

# Mathematical Foundations

Foundations of Statistical Natural Language Processing, chapter2

Presented by Yi-Ting (怡婷)  
CSIE, NTNU

# Outline

- Elementary Probability Theory
  - Probability spaces
  - Conditional probability and independence
  - Bayes' theorem
  - Random variables
  - Expectation and variance
  - Joint and conditional distributions
  - Standard distributions
  - Bayesian statistics
- Essential Information Theory
  - Entropy
  - Joint entropy and conditional entropy
  - Mutual information
  - Relative entropy or Kullback-Leibler divergence

# Elementary Probability Theory

## Probability spaces

- Sample space:  $\Omega$
- Event  $A$  is the subset of  $\Omega$

- Probability function

$$P(A) = \frac{|A|}{|\Omega|}$$

- $P(\Omega)=1$

- Example :

A fair coin tossed 3 times. What is the chance of 2 heads?

- $\Omega = \{HHH, HHT, HTH, HTT, THH, THT, TTH, TTT\}$

- $A = \{HHT, HTH, THH\}$

- So

$$P(A) = \frac{|A|}{|\Omega|} = \frac{3}{8}$$

# Elementary Probability Theory

## Conditional probability and independence

- The **conditional probability** of an event A given that an event B has occurred is

$$P(A | B) = \frac{P(A \cap B)}{P(B)}$$

- Even if  $P(B)=0$  we have that :

$$P(A \cap B) = P(B)P(A|B) = P(A)P(B|A)$$

- The **chain rule** is as follows:

$$P(A_1 \cap \dots \cap A_n) =$$

$$P(A_1)P(A_2 | A_1)P(A_3 | A_1 \cap A_2) \dots P(A_n | \bigcap_{i=1}^{n-1} A_i)$$

# Elementary Probability Theory

## Conditional probability and independence

- Two event A, B are independent of each other if
$$P(A \cap B) = P(A)P(B)$$
- Two event A and B are conditionally independent given C when

$$P(A \cap B | C) = P(A|C)P(B|C)$$

# Elementary Probability Theory

## Bayes' theorem

- Bayes' theorem lets us swap the order of dependence between events.

$$P(B | A) = \frac{P(B \cap A)}{P(A)} = \frac{P(A | B)P(B)}{P(A)}$$

$$\arg \max_B P(B | A) = \arg \max_B \frac{P(B \cap A)}{P(A)} =$$

$$\arg \max_B \frac{P(A | B)P(B)}{P(A)} = \arg \max_B P(A | B)P(B)$$

# Elementary Probability Theory

## Bayes' theorem

- The set  $A$  can be divided into two parts

$$P(A \cap B) = P(A|B)P(B), \quad P(A \cap \bar{B}) = P(A|\bar{B})P(\bar{B})$$

*so we have:*

$$P(A) = P(A \cap B) + P(A \cap \bar{B}) = P(A|B)P(B) + P(A|\bar{B})P(\bar{B})$$

- If we have some group of sets  $B_i$  that partition  $A$ , if  $A \subseteq \bigcup_i B_i$  and the  $B_i$  are disjoint, then

$$P(A) = \sum_i P(A|B_i)P(B_i)$$

# Elementary Probability Theory

## Bayes' theorem

- Bayes' theorem

*if  $A \subseteq \cup_{i=1}^n B_i$ ,  $P(A) > 0$ , and  $B_i \cap B_j = \phi$ , for  $i \neq j$  then:*

$$P(B_j | A) = \frac{P(A | B_j)P(B_j)}{P(A)} = \frac{P(A | B_j)P(B_j)}{\sum_{i=1}^n P(A | B_i)P(B_i)}$$



# Elementary Probability Theory

## Random variables

- Random variables is simply a function

$$X: \Omega \rightarrow \mathbf{R}^n$$

$\mathbf{R}$  is the set of real numbers, commonly with  $n=1$

- A discrete random variable is a function

$$X: \Omega \rightarrow \mathbf{S}$$

where  $\mathbf{S}$  is a countable subset of  $\mathbf{R}$

- A indicator random variable is a function

$$X: \Omega \rightarrow \{0,1\},$$

and  $X$  is also called a Bernoulli trial

# Elementary Probability Theory

## Random variables

- We can define the **probability mass function (pmf)** for a random variable  $X$ , which gives the random variable has different numeric values:

$$\text{pmf } p(x) = p(X = x) = P(A_x)$$

$$\text{where } A_x = \{\omega \in \Omega : X(\omega) = x\}$$

- For a discrete random variable , we have

$$\sum_i p(x_i) = \sum_i P(A_{x_i}) = P(\Omega) = 1$$

# Elementary Probability Theory

## Expectation and variance

- The **expectation** is the mean or average of a random variable
- If  $X$  is a random variable with a pmf  $p(x)$  such that

$$\sum_x |x|p(x) < \infty$$

- Then the expectation is

$$E(X) = \sum_x xp(x)$$

- Example: if  $Y$  is the value of face on one rolling die ,then

$$E(Y) = \sum_{y=1}^6 yp(y) = \frac{1}{6} \sum_{y=1}^6 y = \frac{21}{6} = 3\frac{1}{2}$$

- This is the expected average found by totaling up a large number of throws of the die, and dividing by the number of throws.

