# *Greedy Algorithms*

**Heejin Park**

*Hanyang University*

# Contents

- **Introduction**

- **An activity selection problem**

- **Elements of the greedy strategy**

- **Huffman codes**

# Introduction

- A *greedy algorithm* always makes the choice that looks best at the moment.

- It makes a locally optimal choice in the hope that this choice will lead to a globally optimal solution.

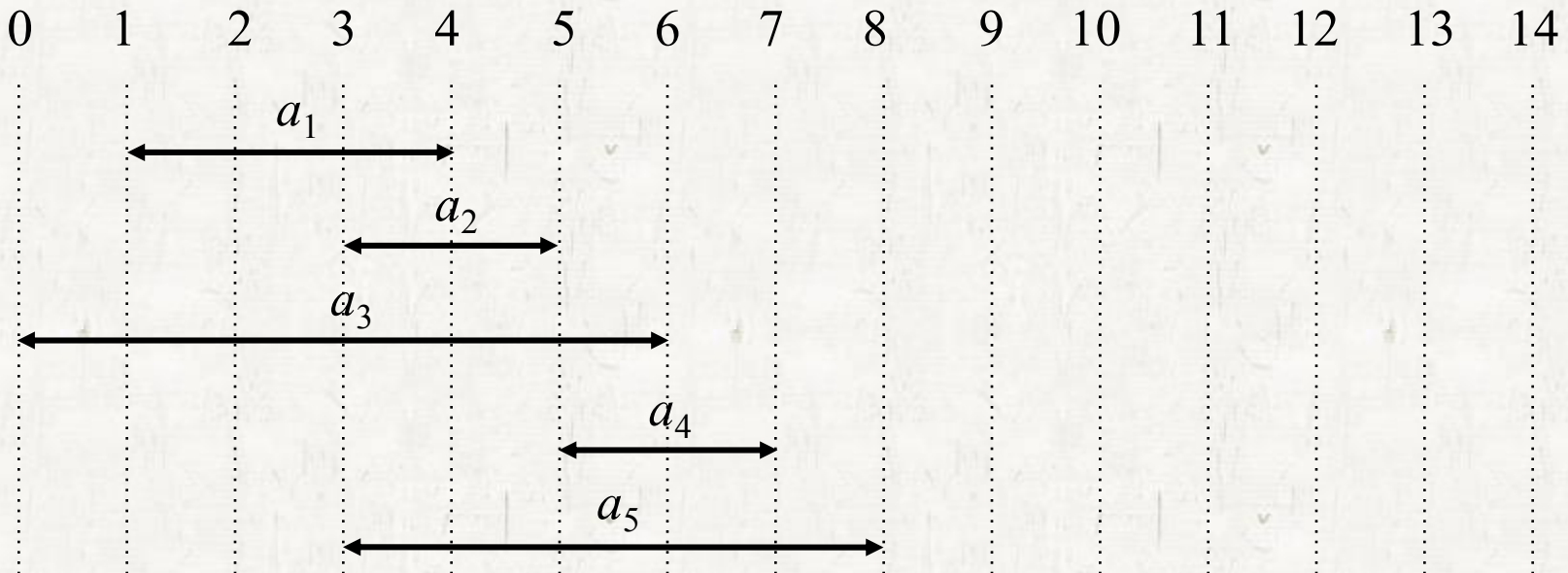- It makes the choice *before* the subproblems are solved.

