



香港城市大學  
City University of Hong Kong

專業 創新 胸懷全球  
Professional • Creative  
For The World

# Information Theory

CS5483 Data Warehousing and Data Mining

## The origin of ...

$$\text{Info}(D) := \sum_k p_k \log \frac{1}{p_k}$$

$$\text{Info}_X(D) := \sum_j \frac{|D_j|}{|D|} \text{Info}(D_j)$$

$$\text{Gain}_X(D) := \text{Info}(D) - \text{Info}_X(D)$$

$$\text{SplitInfo}_X(D) := \sum_j \frac{|D_j|}{|D|} \log \frac{1}{|D_j|/|D|}$$

- $p_k$ : fraction of tuples in  $D$  with class  $k$ .
- $X$ : splitting attribute.
- $D_j$ : Data satisfying  $X = j$ .

# What is information?

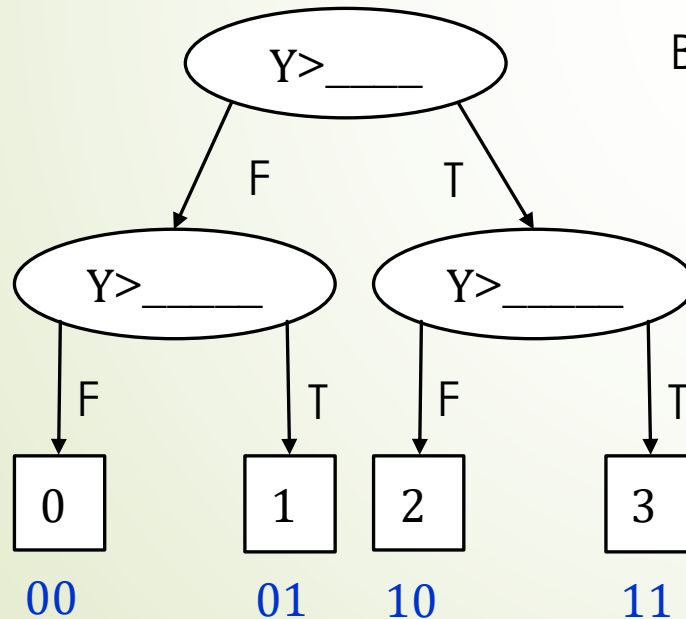
*“Information is the resolution of uncertainty”*



Claude Elwood Shannon  
(1916-2001)

# A game of resolution of uncertainty

- Determine the value of  $Y$  by asking as few T/F questions as possible.
- Given  $Y$  takes value from  $\{0,1\}$ , we can ask whether \_\_\_\_\_. So, \_\_\_\_\_ question.
- What if  $Y$  takes value from  $\{0,1,2,3\}$ ? \_\_\_\_\_ questions.



Basically a decision tree with splitting attribute \_\_\_\_.

