

CSE 5523: Lecture Notes 6

Information Theory

Contents

6.1	A Formal definition of information [Shannon, 1948]	1
6.2	Entropy	2
6.3	Cross entropy and Kullback-Leibler (KL) divergence	2
6.4	Conditional entropy and mutual information	3

6.1 A Formal definition of information [Shannon, 1948]

We can formalize the contribution of learning as **information**.

In this sense, information about a distribution makes it more predictable.

For example, if you think a promotion is $[.50, .50]$, then learn you got it, your distribution is $[0, 1]$.

This is a **bit** of information: the difference between no knowledge and certainty of a Bernoulli trial.

They are the bits you'd use to optimally encode probability-weighted outcomes of a distribution.

For example, with a distribution $[0.5, 0.25, 0.125, 0.125]$ the optimal encoding is *not* this:

event	freq.	code	cost
A	500	00	$500 \times 2 = 1000$ bits
B	250	01	$250 \times 2 = 500$ bits
C	125	10	$125 \times 2 = 250$ bits
D	125	11	$125 \times 2 = 250$ bits
	1000		2000 bits

but rather *this*, with variable length tokens, inversely proportional to the log of the probability:

event	freq.	code	cost
A	500	0	$500 \times 1 = 500$ bits
B	250	10	$250 \times 2 = 500$ bits
C	125	110	$125 \times 3 = 375$ bits
D	125	111	$125 \times 3 = 375$ bits
	1000		1750 bits

To encode 1000 outcomes, you use only 1750 bits instead of 2000!

If an event *always* happens, you give it zero bits – the receiver *already knows* the outcome.

If an event *never* happens, you don't give it a code – you can't send it, but you *won't need to*.

Formally, the information (in bits) of an event is the negative log of its probability:

$$I_{p_1, p_2, \dots}(x) = -\log_2 P_{p_1, p_2, \dots}(x)$$

If an event has probability 1, it has information 0 (use most efficient code imaginable: nothing!).

If an event has probability 0, it has information ∞ (use least efficient code imaginable: everything!).

This is called **self-information** or **surprisal**. It's the information *of the event, given a distribution*.

6.2 Entropy

The expected information of a distribution is then:

$$\begin{aligned} H(X) &= H(P_{p_1, p_2, \dots}(X)) = E_{x \sim P_{p_1, p_2, \dots}(X)} I_{p_1, p_2, \dots}(x) \\ &= E_{x \sim P_{p_1, p_2, \dots}(X)} (-\log_2 P_{p_1, p_2, \dots}(x)) && \text{definition of self-information} \\ &= - \sum_{x \in X} P_{p_1, p_2, \dots}(x) \log_2 P_{p_1, p_2, \dots}(x) && \text{definition of expected value} \end{aligned}$$

This is also called the **entropy** (from Greek 'entropia' roughly meaning 'disorder' or 'chaos').

Indeed, expecting lots of information indicates chaos; expecting no information indicates order.

(Why abbreviate entropy as H? It's a capital Greek eta η , pronounced 'eh', as in 'eh'ntropy.)

And here's the entropy of our promotion distribution, before and after finding out:

$$\begin{aligned} H([.5, .5]) &= .5 \cdot 1 + .5 \cdot 1 = 1 \\ H([0, 1]) &= 0 \cdot \infty + 1 \cdot 0 = 0 \end{aligned}$$

6.3 Cross entropy and Kullback-Leibler (KL) divergence

In defining loss functions for parameters, it's useful to quantify how wrong a distribution Q is.

First, using distribution Q on data distributed according to P has the following information:

$$\begin{aligned} H(P, Q) &= -E_{x \sim P(X)} \log_2 Q(x) \\ &= - \sum_{x \in X} P(x) \log_2 Q(x) && \text{definition of expected value} \end{aligned}$$

(Here we assume P and Q share the same event space, but are not in the same probability space.)

This is called **cross entropy**.

Using a different distribution is *always worse* (optimality of maximum likelihood estimation).

The loss in expected information from using Q instead of P on data distributed according to P is:

$$D_{\text{KL}}(P \parallel Q) = H(P, Q) - H(P)$$

$$\begin{aligned}
&= \left(-\mathbb{E}_{x \sim P(X)} \log_2 Q(x) \right) - \left(-\mathbb{E}_{x \sim P(X)} \log_2 P(x) \right) && \text{definition of (cross) entropy} \\
&= \left(-\sum_{x \in X} P(x) \log_2 Q(x) \right) - \left(-\sum_{x \in X} P(x) \log_2 P(x) \right) && \text{definition of expected value} \\
&= -\sum_{x \in X} P(x) \left(\log_2 Q(x) - \log_2 P(x) \right) && \text{distributive axiom} \\
&= -\sum_{x \in X} P(x) \log_2 \frac{Q(x)}{P(x)} && \text{addition of logs}
\end{aligned}$$

This is called **Kullback-Leibler (KL) divergence** or **relative entropy**.

