

KL divergence or relative entropy

Two pmfs $p(x)$ and $q(x)$:

$$D(p \parallel q) = \sum_{x \in \mathcal{X}} p(x) \log \frac{p(x)}{q(x)} \quad (5)$$

Say $0 \log \frac{0}{q} = 0$, otherwise $p \log \frac{p}{0} = \infty$.

$$D(p \parallel q) = E_p \left(\log \frac{p(X)}{q(X)} \right) \quad (6)$$

$$I(X; Y) = D(p(x, y) \parallel p(x) p(y)) \quad (7)$$

- Measure of how different two probability distributions are
- The average number of bits that are wasted by encoding events from a distribution p with a code based on a not-quite-right distribution q .
- $D(p \parallel q) \geq 0$; $D(p \parallel q) = 0$ iff $p = q$
- Not a metric: not commutative, doesn't satisfy triangle equality

[Slide on $D(p\|q)$ vs $D(q\|p)$]

Cross entropy

- Entropy = uncertainty
- Lower entropy = determining efficient codes
= knowing the structure of the language =
good measure of model quality
- Entropy = measure of surprise
- How surprised we are when w follows h is
pointwise entropy:

$$H(w|h) = -\log_2 p(w|h)$$

$$p(w|h) = 1? \quad p(w|h) = 0$$

- Total surprise:

$$\begin{aligned} H_{\text{total}} &= -\sum_{j=1}^n \log_2 m(w_j | w_1, w_2, \dots, w_{j-1}) \\ &= -\log_2 m(w_1, w_2, \dots, w_n) \end{aligned}$$

Formalizing through cross-entropy

- Our model of language is $q(x)$. How good a model is it?
- Idea: use $D(p \parallel q)$, where p is the correct model.
- Problem: we don't know p .
- But we know roughly what it is like from a corpus
- Cross entropy:

$$H(X, q) = H(X) + D(p \parallel q) \quad (8)$$

$$= - \sum_x p(x) \log q(x)$$

$$= E_p \left(\log \frac{1}{q(x)} \right) \quad (9)$$

- Cross entropy of a language $L = (X_i) \sim p(x)$ according to a model m :

$$H(L, m) = - \lim_{n \rightarrow \infty} \frac{1}{n} \sum_{x_{1n}} p(x_{1n}) \log m(x_{1n})$$

- If the language is ‘nice’:

$$H(L, m) = - \lim_{n \rightarrow \infty} \frac{1}{n} \log m(x_{1n}) \quad (10)$$

I.e., it’s just our average surprise for large n :

$$H(L, m) \approx - \frac{1}{n} \log m(x_{1n}) \quad (11)$$

- Since $H(L)$ is fixed if unknown, minimizing cross-entropy is equivalent to minimizing $D(p \parallel m)$
- Providing: independent test data; assume $L = (X_i)$ is stationary [doesn’t change over time], ergodic [doesn’t get stuck]

Entropy of English text

27 letter alphabet

Model	Cross entropy (bits)
zeroth order	4.76 (log 27)
first order	4.03
second order	2.8
Shannon's experiment	1.3 (1.34)

Perplexity

$$\begin{aligned}\text{perplexity}(x_{1n}, m) &= 2^{H(x_{1n}, m)} \\ &= m(x_{1n})^{-\frac{1}{n}}\end{aligned}$$