**Encode**  $Y$  into a sequence of bits with:  
 $F \leftrightarrow 0$  and  $T \leftrightarrow 1$

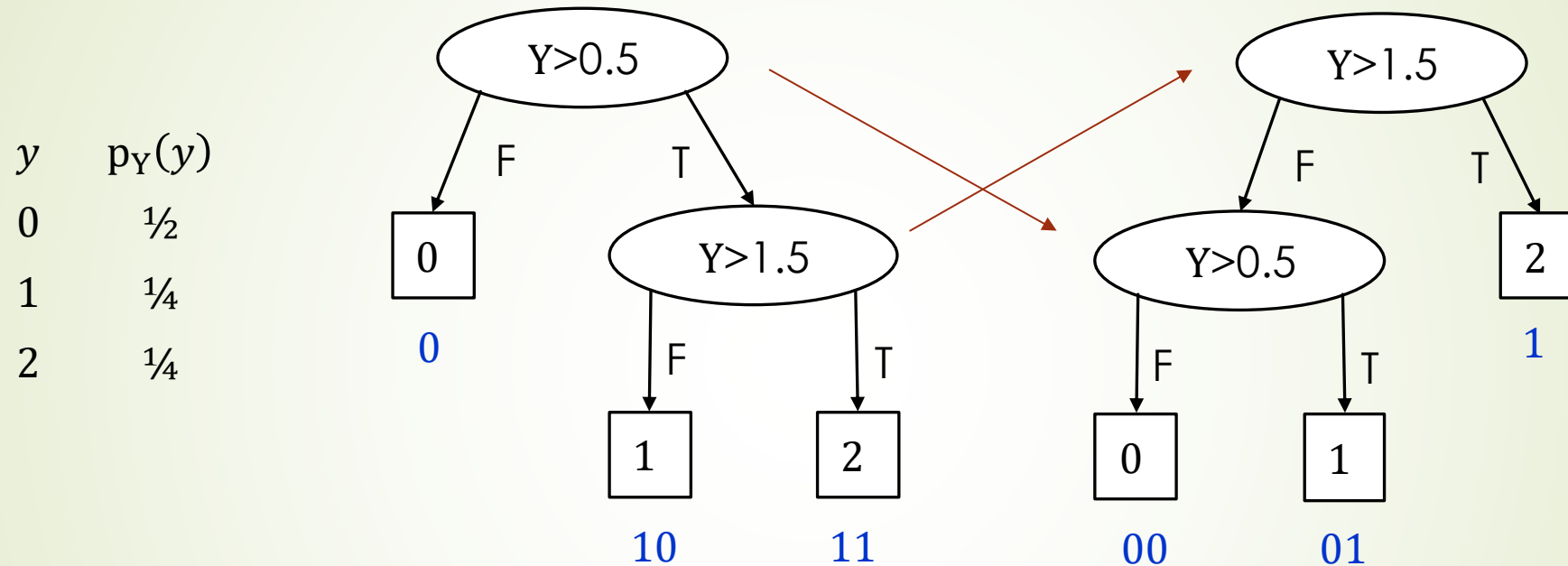
#Questions=#\_\_\_\_\_ required to encode  $Y$ ,  
 i.e., to resolve the *uncertainty* of  $Y$ .

codewords

# A game of resolution of uncertainty

- ▶ What if  $Y$  takes value from  $\{0, 1, \dots, 2^m - 1\}$ ? \_\_\_\_\_ questions.
- ▶ What if  $Y$  takes value from  $\{0, 1, 2\}$ ? \_\_\_\_\_ questions because  $\log_2 3 = 1.585$ .
- ▶ Can we ask only  $\log_2 3 = 1.585$  questions to resolve  $Y$ ?
- ▶ **C** \_\_\_\_\_  $n$  variables
  - ▶ Resolve  $Y_1, Y_2, \dots, Y_n$  together using \_\_\_\_\_ questions.
  - ▶ Information rate  $\frac{1}{n}$  \_\_\_\_\_  $\rightarrow \log_2 3$  questions per variable
- ▶ Can we do even better?

# Expected codeword length

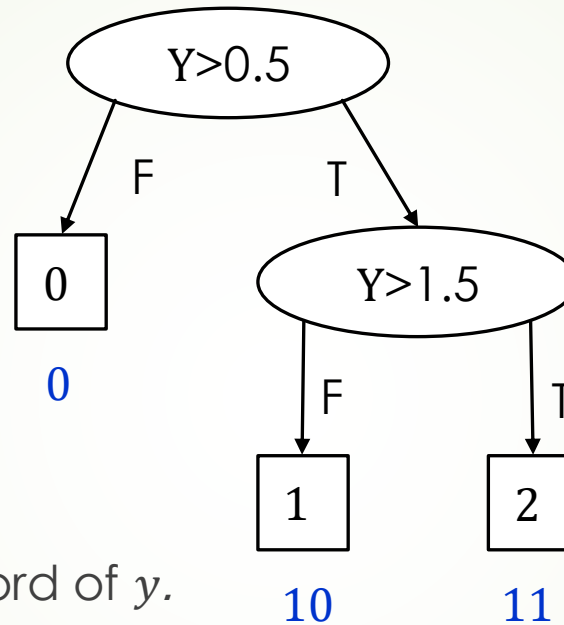


➤ Which encoding is better?