# Elementary Probability Theory

## Expectation and variance

- If  $Y \sim p(y)$  is a random variable, any function  $g(Y)$  defines a new random variable.
- If  $E(g(Y))$  is defined, then

$$E(g(Y)) = \sum_y g(y)p(y)$$

– Example :  $g(Y)=aY+b$ , we see that  $E(g(Y))=aE(Y)+b$

- We also have that  $E(X+Y)=E(X)+E(Y)$
- If  $X$  and  $Y$  are independent, then  $E(XY)=E(X)E(Y)$

# Elementary Probability Theory

## Expectation and variance

- The **variance** is the measure of the random variable tend to be consistent over trials or to vary a lot.
- One measures it by finding out how much on average the variable's values deviate from the variable's expectation

$$\text{Var}(X) = E\left(\left(X - E(X)\right)^2\right) = E(X^2) - E^2(X)$$

- The **standard deviation** of a variable is the square root of the variance.
- In commonly denotes the **mean** is  $\mu$  and the **variance** is  $\sigma^2$  the **standard deviation** is hence written as  $\sigma$

# Elementary Probability Theory

## Expectation and variance

- Proof of variance calculation I

$$\begin{aligned} \text{Var}(X) &= E\left(\left(X - E(X)\right)^2\right) \\ &= E\left(X^2 - 2XE(X) + \left(E(X)\right)^2\right) \\ &= E\left(X^2\right) - E\left(2XE(X)\right) + E\left(\left(E(X)\right)^2\right) \\ &= E\left(X^2\right) - 2E(X)E(X) + \left(E(X)\right)^2 \\ &= E\left(X^2\right) - E^2(X) \end{aligned}$$

# Elementary Probability Theory

## Expectation and variance

- Proof of variance calculation II

$$\begin{aligned} \text{Var}(X) &= E\left(\left(X - E(X)\right)^2\right) = E(X^2) - E^2(X) \\ &= \sum_x p(x^2)x^2 - 2E^2(X) + E^2(X) \\ &= \sum_x p(x)x^2 - 2E(X)\sum_x p(x)x + 1E^2(X) \\ &= \sum_x p(x)x^2 - \sum_x p(x)x2E(X) + \sum_x p(x)E^2(X) \\ &= \sum_x p(x)\left(x^2 - 2xE(X) + E^2(X)\right) \\ &= \sum_x p(x)(x - E(X))^2 \end{aligned}$$

# Elementary Probability Theory

## Joint and conditional distributions

- The joint probability mass function for two discrete random variables  $X, Y$  is

- 

$$p(x, y) = P(X=x, Y=y)$$

- The marginal pmfs, which total up the probability masses for the value of each variable separately

$$p_X(x) = \sum_y p(x, y) \quad , \quad p_Y(y) = \sum_x p(x, y)$$



# Elementary Probability Theory

## Joint and conditional distributions

- If  $X$  and  $Y$  are independent, then

$$p(x,y)=p_X(x)p_Y(y)$$

- Example:

getting two sixes from rolling two dice, since the events are independent, we can compute that:

$$p(Y = 6, Z = 6) = p(Y = 6)p(Z = 6) = \frac{1}{6} \times \frac{1}{6} = \frac{1}{36}$$

- The conditional pmf in terms of the joint distribution

$$p_{X|Y}(x | y) = \frac{p(x, y)}{p_Y(y)} \quad \text{for } y \text{ such that } p_Y(y) > 0$$

- And deduce a chain rule in terms of random variables, like

$$p(w, x, y, z) = p(w)p(x | w)p(y | w, x)p(z | w, x, y) \quad 17$$

# Elementary Probability Theory

## Standard distributions

- Discrete distributions:
  - Binomial distribution
- Continuous distributions:
  - Normal distribution

# Elementary Probability Theory

## Standard distributions

- The **Binomial distribution** results when one has a series of trials with only two outcomes, each trial being independent from all the others.
- The **binomial distributions** gives the number  $r$  of successes out of  $n$  trials and the **probability of success** in any trial is  $p$

$$b(r; n, p) = \binom{n}{r} p^r (1-p)^{n-r} \text{ where } \binom{n}{r} = \frac{n!}{(n-r)!r!} \quad 0 \leq r \leq n$$

