

Entropy and Perplexity

Saturday, February 1, 2020 11:24

$$\text{info}_p(x_i) = -\log_2 p(x_i) \rightarrow \text{rare events} = \text{higher information} = \text{higher surprise}$$

entropy H = weighted information

$$= \sum p(x_i) \cdot -\log_2 p(x_i)$$

$$= -\sum p(x_i) \cdot \log_2 p(x_i)$$

entropy rate H_r = per-unit entropy

$$= H(\underbrace{x_1 \dots x_n}_{\text{a msg}}) / \underbrace{n}_{\text{msg length}}$$

entropy rate of n natural language

$$= \lim_{n \rightarrow \infty} \frac{H(\underbrace{x_1 \dots x_n}_{\text{observed samples}})}{n}$$

relative entropy
(or KL distance, KL divergence)

$D(p \parallel q) = \underbrace{\text{weighted average of differences in "Surprise"}}$

$$= \sum_x [q(x) \cdot (\text{info}_q(x) - \text{info}_p(x))]$$

less surprising

note that we can't build a model better than reality

thus, $\text{info}_p(x) \leq \text{info}_q(x) \rightarrow$ to keep rel. ent positive, we subtract $\text{info}_p(x)$ from $\text{info}_q(x)$

$$= \sum p(x) (-\log_2 q(x) + \log_2 p(x))$$

$$= \sum_x p(x) \cdot \log_2 \frac{p(x)}{q(x)} = \mathbb{E}_{x \in X} \log_2 \frac{p(x)}{q(x)}$$

cross entropy $H(X, q)$

\downarrow model pdf \downarrow true pdf

$$\begin{aligned} H(X, q) &= H(X) + D(p \parallel q) \\ &= -\sum p(x) \log_2 p(x) + \sum p(x) \log_2 \frac{p(x)}{q(x)} \\ &= -\cancel{\sum p(x) \log_2 p(x)} + \cancel{\sum p(x) \cdot \log_2 p(x)} - \sum p(x) \cdot \log_2 q(x) \\ &= -\sum_x p(x) \log_2 q(x) \end{aligned}$$

cross entropy of language L and model m

$$H(L, m) = \lim_{n \rightarrow \infty} \frac{1}{n} \cdot H(L, m)$$

$$= \lim_{n \rightarrow \infty} \frac{1}{n} \sum_{x_1 \dots x_n}^L L(x_1 \dots x_n) \cdot -\log_2 m(x_1 \dots x_n)$$

ergodicity (or the principle of equal a priori probability)

because we've seen enough samples, we take the observation as-is, not averaging

$$= \lim_{n \rightarrow \infty} \frac{1}{n} \log_2 \frac{1}{m(x_1 \dots x_n)}$$

n is sufficiently large, so we just approximate

$$\approx \frac{1}{n} \log_2 \frac{1}{m(x_1 \dots x_n)}$$

perplexity of model m against L

perplexity of model m against L

$$\begin{aligned} \text{pp}(L, m) &= 2^{-H(L, m)} \\ &= 2^{-\frac{1}{n} \log_2 \frac{1}{m(x_1 \dots x_n)}} \\ &= \left(2^{\log_2 \frac{1}{m(x_1 \dots x_n)}} \right)^{\frac{1}{n}} \\ &= \left(\frac{1}{m(x_1 \dots x_n)} \right)^{\frac{1}{n}} \\ &= \sqrt[n]{\frac{1}{m(x_1 \dots x_n)}} \end{aligned}$$

probability of sequence $x_1 \dots x_n$
as per the model

q: How to compute this probability
under bigram models?