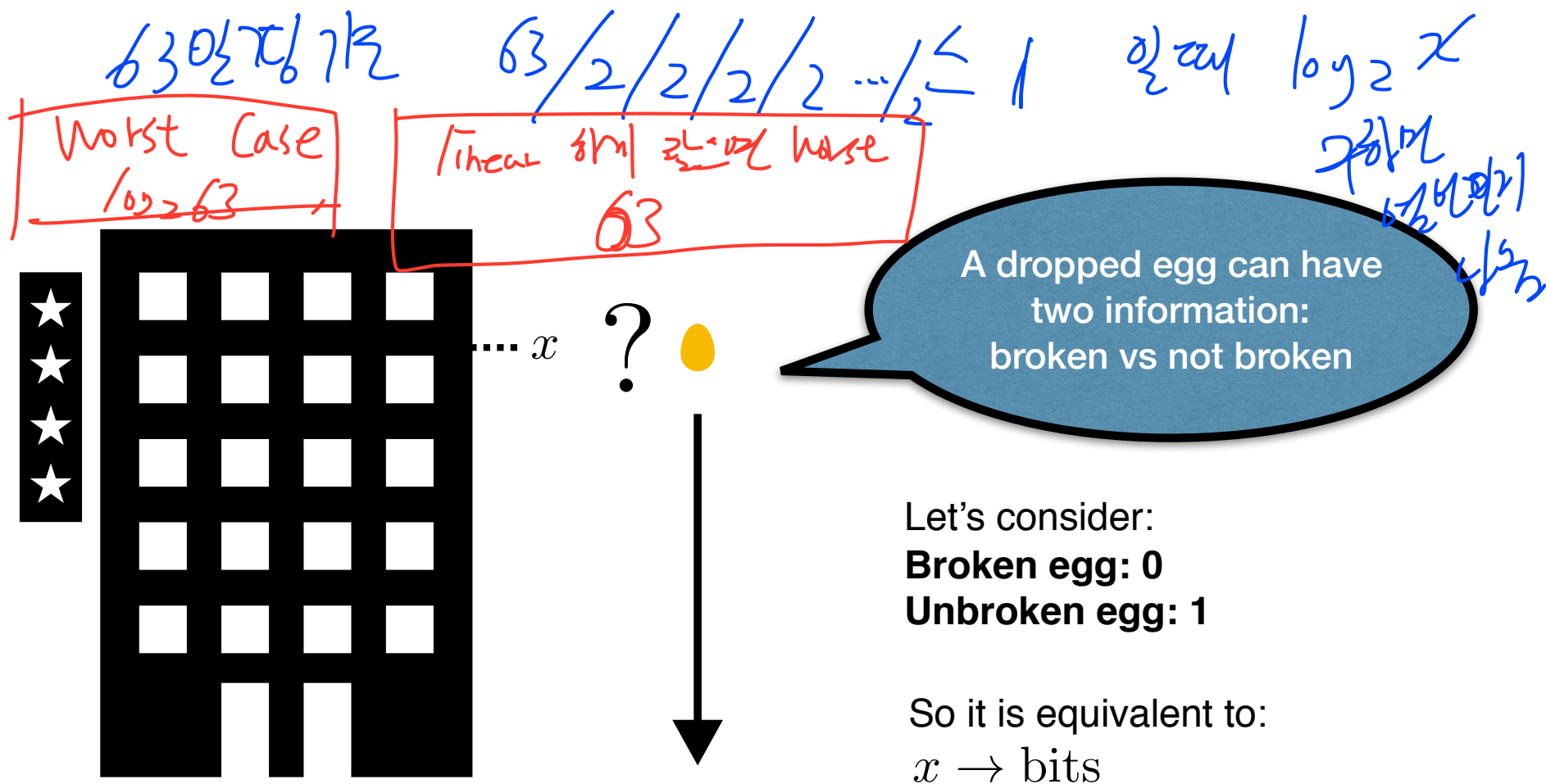


# Introduction to Computer Science & Engineering

## Lecture 3.5: Data Compression Theory

Jeonghun Park

# General Compression Theory



$\log_2 x$

Binary Search

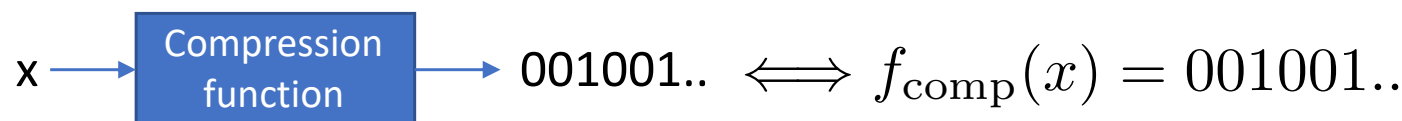
A basic concept of compression

# Data Compression

압축

- Definition

- ▶ Represent a general data (let's call  $x$ ) by using a finite number of ***bits***