- Let  $R$  have as value the number of heads in  $n$  tosses of a coin, where the probability of a head is  $p$

$$p(R = r) = b(r; n, p)$$

# Elementary Probability Theory

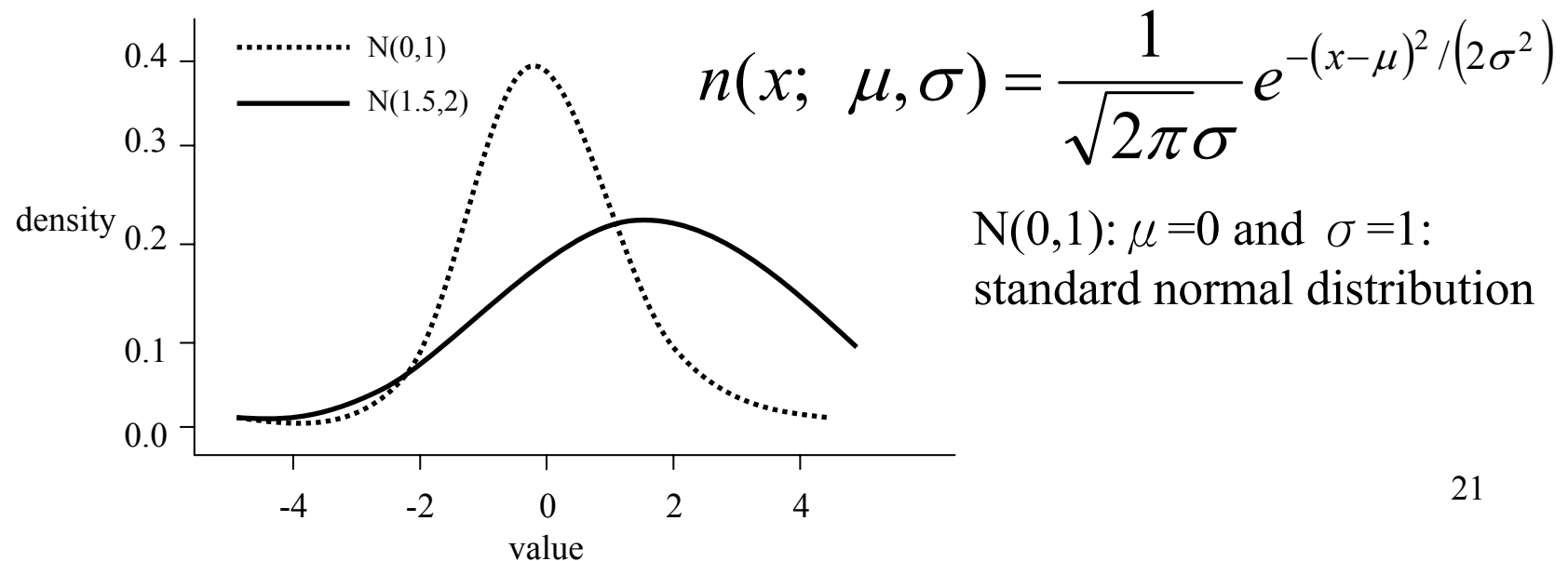
## Standard distributions

- Multinomial distribution
  - The generalization of a binomial trial to the case where each of the trial has **more than two basic outcomes** is called multinomial experiment and modeled by it.
  - A zeroth order n-gram model is a straightforward example of a multinomial distribution.

# Elementary Probability Theory

## Standard distributions

- Normal distribution
  - With two parameters :  $\mu$  : mean (variance)  
 $\sigma$  : standard deviation
  - And the bell curve is given by:



# Elementary Probability Theory

## Bayesian statistics

- Bayesian updating
  - A coin is tossed in times and gets 8 heads then this coin comes down heads 8 times out of 10.
    - This is the **maximum likelihood estimate**
  - But he belief the coin would come down equally head and tails over the long run this is called a **prior belief**
  - Bayesian statistics
    - Measure degree of belief
    - Starting with prior belief
    - updating tem in the face of evidence
    - By use of Bayes' theorem

# Elementary Probability Theory

## Bayesian statistics

- $\mu_m$  be the model that asserts  $P(\text{head}) = m$
- $s$  be a sequence of observations:  
 $i$  heads and  $j$  tails
- For any  $m$ ,  $0 \leq m \leq 1$   
$$P(s | \mu_m) = m^i (1-m)^j$$
- From a frequentist point of view, we wish to find the MLE  
$$\arg \max_m P(s | \mu_m)$$
- We can differentiate the above polynomial then the answer is  $i / (i+j)$ , or 0.8 for the case of 8 heads and 2 tails