# An activity selection problem

- **An activity selection problem**
  - To select a maximum-size subset of mutually compatible activities.
  - For example
    - Given $n$ classes and 1 lecture room,
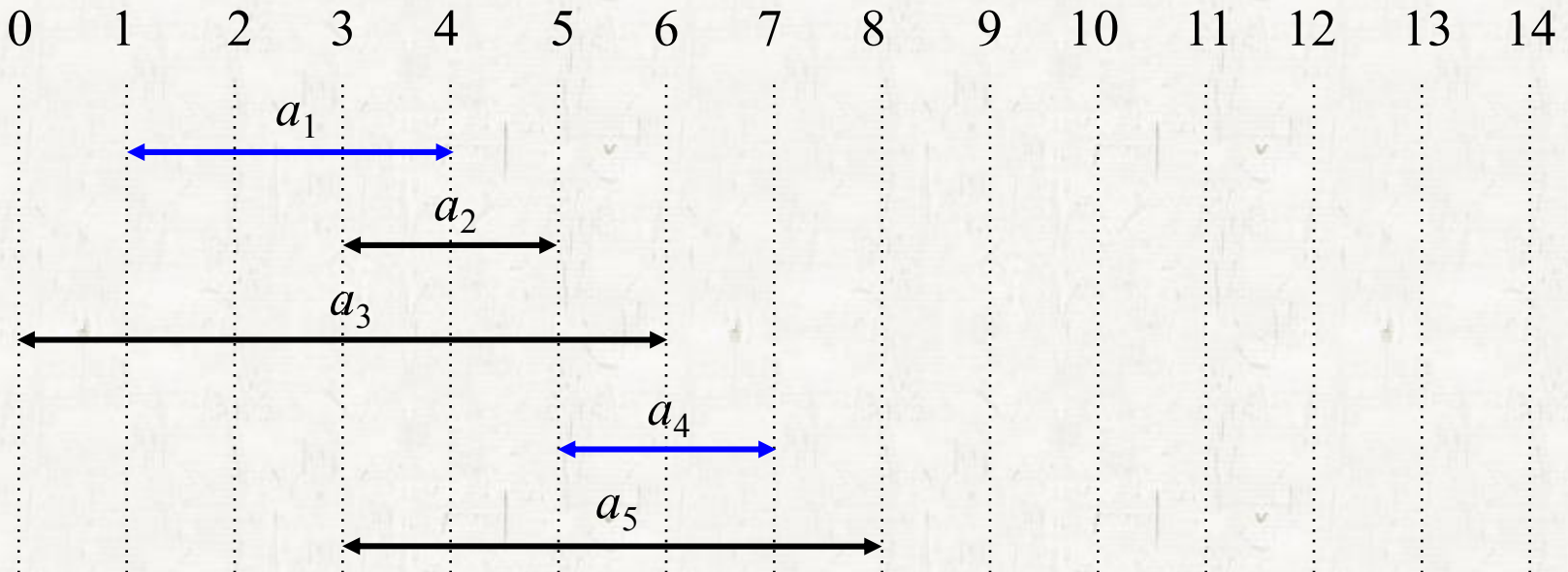    - to select the maximum number of classes

# An activity selection problem

- A set of *activities*: $S = \{a_1, a_2, ..., a_n\}$

- Each activity $a_i$ has its *start time $s_i$* and *finish time $f_i$*.
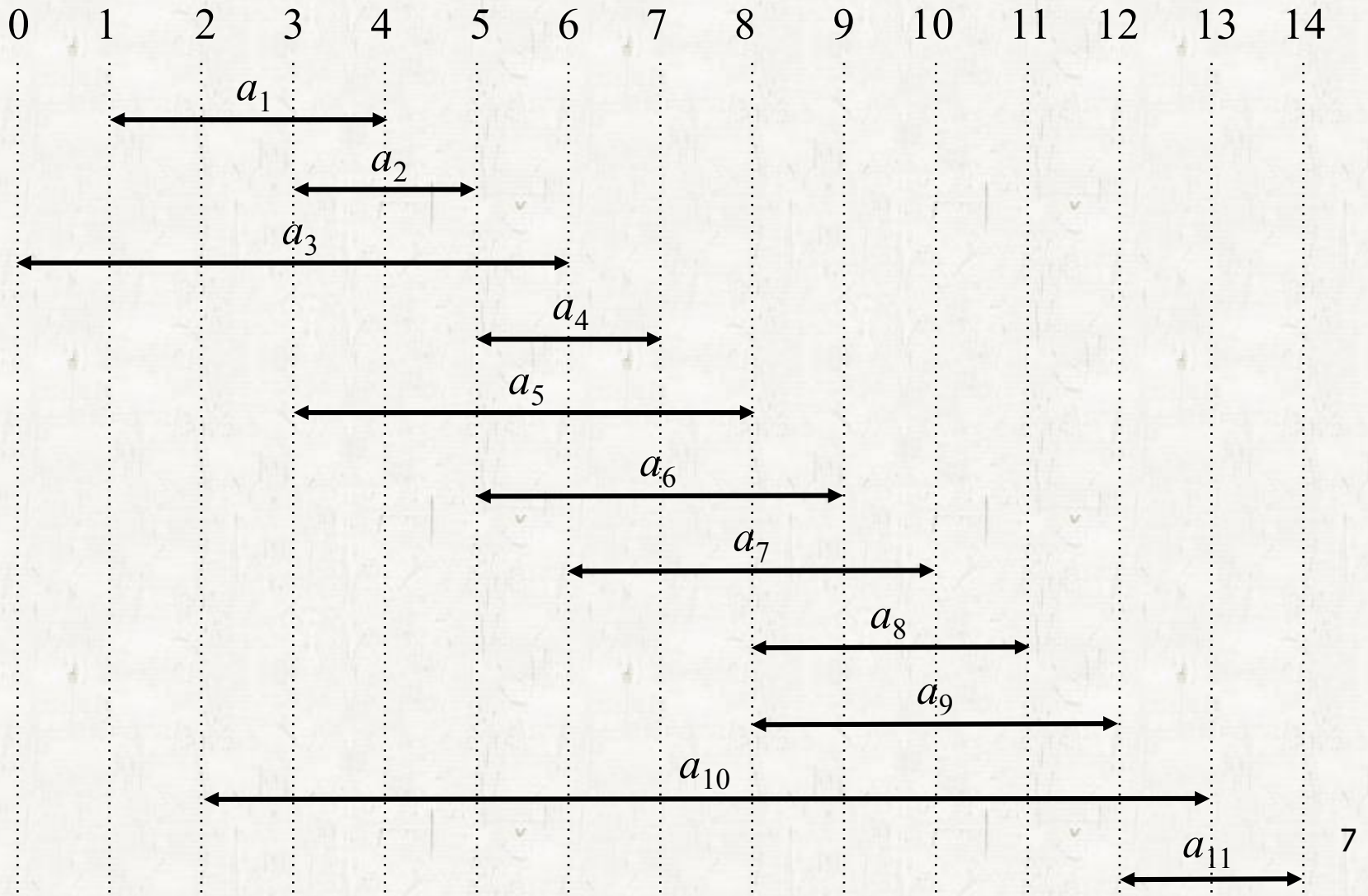  - $0 \le s_i < f_i < \infty$

