} = \lim_{k \rightarrow \infty} \frac{\log k}{k} = 0$$

Prop. : $H(X) \geq 0$ with equality if and only if X is deterministic.

Proof: $\log \frac{1}{p} \geq 0$ for any $p \leq 1$



$H(X) = \sum p_x \log \frac{1}{p_x} \geq 0$ can be 0 only if $p_x \in \{0, 1\}$ for all x .

□

$-x \log x$ is a concave function of the pmf.

Entropy

$$f(t) = -t \log \frac{1}{t} \rightarrow H(X) = \sum_x f(p_x)$$

$$= -t \log t$$

$f(t)$ is strictly concave since

$$f'(t) = -\log t - \log e$$

$$f''(t) = -\frac{\log e}{t}$$

Prop: $H(X) \leq \log |X|$ with equality iff X is uniform.

Proof for $X = \{0, 1\}$:

Take pmf $\{p, 1-p\}$ for $p \in [0, 1]$. We write

$$\begin{aligned} H(X) &= f(p) + f(1-p) \\ &= \frac{1}{2}(f(p) + f(1-p)) + \frac{1}{2}(f(p) + f(1-p)) \\ &\stackrel{\text{Jensen's inequality}}{\leq} f\left(\frac{1}{2}p + \frac{1}{2}(1-p)\right) \cdot 2 \\ &= f\left(\frac{1}{2}\right) \cdot 2 \quad \leftarrow \text{recall } f(t) = -t \log t \\ &= -\frac{1}{2} \log \frac{1}{2} \cdot 2 = 1 \end{aligned}$$

1.2 Conditional entropy and mutual information

Joint entropy

Two rvs X and Y with joint pmf $P_{XY}(x, y) = p_{xy}$

$$H(XY) = E[S(X, Y)] = \sum_{(x,y) \in X \times Y} p_{xy} \log \frac{1}{p_{xy}}$$

$$= \sum_x \sum_y p_{xy} \log \frac{1}{p_{xy}}$$

Note : if X and Y are independent, then

$$\begin{aligned} H(XY) &= E[S(X, Y)] = E[S(X) + S(Y)] \\ &= E[S(X)] + E[S(Y)] \\ &= H(X) + H(Y) \end{aligned}$$

Example where $H(XY) = H(X) = H(Y)$?

$$X=Y, \text{ e.g. } P_{xy} = \sum_{\substack{x=y \\ x,y}} 1 - \sum_{\substack{x \neq y \\ x,y}} 0$$

$$\text{with } x, y \in \{0, 1\}$$

$$P_{y|x} = \sum_{\substack{y=x \\ y \neq x}} 1 - \sum_{\substack{y \neq x \\ y \neq x}} 0$$

Conditional entropy of Y given X is defined as

$$H(Y|X) = E\left[\log \frac{1}{P_{Y|X}}\right] = \sum_x \sum_y p_{xy} \log \frac{1}{P_{Y|X}}$$

$$= \sum_x p_x \underbrace{\sum_y p_{xy} \log \frac{1}{P_{Y|X}}}_{\text{---}}$$

$$= \overbrace{H(Y_x)} = H(Y|X=x)$$

$$= \underline{\mathbb{E}[H(Y_x)]}$$

$$\Rightarrow 0 \leq H(Y|X) \leq \log |Y|$$

Prop. (Chain rule) $H(XY) = H(X) + H(Y|X)$

Proof. $H(XY) = \sum_x \sum_y p_{xy} \log \frac{1}{p_{xy}} \quad \leftarrow p_{xy} = p_x p_{y|x}$

$$= \underbrace{\sum_x \sum_y p_{xy} \log \frac{1}{p_x}}_{H(X)} + \underbrace{\sum_x \sum_y p_{xy} \log \frac{1}{p_{y|x}}}_{H(Y|X)}$$

□

Often $H(Y|X) = H(XY) - H(X)$ is used as the definition of $H(Y|X)$!

Prop : $H(Y|X) \leq H(Y) \iff H(XY) \leq H(X) + H(Y)$

↑
sub-additivity of
entropy

Proof: Only need to show

$$H(Y|X) = \mathbb{E} \left[\sum_y p_{y|x} \log \frac{1}{p_{y|x}} \right]$$

$$\begin{aligned}
 & H(Y_X) \\
 &= \sum_Y \overbrace{E[f(P_{Y|X})]}^{\text{Jensen's inequality}} \\
 &\leq \sum_Y f(\underbrace{E[P_{Y|X}]}_{P_Y}) \\
 &= \sum_Y f(P_Y) = H(Y)
 \end{aligned}$$

Strong sub-additivity : $H(X|YZ) \leq H(X|Z)$

Mutual information

$$\begin{aligned}
 \text{Define } I(X:Y) &= H(X) + H(Y) - H(XY) \\
 &= H(Y) - H(Y|X) \\
 &\quad \xrightarrow{\text{from chain rule}} H(X) - H(X|Y)
 \end{aligned}$$