# Elementary Probability Theory

## Bayesian statistics

- Assume one's prior belief is modeled by

$$P(\mu_m) = 6m(1-m)$$

because this distribution is centered on 1/2

- By bayes' theorem

$$\begin{aligned} P(\mu_m | s) &= \frac{P(s | \mu_m)P(\mu_m)}{P(s)} \\ &= \frac{m^i (1-m)^j \times 6m(1-m)}{P(s)} \\ &= \frac{6m^{i+1} (1-m)^{j+1}}{P(s)} \end{aligned}$$



# Elementary Probability Theory

## Bayesian statistics

- $P(s)$  is the prior probability of  $s$
- $s$  doesn't depend on  $\mu_m$  so we can ignore it
- Then we can determine the case for 8 heads and 2 tails

$$\begin{aligned}\arg \max_m P(\mu_m | s) &= \frac{6m^{i+1}(1-m)^{j+1}}{P(s)} \\ &= \arg \max_m 6m^{i+1}(1-m)^{j+1} = \arg \max_m 6m^{8+1}(1-m)^{2+1} \\ &= \arg \max_m 6m^9(1-m)^3 = \frac{9}{9+3} = \frac{3}{4}\end{aligned}$$

- We have moved a long way in the direction of believing that the coin is biased, but we haven't moved all the way to 0.8

# Elementary Probability Theory

## Bayesian statistics

- Marginal probability
  - Adding up all the  $P(s | \mu_m)$  weighted by the probability of  $\mu_m$
- For the continuous case

$$\begin{aligned}P(s) &= \int_0^1 P(s | \mu_m) P(\mu_m) dm \\ &= \int_0^1 6m^{i+1} (1-m)^{j+1} dm \\ &= \frac{6(i+1)!(j+1)!}{(i+j+3)!}\end{aligned}$$

- It is a normalization factor, for  $P(\mu_m | s)$  is actually a probability function

# Elementary Probability Theory

## Bayesian statistics

- Bayesian decision theory
  - To evaluate which model better explains some data
- Example:  
comparing two models  $\nu$  and  $\mu$ 
  - Tossing two fair coins and called out “tails” if both of them come down tails this is called theory  $\nu$  and the theory  $\mu$  above

$$\text{we have } P(s | \nu) = \left(\frac{3}{4}\right)^i \left(\frac{1}{4}\right)^j \text{ and } P(\mu) = P(\nu) = \frac{1}{2}$$

$$P(\mu | s) = \frac{P(s | \mu)P(\mu)}{P(s)}, P(\nu | s) = \frac{P(s | \nu)P(\nu)}{P(s)}$$

# Elementary Probability Theory

## Bayesian statistics

- Bayesian decision theory

$$\begin{aligned} \frac{P(\mu | s)}{P(\nu | s)} &= \frac{P(s | \mu)P(\mu)}{P(s)} \times \frac{P(s)}{P(s | \nu)P(\nu)} \\ &= \frac{P(s | \mu)P(\mu)}{P(s | \nu)P(\nu)} = \frac{\frac{6(i+1)!(j+1)!}{(i+j+3)!}}{\left(\frac{3}{4}\right)^i \left(\frac{1}{4}\right)^j} = \frac{\frac{6(8+1)!(2+1)!}{(8+2+3)!}}{\left(\frac{3}{4}\right)^8 \left(\frac{1}{4}\right)^2} = 0.33 \end{aligned}$$

- The quantity we are now describing as  $P(s | \mu)$  is the quantity that we wrote as just  $P(s)$
- If the ratio is greater than 1, we should prefer  $\mu$

# Outline

- Essential Information Theory
  - Entropy
  - Joint entropy and conditional entropy
  - Mutual information
  - Relative entropy or Kullback-Leibler divergence

# Essential Information Theory

## Entropy

- Entropy measures the amount of information in a random variable. It is normally measured in bits.

$$H(X) = - \sum_{x \in X} p(x) \log_2 p(x)$$

- We define

$$0 \log_2 0 = 0$$

# Essential Information Theory

## Entropy

- Example:

Suppose you are reporting the result of rolling an 8-sided die. Then the entropy is:

$$\begin{aligned} H(X) &= -\sum_{i=1}^8 p(i) \log p(i) = -\sum_{i=1}^8 \frac{1}{8} \log \frac{1}{8} \\ &= -\log \frac{1}{8} = \log 8 = 3 \text{ bits} \end{aligned}$$

# Essential Information Theory

## Entropy

- Entropy:
  - The average number of bits used for identifying the transmission of the information
  - We hope the entropy is lower in the system