- Left/right because the expected codeword is larger/smaller.
- In general, assign shorter/longer codeword for more probable  $y$ .

# Entropy coding

$y$	$p_Y(y)$	$\log \frac{1}{p_Y(y)}$
0	$\frac{1}{2}$	—
1	$\frac{1}{4}$	—
2	$\frac{1}{4}$	—



- $l(y) :=$  length of the codeword of  $y$ .
- [Variable-length coding]  $Y$  can be encoded with

$$l(y) = \left\lceil \log \frac{1}{p_Y(y)} \right\rceil$$

- Expected codeword length is \_\_\_\_\_ .



# Entropy coding with concatenation

- Encode i.i.d. sequence  $Y^n := (Y_1, Y_2, \dots, Y_n)$ :

$$\begin{aligned}
 l(y_1, y_2, \dots, y_n) &= \left\lceil \log_2 \frac{1}{p_{Y_1, \dots, Y_n}(y_1, \dots, y_n)} \right\rceil \\
 &= \left\lceil \log_2 \frac{1}{p_Y(y_1)p_Y(y_2) \dots p_Y(y_n)} \right\rceil \\
 &= \left\lceil \sum_{i=1}^n \log_2 \frac{1}{p_Y(y_i)} \right\rceil \\
 \frac{1}{n} l(y_1, y_2, \dots, y_n) &= \frac{1}{n} \left\lceil \sum_{i=1}^n \log_2 \frac{1}{p_Y(y_i)} \right\rceil \\
 &= \frac{1}{n} \sum_{i=1}^n \log_2 \frac{1}{p_Y(y_i)} + o\left(\frac{1}{n}\right) \\
 &\rightarrow E \left[ \log_2 \frac{1}{p_Y(Y)} \right] = H(Y)
 \end{aligned}$$

- Length is minimal by Kraft's inequality.

$y_1$	$y_2$	$p_Y(y_1)$	$p_Y(y_2)$	$p_{Y_1 Y_2}(y_1, y_2)$	$\log_2 \frac{1}{p_{Y_1 Y_2}(y_1, y_2)}$
0	0	$2^{-1}$	$2^{-1}$	$2^{-2}$	2
	1		$2^{-2}$	$2^{-3}$	3
	2		$2^{-2}$	$2^{-3}$	3
1	0	$2^{-2}$	$2^{-1}$	$2^{-3}$	3
	1		$2^{-2}$	$2^{-4}$	4
	2		$2^{-2}$	$2^{-4}$	4
2	0	$2^{-2}$	$2^{-1}$	$2^{-3}$	3
	1		$2^{-2}$	$2^{-4}$	4
	2		$2^{-2}$	$2^{-4}$	4



# Conditional entropy and mutual information

- Join entropy:

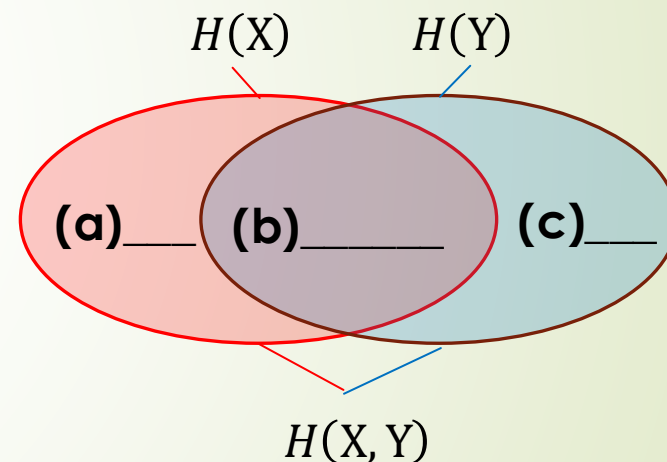
$$H(X, Y) := E \left[ \log \frac{1}{p_{XY}(X, Y)} \right]$$