Prop. $0 \leq I(X:Y) \leq \min \{ \log |X|, \log |Y| \}$

Conditional mutual information

$$I(X:Y|Z) = \sum_z P_z \overbrace{I(X:Y|Z=z)}^{\uparrow \dots}$$

mutual information
evaluated for $P_{Y|Z}$

$$\begin{aligned} &= H(Y|Z) - H(Y|XZ) \\ &= H(X|Z) - H(X|YZ) \end{aligned}$$

chain rule : $I(X:YZ) = I(X:Y) + I(X:Z|Y)$

If $X-Z-Y$ form a Markov chain, then

$$I(X:Y|Z) = 0.$$

data-processing
inequality

Prop. : Let $X-Y-Z$ form a Markov chain.

Then $I(X:Y) \geq I(X:Z)$.

Proof: $I(X:Y) = I(X:YZ)$

due to chain rule and Markov condition

It remains to show $I(X:YZ) \geq I(X:Z)$

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

$$\Leftrightarrow H(X|Z) \geq H(X|YZ)$$

↑
strong sub-additivity

Relative entropy

Def. Let P and Q be two pmf on an alphabet \mathcal{X} .
 The relative entropy of P with regards to Q is
 defined as

$$D(P \parallel Q) = \sum_{\substack{x \in \mathcal{X} \\ P(x) > 0}} P(x) \log \frac{P(x)}{Q(x)}$$

If $P(x) > 0 \Rightarrow Q(x) > 0$ for all $x \in \mathcal{X}$, otherwise

$$D(P \parallel Q) = +\infty.$$

We can see relative entropy as an expectation
 of the log-likelihood ratio

$$Z(x) = \log \frac{P(x)}{Q(x)}$$

$$D(P \parallel Q) = \mathbb{E}_{x \sim P}[Z(x)].$$

$$\text{Prop: } I(X:Y) = D(P_{XY} \parallel P_X \times P_Y)$$

$$H(X) = \log |\mathcal{X}| - D(P_X \parallel U_X)$$

$$H(X|Y) = \log |\mathcal{Z}|$$

$$- D(P_{XY} \parallel U_X \times P_Y)$$

$$\begin{aligned} \uparrow & \text{uniform distribution} \\ U_X(x) &= \frac{1}{|\mathcal{Z}|} \text{ for all } x \end{aligned}$$

Proof: Homework

Prop : $D(P||Q) \geq 0$ with equality only if $P=Q$.

Proof : $g(x) = -\log x$ is strictly convex

$$\begin{aligned}
 D(P||Q) &= \sum_x \underbrace{P(x)}_{P(x)>0} \log \frac{P(x)}{Q(x)} \\
 &= \sum_{\substack{x \\ P(x)>0}} P(x) g\left(\frac{Q(x)}{P(x)}\right) \\
 &\geq g\left(\sum_{\substack{x \\ P(x)>0}} P(x) \frac{Q(x)}{P(x)}\right) \\
 &= g\left(\sum_{\substack{x \\ P(x)>0}} Q(x)\right) \\
 &\geq g\left(\sum_x Q(x)\right) = g(1) = 0 \quad \square
 \end{aligned}$$

Prop (Data-processing inequality)

Let P_X and Q_X be two pmfs on \mathcal{X} (the input distributions). Let $P_{Y|X}$ be a conditional pmf (the channel / noise channel).

Define $P_Y(y) = \sum_{x \in \mathcal{X}} \underbrace{P_{Y|X}(y|x)}_{P_{XY}} \underbrace{P_X(x)}_{P_X}$

$Q_Y(y) = \sum_{x \in \mathcal{X}} \underbrace{P_{Y|X}(y|x)}_{Q_{XY}} \underbrace{Q_X(x)}_{Q_X}$

(output distributions)

Thus, the data-processing inequality (DPI) states

$$D(P_x \| Q_x) \geq D(P_y \| Q_y)$$

Proof: $D(P_{xy} \| Q_{xy}) - D(P_x \| Q_x)$ ②

$$= \left(\sum_{x,y} P_{xy} \log \frac{P_{xy}}{Q_{xy}} \right) - \left(\sum_x P_x \log \frac{P_x}{Q_x} \right)$$

$$= \sum_{x,y} P_{xy} \left(\log \frac{P_{xy}}{Q_{xy}} - \log \frac{P_x}{Q_x} \right)$$

$$= \sum_x P_x \sum_y P_{yx} \log \frac{P_{yx}}{Q_{yx}}$$

$$= \sum_x P_x D(P_{y|x=x} \| Q_{y|x=x}) = 0 !$$

$D(P_{xy} \| Q_{xy}) - D(P_y \| Q_y)$ ①

$$= \sum_y P_y D(P_{x|y=y} \| Q_{x|y=y}) \geq 0$$

$$\textcircled{1} - \textcircled{2} \geq 0$$

"

$$D(P_x \| Q_x) - D(P_y \| Q_y)$$

□