# Essential Information Theory

## Entropy

- Properties of Entropy:

$$H(X) = - \sum_{x \in X} p(x) \log_2 p(x)$$

$$= \sum_{x \in X} p(x) \log_2 \frac{1}{p(x)}$$

$$= E \left( \log \frac{1}{p(x)} \right)$$

# Essential Information Theory

## Joint Entropy and Conditional Entropy

- Joint Entropy:

$$H(X, Y) = - \sum_{x \in X} \sum_{y \in Y} p(x, y) \log p(x, y)$$

- Conditional Entropy:

$$H(Y | X) = - \sum_{x \in X} \sum_{y \in Y} p(y, x) \log p(y | x)$$

# Essential Information Theory

## Joint Entropy and Conditional Entropy

- Proof of Conditional Entropy:

$$\begin{aligned} H(Y | X) &= \sum_{x \in X} p(x) H(Y | X = x) \\ &= \sum_{x \in X} p(x) \left[ - \sum_{y \in Y} p(y | x) \log p(y | x) \right] \\ &= - \sum_{x \in X} \sum_{y \in Y} p(y, x) \log p(y | x) \end{aligned}$$

# Essential Information Theory

## Joint Entropy and Conditional Entropy

- Chain rule for Entropy:

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

- Proof:

$$\begin{aligned} H(X, Y) &= -\sum_{x \in X} \sum_{y \in Y} p(x, y) \log p(x, y) \\ &= -\sum_{x \in X} \sum_{y \in Y} p(x, y) \log(p(y | x) p(x)) \\ &= -\sum_{x \in X} \sum_{y \in Y} p(x, y) (\log p(y | x) + \log p(x)) \\ &= -\sum_{x \in X} \sum_{y \in Y} p(x, y) \log p(y | x) - \sum_{x \in X} \sum_{y \in Y} p(x, y) \log p(x) \\ &= H(Y | X) + H(X) \end{aligned}$$

# Essential Information Theory

## Entropy rate

- Per-letter or per-word entropy
- For a message of length  $n$  the entropy rate

$$H_{rate} = \frac{1}{n} H(X_{1n}) = -\frac{1}{n} \sum_{X_{1n}} p(X_{1n}) \log p(X_{1n})$$

- Assume that a language is a stochastic process consisting of a sequence of tokens  $L=(X_i)$

$$H_{rate}(L) = \lim_{n \rightarrow \infty} \frac{1}{n} H(X_1, X_2, \dots, X_n)$$

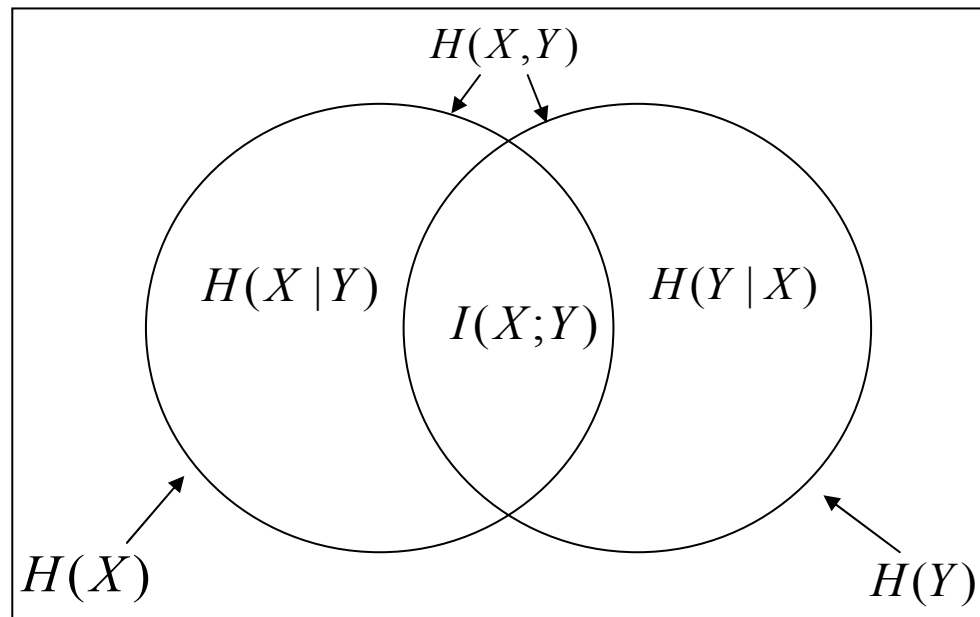
# Essential Information Theory

## Mutual Information

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

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

This difference is called the mutual information between X and Y



# Essential Information Theory

## Mutual Information

- The likeness of Information ◦
- This difference is called the *mutual information* between X and Y.
- The amount of information one random variable contains about another.
- It is 0 only when two variables are independent.  
The mutual Information is 0 for two independent events ◦

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

# Essential Information Theory

## Mutual Information

- How to simply calculate Mutual Information ?

