

Information Theory and Learning

Information Theory and Games

Basic Concepts

Definition *Entropy is a measure of uncertainty of a random variable. Let X be a discrete random variable with alphabet \mathcal{X} .*

$$p(x) = \Pr[X = x], \quad \text{where } x \in \mathcal{X}.$$

The entropy $H(X)$ of the discrete random variable X is defined as

$$\begin{aligned} H(X) &= \mathbb{E}_p \lg \frac{1}{p(X)} \\ &= - \sum_{x \in \mathcal{X}} p(x) \lg p(x). \end{aligned}$$

Facts

1. $H(X) \geq 0$. Entropy is always nonnegative. $0 \leq p(x) \leq 1$; $-\lg p(x) \geq 0$. Hence, $E_p \lg(1/p(x)) \geq 0$.)
2. $H(X) \leq \lg |\mathcal{X}|$. Consider the uniform distribution $u(x)$. $\forall_{x \in \mathcal{X}} u(x) = 1/|\mathcal{X}|$. $H(u) = \sum_x (1/|\mathcal{X}|) \lg |\mathcal{X}| = \lg |\mathcal{X}|$.

Game Theory & Learning: LECTURE 11

3. $H(X)$ = Average number of bits required to encode the discrete random variable X .

Joint & Conditional Entropy

(X, Y) = A pair of discrete random variables with joint distribution $p(x, y)$.

Joint Entropy =

$$\begin{aligned} H(X, Y) &= \mathbb{E}_p \lg \frac{1}{p(X, Y)} \\ &= - \sum_{x \in \mathcal{X}} \sum_{y \in \mathcal{Y}} p(x, y) \lg p(x, y). \end{aligned}$$

Conditional Entropy =

$$\begin{aligned} H(Y|X) &= \mathbb{E}_p \lg \frac{1}{p(Y|X)} \\ &= - \sum_{x \in \mathcal{X}} \sum_{y \in \mathcal{Y}} p(x, y) \lg p(y|x) \\ &= - \sum_{x \in \mathcal{X}} p(x) \sum_{y \in \mathcal{Y}} p(y|x) \lg p(y|x) \\ &= - \sum_{x \in \mathcal{X}} p(x) H(Y|x). \end{aligned}$$

Chain Rule

$$\begin{aligned} p(X, Y) &= p(X) p(Y|X) && \text{Bayes' Rule} \\ \lg p(X, Y) &= \lg p(X) + \lg p(Y|X) \\ \lg \frac{1}{p(X, Y)} &= \lg \frac{1}{p(X)} + \lg \frac{1}{p(Y|X)} \\ \mathbb{E}_p \lg \frac{1}{p(X, Y)} &= \mathbb{E}_p \lg \frac{1}{p(X)} + \mathbb{E}_p \lg \frac{1}{p(Y|X)} && \text{Linearity of Expectation} \\ H(X, Y) &= H(X) + H(Y|X). \end{aligned}$$

Corollary

Game Theory & Learning: LECTURE 11

1. $H(X, Y|Z) = H(X|Z) + H(Y|X, Z)$.
2. $H(X) + H(Y|X) = H(Y) + H(X|Y)$
 $\Rightarrow H(X) - H(X|Y) = H(Y) - H(Y|X)$.
3. Note that $H(X|Y) \neq H(Y|X)$.

Relative Entropy & Mutual Information

Definition **Relative Entropy**—Also, **Kullback-Liebler Distance** between two probability mass functions $p(x)$ and $q(x)$.

$$D(p||q) = \mathbb{E}_p \lg \frac{p(x)}{q(x)} = - \sum_x p(x) \lg \frac{q(x)}{p(x)}.$$

Note that $D(p||p) = 0$. If $u(x) = \frac{1}{|\mathcal{X}|}$, for all x . Then $D(p||u)$ is

$$D(p||u) = - \sum p(x) \lg \frac{1}{p(x)} + \sum p(x) \lg |\mathcal{X}| = \lg |\mathcal{X}| - H(X).$$

Definition **Mutual Information**

Let X and Y be two discrete random variables with a joint probability mass function $p(x, y)$, and with marginal probability mass functions

$$p(x) = \sum_{y \in \mathcal{Y}} p(x, y) \quad \& \quad p(y) = \sum_{x \in \mathcal{X}} p(x, y).$$

Mutual Information,

$$\begin{aligned} I(X; Y) &= D\left(p(x, y) \parallel p(x)p(y)\right) \\ &= \mathbb{E}_{p(x, y)} \lg \frac{p(x, y)}{p(x)p(y)} \\ &= - \sum_{x \in \mathcal{X}} \sum_{y \in \mathcal{Y}} p(x, y) \lg \frac{p(x, y)}{p(x)p(y)} \end{aligned}$$

Game Theory & Learning: LECTURE 11

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

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

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

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

Chain Rules for Entropy, Relative Entropy and Mutual Information

$$\begin{aligned} H(X_1, X_2, \dots, X_n) &= H(X_1) + H(X_2|X_1) + \dots \\ &\quad + H(X_n|X_1, \dots, X_{n-1}) \\ &= \sum_{i=1}^n H(X_i|X_1, \dots, X_{i-1}). \end{aligned}$$

$$\begin{aligned} I(X_1, X_2, \dots, X_n; Y) &= H(X_1, \dots, X_n) - H(X_1, \dots, X_n|Y) \\ &= \sum_{i=1}^n H(X_i|X_1, \dots, X_i) - \sum_{i=1}^n H(X_i|X_1, \dots, X_i, Y) \\ &= \sum_{i=1}^n H(X_i|X_1, \dots, X_i) - H(X_i|X_1, \dots, X_i, Y) \\ &= \sum_{i=1}^n I(X_i; Y|X_1, \dots, X_{i-1}). \end{aligned}$$