- Conditional entropy:

$$\begin{aligned} H(Y|X) &:= H(X, Y) - H(X) \\ &= E \left[ \log \frac{1}{p_{XY}(X, Y)} - \log \frac{1}{p_X(X)} \right] \\ &= E \left[ \log \frac{p_X(X)}{p_{XY}(X, Y)} \right] \\ &= E \left[ \log \frac{1}{p_{Y|X}(Y|X)} \right] \end{aligned}$$

- Mutual Information:  $I(X; Y) := H(Y) - H(Y|X)$

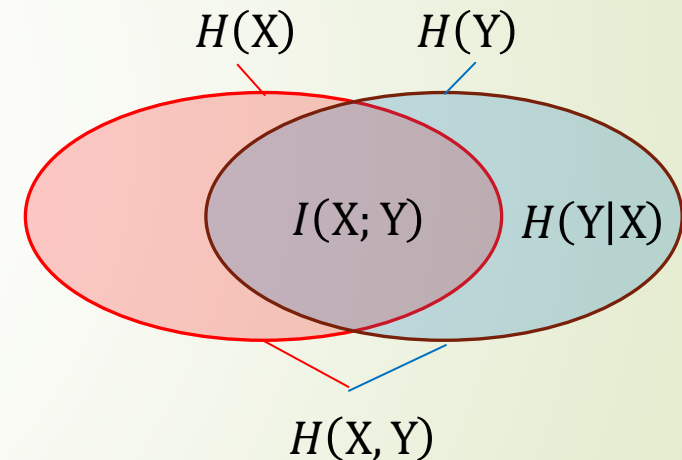
Venn diagram interpretation



# Conditional entropy

- $\text{Info}_X(D) := \sum_j \frac{|D_j|}{|D|} \text{Info}(D_j)$   
 $= \sum_j p_X(j) \sum_y p_{Y|X}(y|j) \log \frac{1}{p_{Y|X}(y|j)}$   
 $= E \left[ \log \frac{1}{p_{Y|X}(Y|X)} \right] = \underline{\hspace{2cm}}$
- $\text{Gain}_X(D) := \text{Info}(D) - \text{Info}_X(D)$   
 $= H(Y) - H(Y|X) = \underline{\hspace{2cm}}$
- $\text{Gain}_X(D) / \text{SplitInfo}_X(D) = \underline{\hspace{2cm}}$

Venn diagram interpretation



# Properties of information

## Chain rule:

$$I(\underbrace{X_1, X_2}_{\text{joint}}, \underbrace{Y}_{\text{target}}) = \underbrace{I(X_1; Y)}_{\text{mutual info}} + \underbrace{I(X_2; Y|X_1)}_{\text{conditional mutual info}} = \underbrace{I(X_2; Y)}_{\text{mutual info}} + \underbrace{I(X_1; Y|X_2)}_{\text{conditional mutual info}}$$

## Conditioning reduces entropy:

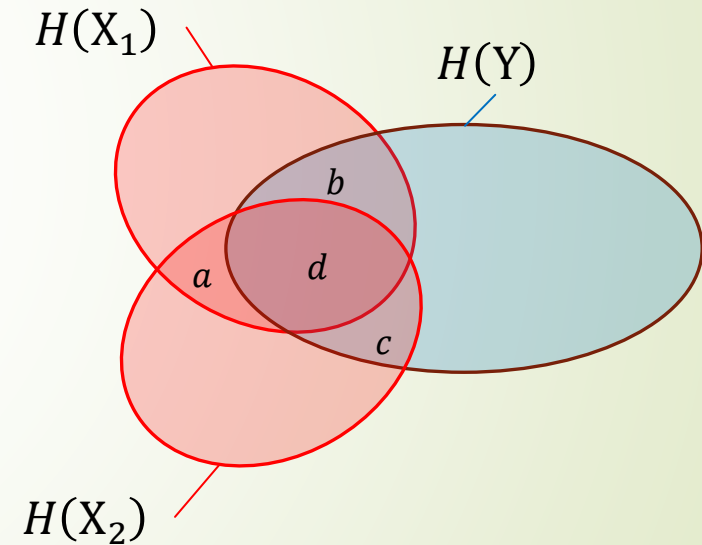
$$H(Y) \geq H(Y|X_1) \geq H(Y|X_1, X_2)$$