- Target

- ▶ No confusion! (= no duplication!)
  - ← *호의 많은 value 중 찾아야 함*
- ▶ Design a compression function ***cleverly***, so that the output bits' length is minimized (= decrease the number of eggs)
  - 명리하게.*

# Source Coding Basic

데이터를 이분화  
Compression

- We call the input  $x$  as “source,” and the output bits (let’s write  $b = (001001..)$ ) as “code”
- Let’s assume that the input  $x$  is as follows:
  - ▶  $x$  can be 1,2,3, or 4
  - ▶ The probability is:  $\mathbb{P}[x = 1] = \frac{1}{2}$
  - $\mathbb{P}[x = 2] = \frac{1}{4}$
  - $\mathbb{P}[x = 3] = \frac{1}{8}$
  - $\mathbb{P}[x = 4] = \frac{1}{8}$
  - ▶ What is an efficient compression function?

# Source Coding Basic

- Let's consider this
  - ▶ Okay, the number of the sources is 4 (1,2,3,4)
  - ▶ Then, we may design a compression function as:

$$f_{\text{comp}}(x = 1) = \overset{\text{length}}{\textcircled{00}} = \frac{1}{4}$$

$$f_{\text{comp}}(x = 2) = 01 = \frac{1}{4}$$

$$f_{\text{comp}}(x = 3) = 10 = \frac{1}{4}$$

$$f_{\text{comp}}(x = 4) = 11 = \frac{1}{4}$$

Uniform 한 경우  
length

$$\left(\frac{1}{4} \times 2\right) + \cancel{1} + \cancel{1} + \cancel{1} = 2$$

(bits / length)

- ▶ Check point: There is no confusion
- ▶ We just complete to design **our first compression function**

	length
0	1
0 1	2
1 0	2
1 0 0	3

# Efficiency

- What is the expected number of bits length of our first design?

► This is easy. The bits length is **2**

- Can we do better?

► Let's consider the following:

$$f_{\text{comp}}(x = 1) = 0$$

$$f_{\text{comp}}(x = 2) = 10$$

$$f_{\text{comp}}(x = 3) = 110$$

$$f_{\text{comp}}(x = 4) = 111$$

general 하게  
답은 아니다.

(평균 length를 알려면  
각각 사용하는 거에 따른  
가중치를 곱해서  
평균을 구하면 된다.)

► The expected length is  $1 \times \frac{1}{2} + 2 \times \frac{1}{4} + 3 \times \frac{1}{8} + 3 \times \frac{1}{8} = \mathbf{1.75}$

$l_1$     $l_2$     $l_3$     $l_4$

# Fundamentals

- The insights behind this technique is:
  - ▶ Give a few bits to sources that **less appear**
- One may wonder..
  - ▶ Is there the optimal bits length for a given source?
  - ▶ And if so, how to design that?
- The above questions are fundamental questions where we try to find the answers

해결하려면 한쪽의  
길이/2의 값을  
나눠야

# Kraft Inequality

↳ 무조건 만족해야함

- Assume that we have  $m$  number of sources. Then also assume that our codes' length is  $l_1, l_2, \dots, l_m$ .
- Then the following should be satisfied

$$\sum_{i=1}^m 2^{-l_i} \leq 1$$

$$l_1 = 1$$

$$l_2 = 2$$

$$l_3 = 3$$

$$l_4 = 3$$

$$\sum_{i=1}^m 2^{-l_i}$$

$$2^{-1} + 2^{-2} + 2^{-3} + 2^{-3}$$

$$= \frac{1}{2} + \frac{1}{4} + \frac{1}{8} + \frac{1}{8} = 1$$

장/2의 값이 1이므로 Best Case.