Game Theory & Learning: LECTURE 11

$$\begin{aligned} & D\left(p(x, y) \parallel q(x, y)\right) \\ &= \sum_x \sum_y p(x, y) \lg \frac{p(x, y)}{q(x, y)} \\ &= \sum_x \sum_y p(x, y) \lg \frac{p(x) p(y|x)}{q(x) q(y|x)} \\ &= \sum_x \sum_y p(x, y) \lg \frac{p(x)}{q(x)} + \sum_x \sum_y p(x, y) \lg \frac{p(y|x)}{q(y|x)} \\ &= \sum_x p(x) \lg \frac{p(x)}{q(x)} + \sum_y p(y|x) \lg \frac{p(y|x)}{q(y|x)} \\ &= D\left(p(x) \parallel q(x)\right) + D\left(p(y|x) \parallel q(y|x)\right). \end{aligned}$$

Information Inequality

$$\begin{aligned} -D(p \parallel q) &= \sum_x p(x) \lg \frac{q(x)}{p(x)} \quad \lg \text{ is a concave function} \\ &\leq \lg \sum_x p(x) \frac{q(x)}{p(x)} \leq \lg \sum_x q(x) = \lg 1 = 0. \end{aligned}$$

Theorem $D(p \parallel q) \geq 0$ (with equality iff $p(x) = q(x)$ for all x .)

Corollary

$$I(X; Y) = D\left(p(x, y) \parallel p(x)p(y)\right) \geq 0,$$

(with equality iff X and Y are independent, i.e., $p(x, y) = p(x)p(y)$ for all x and y .)

Let $u(x) = \frac{1}{|\mathcal{X}|}$.

$$D(p \parallel u) = \lg |\mathcal{X}| - H(X) \geq 0.$$

Game Theory & Learning: LECTURE 11

Hence,

$$H(X) \leq \lg |\mathcal{X}|,$$

(with equality iff X has a uniform distribution over \mathcal{X} .)

$$I(X; Y) = H(X) - H(X|Y) \geq 0.$$

Theorem

$$H(X|Y) \leq H(X).$$

Conditioning reduces entropy.

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

Corollary

$$H(X_1, \dots, X_n) \leq \sum_{i=1}^n H(X_i),$$

with equality iff X_i 's are independent.

Stationary Markov Process

- *Markovian*

$$Pr[X_n | X_1, \dots, X_i] = Pr[X_n | X_i], \quad i \leq n.$$

- *Stationary*

$$Pr[X_n | X_1, \dots, X_i] = Pr[X_{n+1} | X_2, \dots, X_{i+1}].$$

Game Theory & Learning: LECTURE 11

$$\begin{aligned} H(X_n|X_1) &\geq H(X_n|X_1, X_2) && \text{conditioning reduces entropy} \\ &= H(X_n|X_2) && \text{Markov} \\ &= H(X_{n-1}|X_1) && \text{Stationary .} \end{aligned}$$

2nd Law of Thermodynamics

Theorem *Conditional entropy $H(X_n|X_1)$ increases with time n for a stationary Markov process.*

Relative entropy $D(\pi_n || \pi'_n)$ decreases with time n .

Let π_n and π'_n be two postulated probability distributions on the state space of a Markov Process. At time $n + 1$, the distribution changes to π_{n+1} and π'_{n+1} , governed by the transition probabilities $r(x_n, x_{n+1})$.

Thus

$$\begin{aligned} p(x_n, x_{n+1}) &= p(x_n)r(x_n, x_{n+1}) \\ &= p(x_n)p(x_{n+1}|x_n) \end{aligned}$$

similarly,

$$\begin{aligned} q(x_n, x_{n+1}) &= q(x_n)r(x_n, x_{n+1}) \\ &= q(x_n)q(x_{n+1}|x_n) \end{aligned}$$

Thus, we have

$$\begin{aligned} &D\left(p(x_n, x_{n+1}) \parallel q(x_n, x_{n+1})\right) \\ &= D\left(p(x_n) \parallel q(x_n)\right) + D\left(p(x_{n+1}|x_n) \parallel q(x_{n+1}|x_n)\right) \\ &= D\left(p(x_n) \parallel q(x_n)\right). \end{aligned}$$

And

$$D\left(p(x_n, x_{n+1}) \parallel q(x_n, x_{n+1})\right)$$

Game Theory & Learning: LECTURE 11

$$\begin{aligned} &= D\left(p(x_{n+1}) \parallel q(x_{n+1})\right) + D\left(p(x_n|x_{n+1}) \parallel q(x_n|x_{n+1})\right) \\ &\geq D\left(p(x_{n+1}) \parallel q(x_{n+1})\right). \end{aligned}$$

We conclude that

$$D\left(p(x_n) \parallel q(x_n)\right) \geq D\left(p(x_{n+1}) \parallel q(x_{n+1})\right).$$

Thus the relative entropy for this system must decrease:

$$\begin{aligned} D(\pi_1 \parallel \pi'_1) &\geq D(\pi_2 \parallel \pi'_2) \geq \dots \\ &\geq D(\pi_n \parallel \pi'_n) \geq D(\pi_{n+1} \parallel \pi'_{n+1}) \geq \dots \rightarrow 0. \end{aligned}$$