$$\text{Equiv.}, I(\text{_____}), I(\text{_____}) \geq 0$$

## Data processing inequality:

$$I(X_2; Y|X_1) = 0 \Rightarrow I(X_1; Y) \leq/\geq I(X_2; Y)$$

## ITIP (Information Theory Inequality Prover)



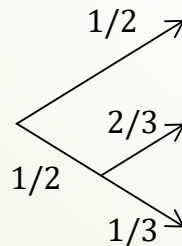
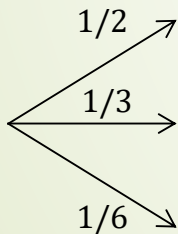
# Reference

- C. E. Shannon: A mathematical theory of communication. Bell System Technical Journal, vol. 27, pp. 379–423 and 623–656, July and October 1948
- Cover, T., & Thomas, J. (2006). Elements of information theory (2nd ed.). Hoboken, N.J.: Wiley-Interscience.
- Source coding with known distribution
  - Entropy coding
  - Huffman coding
- Universal source coding
  - Lempel-Ziv algorithms (LZ77 and LZ78)
  - Deflate algorithm (Huffman coding + LZ77) commonly used in ZIP, gzip, and PNG.

# Entropy is the only function $H$ that satisfies...

1.  $H$  is a continuous function of a probability distribution  $[p_1, p_2, \dots, p_n]$ .
2. If  $p_i = \frac{1}{n}$  for all  $i$ , then  $H$  is an increasing function of  $n$ .
3. If a choice is broken down into successive choices, the original  $H$  should be a weighted sum of the individual values of  $H$ . E.g.,

$$H\left(\frac{1}{2}, \frac{1}{3}, \frac{1}{6}\right) = H\left(\frac{1}{2}, \frac{1}{2}\right) + \frac{1}{2}H\left(\frac{2}{3}, \frac{1}{3}\right)$$



$H = K \sum_i p_i \log \frac{1}{p_i}$  unique up to a scaling factor  $K > 0$ .

# Why measure information in “bits”?

*"The choice of a logarithmic base corresponds to the choice of a unit for measuring information. If the base 2 is used the resulting units may be called **binary digits, or more briefly bits**, a word suggested by J. W. Tukey."*

Shannon July 1948



John Wilder Tukey  
(1915-2000)



# Why the name “entropy”?

*“My greatest concern was what to call it. I thought of calling it ‘information’, but the word was overly used, so I decided to call it ‘uncertainty’. When I discussed it with John von Neumann, he had a better idea. Von Neumann told me, ‘You should call it entropy, for two reasons. In the first place your uncertainty function **has been used in statistical mechanics under that name**. In the second place, and more importantly, **no one knows what entropy really is**, so in a debate you will always have the advantage.’”*

Shannon April 1961



John von Neumann  
(1903-1957)



# Information theory and machine learning

*"Turing and I had an awful lot in common, and we would talk about that kind of question. He had already written his famous paper about Turing Machines, so called, as they call them now, Turing Machines. They didn't call them that then. And we spent much time discussing the concepts of what's in the human brain. How the brain is built, how it works and what can be done with machines and **whether you can do anything with machines that you can do with the human brain** and so on. And that kind of thing. And I had talked to him several times about my notions on Information Theory, I know, and he was interested in those."*

Shannon July 1982



Alan Mathison Turing  
(1912–1954)

# Concerning learning machines

*"If, as is usual with computing machines, the operators of the machine itself could alter its instructions, there is the possibility that a **learning process could by this means completely alter the programme** in the machine..."*

*In many types of investigation, e.g., in the theory of information, it is a useful **assumption that 'computing costs nothing'**. It is important however **not to let this assumption become a belief**. In particular when one is considering brains and computers the assumption, and the theories based on it, are not applicable..."*

Turing contribution at conference on Information theory 1950



Alan Mathison Turing  
(1912–1954)