# Proof

ex) 00 쓰면 000이나 001은 갖지 못함  
 $\Rightarrow$  Prefix Condition

- Consider a binary tree, where each branch represents each code. Then, by the **prefix condition** (for unique decodability), no code is an ancestor of any other code on this tree. Hence, each code eliminates its descendants as possible codes.
- Let  $l_{\max}$  be the length of the longest codes. A code at level  $l_i$  has  $2^{l_{\max}-l_i}$  descendants at level  $l_{\max}$   
 $\hookrightarrow$  disjoint region
- Since each of the above mentioned descendant must be disjoint, which implies that

$$\sum_{i=1}^m 2^{l_{\max}-l_i} \leq \underline{2^{l_{\max}}}$$

$l_{\max}$ 은 뭐라도 상관 X

# Some Notes


- Kraft inequality is a clue of the optimal compression function
- If there is a set of codes, and this code length satisfies Kraft inequality, this is the optimal code!

$$\sum_{i=1}^m 2^{-\ell_i} = 1$$

# Optimal Codes

- The expected **code length** is calculated as

$$L = \sum_{i=1}^m p_i \ell_i$$


 Probability of the  $i$ -th source

- And from Kraft inequality, we have  $\sum_{i=1}^m 2^{-\ell_i} \leq 1$
- Jointly considering the above two conditions, we obtain the optimal code length is  $(\ell_i)^* = \log_2 p_i^{-1}$

Also fits our intuition!

# Achievability

- Okay, now we know the optimal code length.
- But how to achieve..? This is a different story.

# Huffman Coding

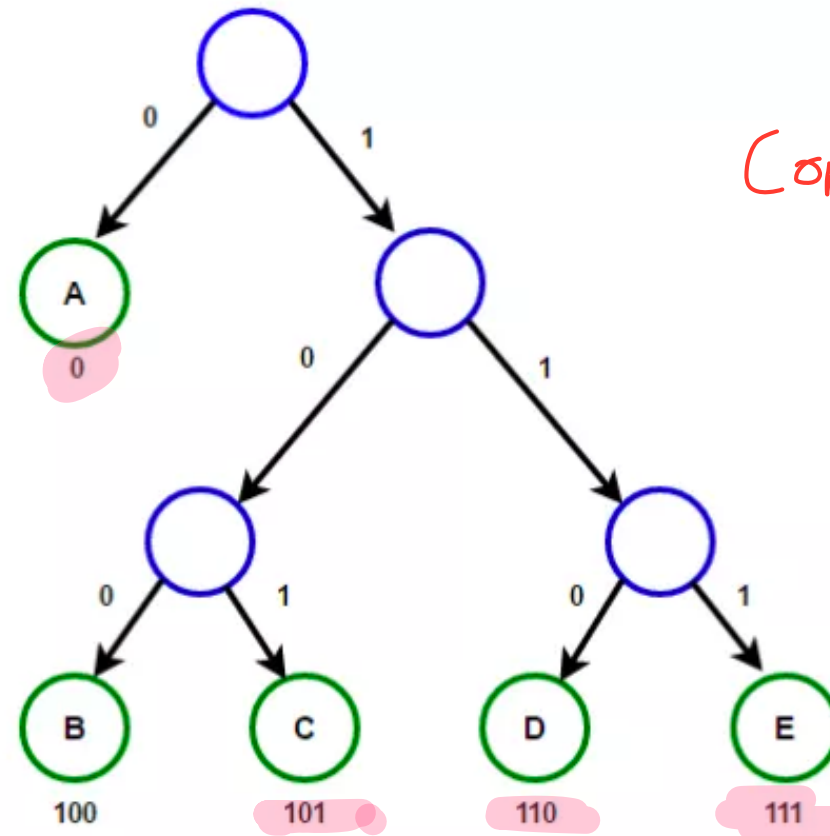
⇒ 파일 용량 줄이기 실험

(Copy, PNG)

- David Huffman's student project result
  - ▶ During his Ph.D. course!
  - ▶ This achieves the optimal code length