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

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

$$= \sum_x p(x) \log \frac{1}{p(x)} + \sum_y p(y) \log \frac{1}{p(y)} + \sum_{x,y} p(x,y) \log p(x,y)$$

$$= \sum_{x,y} p(x,y) \log \frac{1}{p(x)} + \sum_{x,y} p(x,y) \log \frac{1}{p(y)} + \sum_{x,y} p(x,y) \log p(x,y)$$

$$= \sum_{x,y} p(x,y) \left[ \log \frac{1}{p(x)} + \log \frac{1}{p(y)} + \log p(x,y) \right]$$

$$= \sum_{x,y} p(x,y) \log \frac{p(x,y)}{p(x)p(y)}$$



# Essential Information Theory

## Mutual Information

- Conditional mutual information

$$I(X; Y | Z) = I((X; Y) | Z) = H(X | Z) - H(X | Y, Z)$$

- Chain rule

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

# Essential Information Theory

## Mutual Information

- Define the *pointwise mutual information* between two particular points.

$$I(x, y) = \log \frac{p(x, y)}{p(x)p(y)}$$

This has sometimes been used as a measure of association between elements.

# Essential Information Theory

## Relative Entropy or Kullback-Leibler divergence

- For two probability mass functions,  $p(x)$  ,  $q(x)$  their relative entropy is given by:

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

*define*  $0 \log \frac{0}{q} = 0$  and  $p \log \frac{p}{0} = \infty$

# Essential Information Theory

## Relative Entropy or Kullback-Leibler divergence

- Meaning : It is 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  $q$ .
- Some authors use the name “**KL distance**”, but note that relative entropy isn't a metric (it doesn't satisfy the triangle inequality)

# Essential Information Theory

## Relative Entropy or Kullback-Leibler divergence

Properties of KL-divergence:

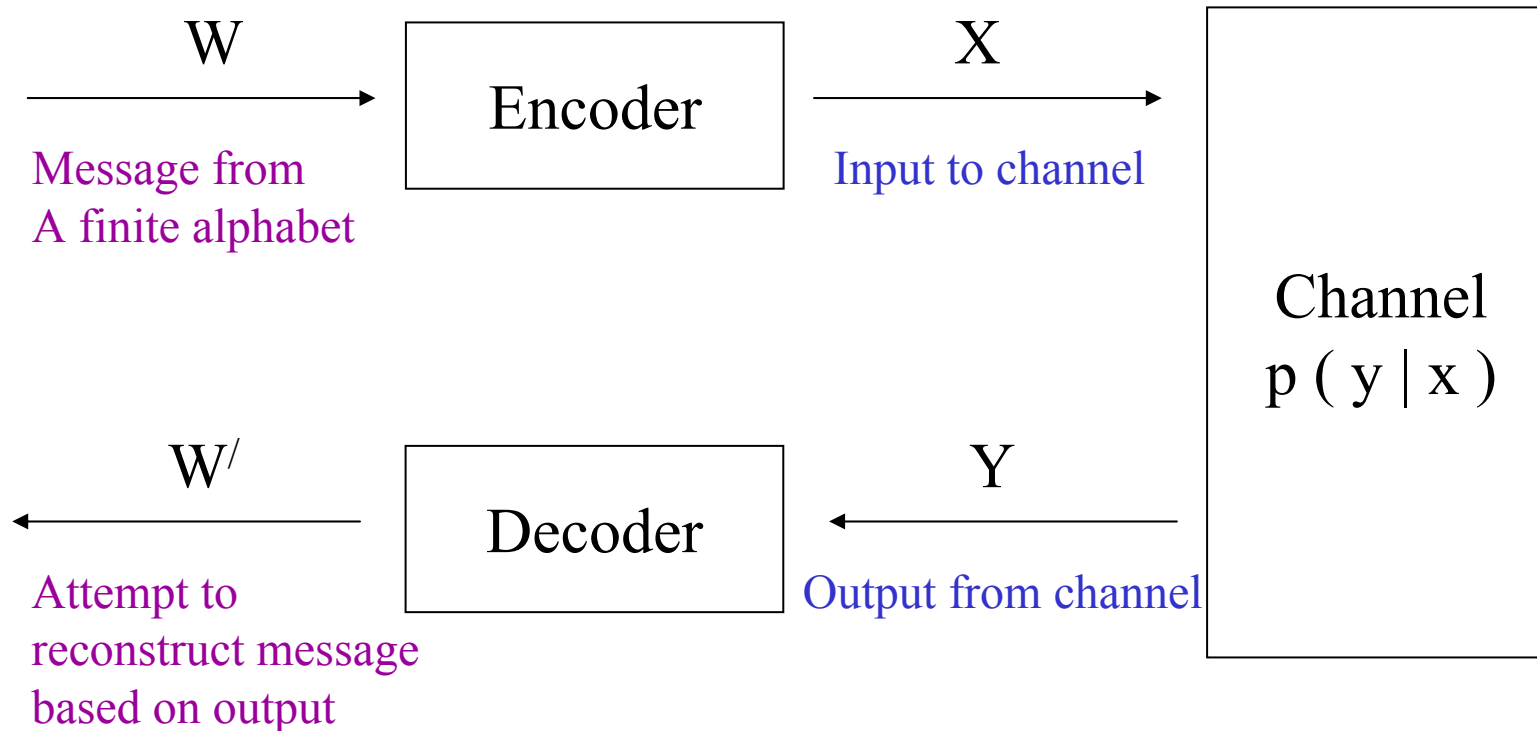
$$I(X;Y) = \sum_{x,y} p(x,y) \log \frac{p(x,y)}{p(x)p(y)}$$
$$= D(p(x,y) \parallel p(x)p(y))$$