It's zero (log of one) when the distributions match, and positive when they don't.

6.4 Conditional entropy and mutual information

Sometimes it's valuable to see how much information two variables share.

First, loss in expected information from using $P(X)$ instead of $P(X, Y)$ on distribution $P(X, Y)$ is:

$$\begin{aligned}
H(Y|X) &= D_{\text{KL}}(P(X, Y) \parallel P(X)) \\
&= \left(-\mathbb{E}_{x,y \sim P(X,Y)} \log_2 P(x) \right) - \left(-\mathbb{E}_{x,y \sim P(X,Y)} \log_2 P(x, y) \right) && \text{def. of KL divergence} \\
&= \left(-\sum_{x,y \in X \times Y} P(x, y) \log_2 P(x) \right) - \left(-\sum_{x,y \in X \times Y} P(x, y) \log_2 P(x, y) \right) && \text{def. of expected value} \\
&= -\sum_{x,y \in X \times Y} P(x, y) \left(\log_2 P(x) - \log_2 P(x, y) \right) && \text{distributive axiom} \\
&= -\sum_{x,y \in X \times Y} P(x, y) \log_2 \frac{P(x)}{P(x, y)} && \text{addition of logs} \\
&= \sum_{x,y \in X \times Y} P(x, y) \log_2 \frac{P(x, y)}{P(x)} && \text{log of inverse} \\
&= \sum_{x,y \in X \times Y} P(x, y) \log_2 P(y|x) && \text{addition of logs}
\end{aligned}$$

(Here we assume $P(X, Y)$ and $P(X)$ are in the same probability space.)

This is called **conditional entropy**.

When $P(X)$ is predictive of $P(X, Y)$ (e.g. X and Y are correlated), this loss is small; otherwise big.

Loss in expected information from using $P(X) \cdot P(Y)$ instead of $P(X, Y)$ on distribution $P(X, Y)$ is:

$$\begin{aligned}
I(X; Y) &= D_{\text{KL}}(P(X, Y) \parallel P(X) \cdot P(Y)) \\
&= \left(-\mathbb{E}_{x,y \sim P(X,Y)} \log_2 P(x) \cdot P(y) \right) - \left(-\mathbb{E}_{x,y \sim P(X,Y)} \log_2 P(x, y) \right) && \text{def. of KL divergence}
\end{aligned}$$

$$\begin{aligned}
&= \left(- \sum_{x,y \in X \times Y} P(x,y) \log_2 P(x) \cdot P(y) \right) - \left(- \sum_{x,y \in X \times Y} P(x,y) \log_2 P(x,y) \right) && \text{def. of expected value} \\
&= - \sum_{x,y \in X \times Y} P(x,y) \left(\log_2 P(x) \cdot P(y) - \log_2 P(x,y) \right) && \text{distributive axiom} \\
&= - \sum_{x,y \in X \times Y} P(x,y) \log_2 \frac{P(x) \cdot P(y)}{P(x,y)} && \text{addition of logs} \\
&= - \sum_{x,y \in X \times Y} P(x,y) \left(\log_2 P(x) + \log_2 \frac{P(y)}{P(x,y)} \right) && \text{addition of logs} \\
&= \left(- \sum_{x,y \in X \times Y} P(x,y) \log_2 P(x) \right) + \left(- \sum_{x,y \in X \times Y} P(x,y) \log_2 \frac{P(y)}{P(x,y)} \right) && \text{distributive axiom} \\
&= \left(- \sum_{x,y \in X \times Y} P(x,y) \log_2 P(x) \right) - \left(- \sum_{x,y \in X \times Y} P(x,y) \log_2 \frac{P(x,y)}{P(y)} \right) && \text{log of inverse} \\
&= \left(- \sum_{x \in X} P(x) \log_2 P(x) \right) - \left(- \sum_{x,y \in X \times Y} P(x,y) \log_2 \frac{P(x,y)}{P(y)} \right) && \text{marginalization} \\
&= H(X) - H(X|Y) && \text{def. of (conditional) entropy}
\end{aligned}$$

This is called **mutual information**. Unlike conditional entropy, it is symmetric.

When X and Y are independent, it's low; otherwise it's high.

(Note this can differ from conditional entropy, e.g. if X is more fine-grained than Y .)

References

[Shannon, 1948] Shannon, C. (1948). A mathematical theory of communication. *Bell System Technical Journal*, 27:379–423, 623–656.