# Huffman Coding

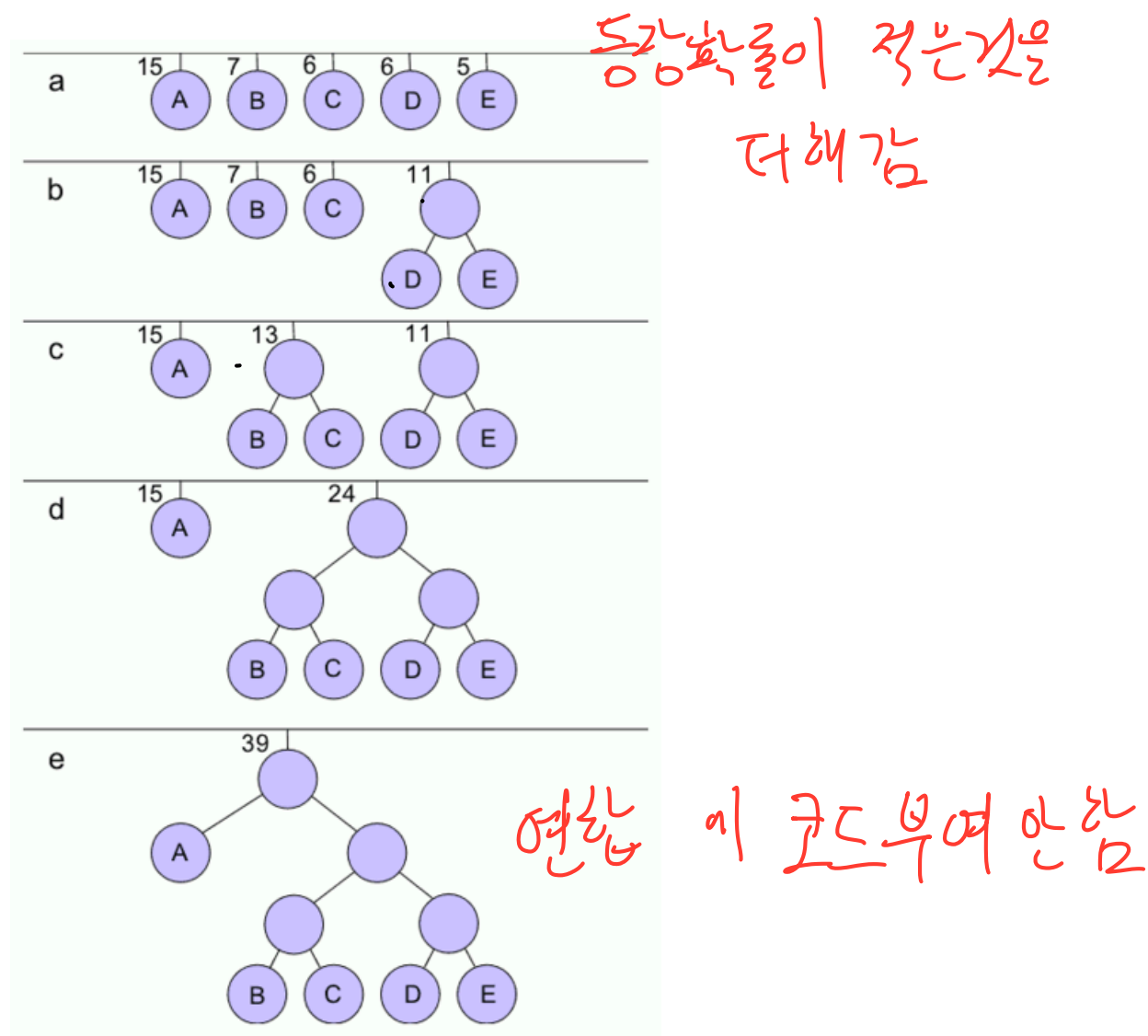
- Huffman coding's key is Huffman tree



Combined 된 연산체인  
부여 X

# Huffman Coding

더하고 ordering 다시  
... 빈칸



# How to Construct Huffman Tree

- Algorithm description

1. Order the sources by their appear probability

2. Add the most less appearing sources

- They will make a branch

3. Go to step 1 and repeat

0, 1 ~~0.25~~ 0.3 0.4 0.25

0, 1, 2, 3 0.25 0.3 0.4 0.25

(최소인 optimal 을  
찾는다)

Codeword				
Length	Codeword	X	Probability	
2	01	1	0.25	0.3
2	<u>10</u>	2	0.25	0.25
2	<u>11</u>	3	0.2	0.25
3	<u>000</u>	4	0.15	0.2
3	<u>001</u>	5	0.15	