Define the Conditional Relative Entropy:

$$D(p(y|x) \parallel q(y|x)) = \sum_x p(x) \sum_y p(y|x) \log \frac{p(y|x)}{q(y|x)}$$

# Essential Information Theory

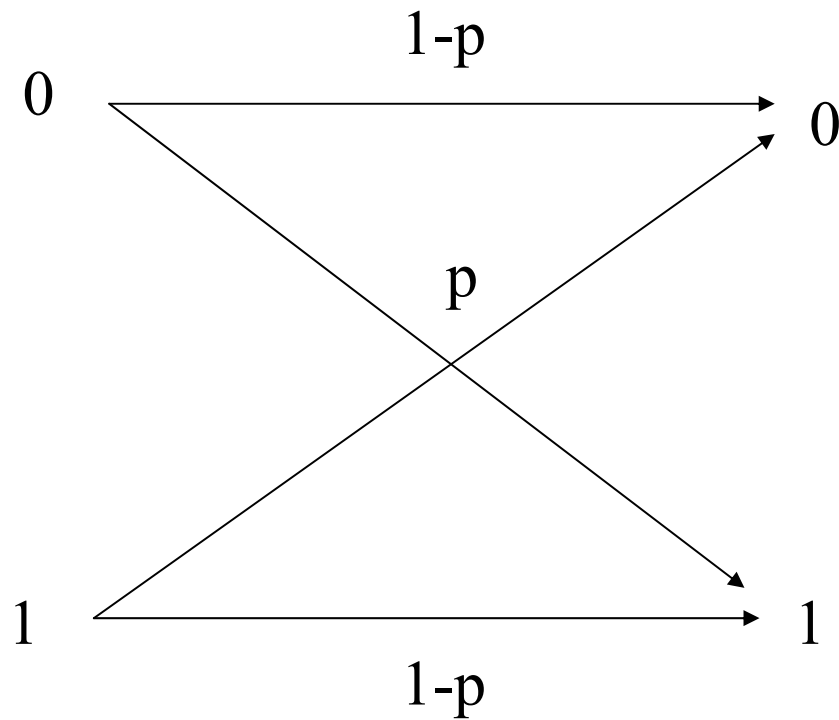
## The noisy channel model



The noisy channel model

# Essential Information Theory

## The noisy channel model



A binary symmetric channel

# Essential Information Theory

## The noisy channel model

- Capacity

- The channel capacity describes the rate at which one can transmit information through the channel with an arbitrarily **low probability of being unable to recover the input from the output.**

$$C = \max_{p(X)} I(X; Y) \quad \text{if } p = 0 \text{ or } p = 1 \Rightarrow C = 1$$

$$= \max_{p(X)} H(Y) - H(Y | X) \quad \text{if } p = \frac{1}{2} \Rightarrow C = 0$$

$$= \max_{p(X)} H(Y) - H(p) = 1 - H(p) \quad 0 < C \leq 1$$

The capacity is used to measure the likeness of X and Y  
If the mutual information is 1 then the X and Y are the same or  
bits are inverted completely



# Essential Information Theory

## The noisy channel model

$$C = \max_{p(X)} I(X; Y)$$

$$= \max_{p(X)} H(Y) - H(Y | X)$$

$$= \max_{p(X)} H(Y) - H(p) = 1 - H(p) \quad 0 < C \leq 1$$

$$H(Y) = - \sum_{y \in Y} p(y) \log_2 p(y) = -\frac{1}{2} \log_2 \frac{1}{2} + -\frac{1}{2} \log_2 \frac{1}{2} = 1$$

