# Information theory

## LING 572
## Fei Xia
## Week 1: 1/05/2010

# Information theory

- Reading: M&S 2.2

- It is the use of probability theory to quantify and measure "information".

- Basic concepts:
  - Entropy
  - Cross entropy and relative entropy
  - Joint entropy and conditional entropy
  - Entropy of the language and perplexity
  - Mutual information

# Entropy

- Entropy is a measure of the uncertainty associated with a distribution.

$$H(X) = -\sum_{x} p(x) \log p(x)$$

- The lower bound on the number of bits that it takes to transmit messages.

- An example:
  - Display the results of horse races.
  - Goal: minimize the number of bits to encode the results.

# An example

- Uniform distribution: $p_i = 1/8$.

$$H(X) = -8 * (\frac{1}{8} \log_2 \frac{1}{8}) = 3 \; bits$$

- Non-uniform distribution: (1/2,1/4,1/8, 1/16, 1/64, 1/64, 1/64, 1/64)

$$H(X) = -(\frac{1}{2} \log \frac{1}{2} + \frac{1}{4} \log \frac{1}{4} + \frac{1}{8} \log \frac{1}{8} + \frac{1}{16} \log \frac{1}{16} + 4 * \frac{1}{64} \log \frac{1}{64}) = 2 \; bits$$

(0, 10, 110, 1110, 111100, 111101, 111110, 111111)

➔Uniform distribution has higher entropy.
➔MaxEnt: make the distribution as "uniform" as possible.

# Cross Entropy

- Entropy:

$$H(X) = -\sum_x p(x) \log p(x)$$

- Cross Entropy:

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

- Cross entropy is a distance measure between p(x) and q(x): p(x) is the true probability; q(x) is our estimate of p(x).

$$H_c(X) \geq H(X)$$

# Relative Entropy

- Also called Kullback-Leibler divergence:

$$KL(p \| q) = \sum p(x) \log_2 \frac{p(x)}{q(x)} = H_c(X) - H(X)$$

- Another "distance" measure between probability functions p and q.

- KL divergence is asymmetric (not a true distance):

$$KL(p, q) \neq KL(q, p)$$

# Joint and conditional entropy

- Joint entropy:

$$H(X,Y) = -\sum_x \sum_y p(x,y) \log p(x,y)$$

- Conditional entropy:

$$H(Y \mid X) = \sum_x p(x) H(Y \mid X = x)$$

$$= -\sum_x p(x) \sum_y p(y \mid x) \log p(y \mid x)$$

$$= -\sum_x \sum_y p(x,y) \log p(y \mid x)$$

$$= H(X,Y) - H(X)$$

# Entropy of a language (per-word entropy)

- The entropy of a language L:

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

- If we make certain assumptions that the language is "nice", then the cross entropy can be calculated as:

$$H(L) = -\lim_{n \to \infty} \frac{\log p(x_{1n})}{n} \approx -\frac{\log p(x_{1n})}{n}$$

# Per-word entropy (cont)

- $p(x_{1n})$ can be calculated by n-gram models

- Ex: unigram model

$$p(x_{1n}) = \prod_i p(x_i)$$

$$log\ p(x_{1n}) = \sum_i log\ p(x_i)$$

# Perplexity

- Perplexity is $2^H$.

- Perplexity is the weighted average number of choices a random variable has to make.

# Mutual information

- It measures how much is in common between X and Y:

$$I(X;Y) = \sum_x \sum_y p(x,y) \log \frac{p(x,y)}{p(x)p(y)}$$

$$= H(X) + H(Y) - H(X,Y)$$

$$= I(Y;X)$$

- I(X;Y)=KL(p(x,y)||p(x)p(y))

- I(X;Y) = I(Y;X)

# Summary on Information theory

- Reading: M&S 2.2

- It is the use of probability theory to quantify and measure "information".

- Basic concepts:
  - Entropy
  - Cross entropy and relative entropy
  - Joint entropy and conditional entropy
  - Entropy of the language and perplexity
  - Mutual information