# Essential Information Theory

## The noisy channel model

Application: (In speech recognition)

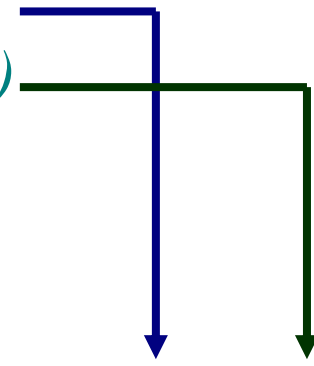
*Input:* word sequences

*Output:* observed speech signal

*P(input):* probability of word sequences

*P(output|input):* acoustic model (channel prob.)

Bayes' theorem

$$\hat{I} = \arg \max_i p(i | o) = \arg \max_i \frac{p(i)p(o | i)}{p(o)} = \arg \max_i \boxed{p(i)} \boxed{p(o | i)}$$


# Essential Information Theory

## Cross entropy

- If a model captures more of the structure of a language, then the entropy of the model should be lower
- Entropy is a measure of the quality of our models

p	t	k	a	i	u	p	t	k	a	i	u
1/8	1/4	1/8	1/4	1/8	1/8	100	00	101	01	110	111

$$\begin{aligned} H(P) &= - \sum_{i \in \{p,t,k,a,i,u\}} P(i) \log P(i) \\ &= - \left[ 4 \times \frac{1}{8} \log \frac{1}{8} + 2 \times \frac{1}{4} \log \frac{1}{4} \right] = 2 \frac{1}{2} \text{ bits} \end{aligned}$$

# Essential Information Theory

## Cross entropy

- Cross entropy:
  - The *cross entropy* between a random variable  $X$  with true probability distribution  $p(X)$  and another pmf  $q$  (normally a model of  $p$ ) is given by:

$$\begin{aligned} H(X, q) &= H(X) + D(p \parallel q) \\ &= \sum_{x \in X} p(x) \log \frac{1}{p(x)} + \sum_{x \in X} p(x) \log \frac{p(x)}{q(x)} \\ &= \sum_{x \in X} p(x) \left[ \log \frac{1}{p(x)} + \log \frac{p(x)}{q(x)} \right] \\ &= \sum_{x \in X} p(x) \left[ \log \frac{1}{q(x)} \right] = - \sum_{x \in X} p(x) \log q(x) \end{aligned}$$

# Essential Information Theory

## Cross entropy

Cross entropy of a language :

*suppose*

*Language  $L = (X_i) \sim p(x)$  according to a model  $m$  by*

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

*We cannot calculate this quantity **without knowing  $p$** . But if we make certain assumptions that the language is 'nice,' then the **cross entropy** for the language can be calculated as:*

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

# Essential Information Theory

## Cross entropy

- Expectation is a weighted average over all possible sequence
- If we have seen a huge amount of the language, what we have seen is “typical”
- We no longer need to average over all samples of the language
- The value for the entropy rate given by this particular sample will be roughly right

$$\begin{aligned} H(L, m) &= -\lim_{n \rightarrow \infty} \frac{1}{n} \sum_{x_{1n}} p(x_{1n}) \log m(x_{1n}) = -\lim_{n \rightarrow \infty} \frac{1}{n} E(\log m(x_{1n})) \\ &= \lim_{n \rightarrow \infty} \frac{1}{n} E\left(\log \frac{1}{m(x_{1n})}\right) \\ &= \lim_{n \rightarrow \infty} \frac{1}{n} \log \frac{1}{m(x_{1n})} \approx -\frac{1}{n} \log m(x_{1n}) \end{aligned}$$

# Essential Information Theory

## Cross entropy

- Cross entropy of a language :
  - We do not actually attempt to calculate the limit, but approximate it by calculating for a sufficiently large  $n$ :

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

- This measure is just the figure for our average surprise.
- *Our goal will be to try to minimize this number. Because  $H(X)$  is fixed, this is equivalent to minimizing the relative entropy, which is a measure of how much our probability distribution departs from actual language use.*

# Essential Information Theory

## Perplexity

*In the speech recognition community, people tend to refer to **perplexity** rather than **cross entropy**. The relationship between the two is simple:*

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

*Why we use perplexity not cross entropy?*

*Because it is much easier to impress funding bodies by saying that “we’ve managed to reduce perplexity from 950 to only 540” than by saying that “we’ve reduced cross entropy from 9.9 to 9.1 bits.”*