# An activity selection problem

- Activity $a_i$ takes place during $[s_i, f_i)$

- Activities $a_i$ and $a_j$ are ***compatible***
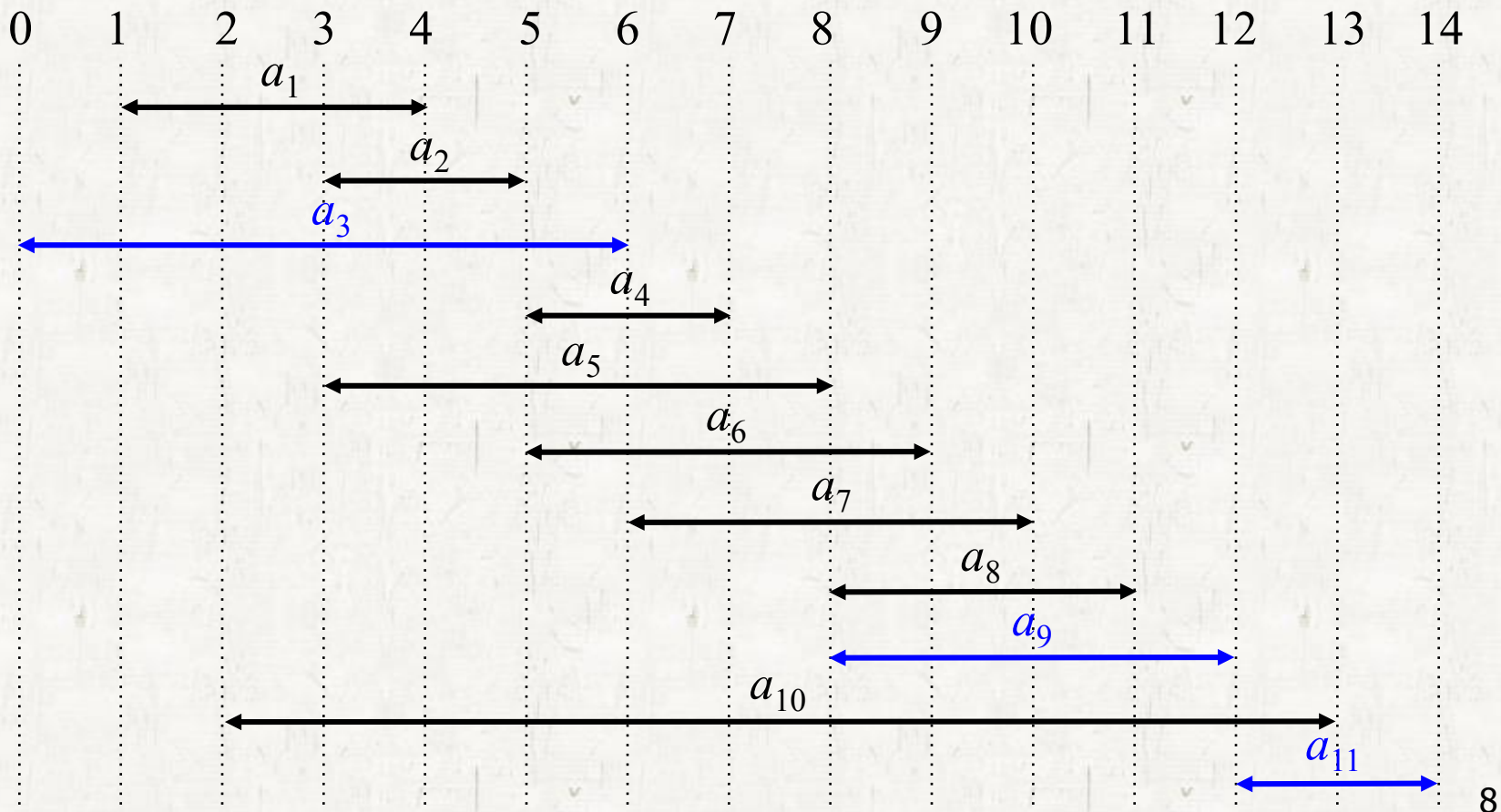  if the intervals $[s_i, f_i)$ and $[s_j f_j)$ do not overlap.

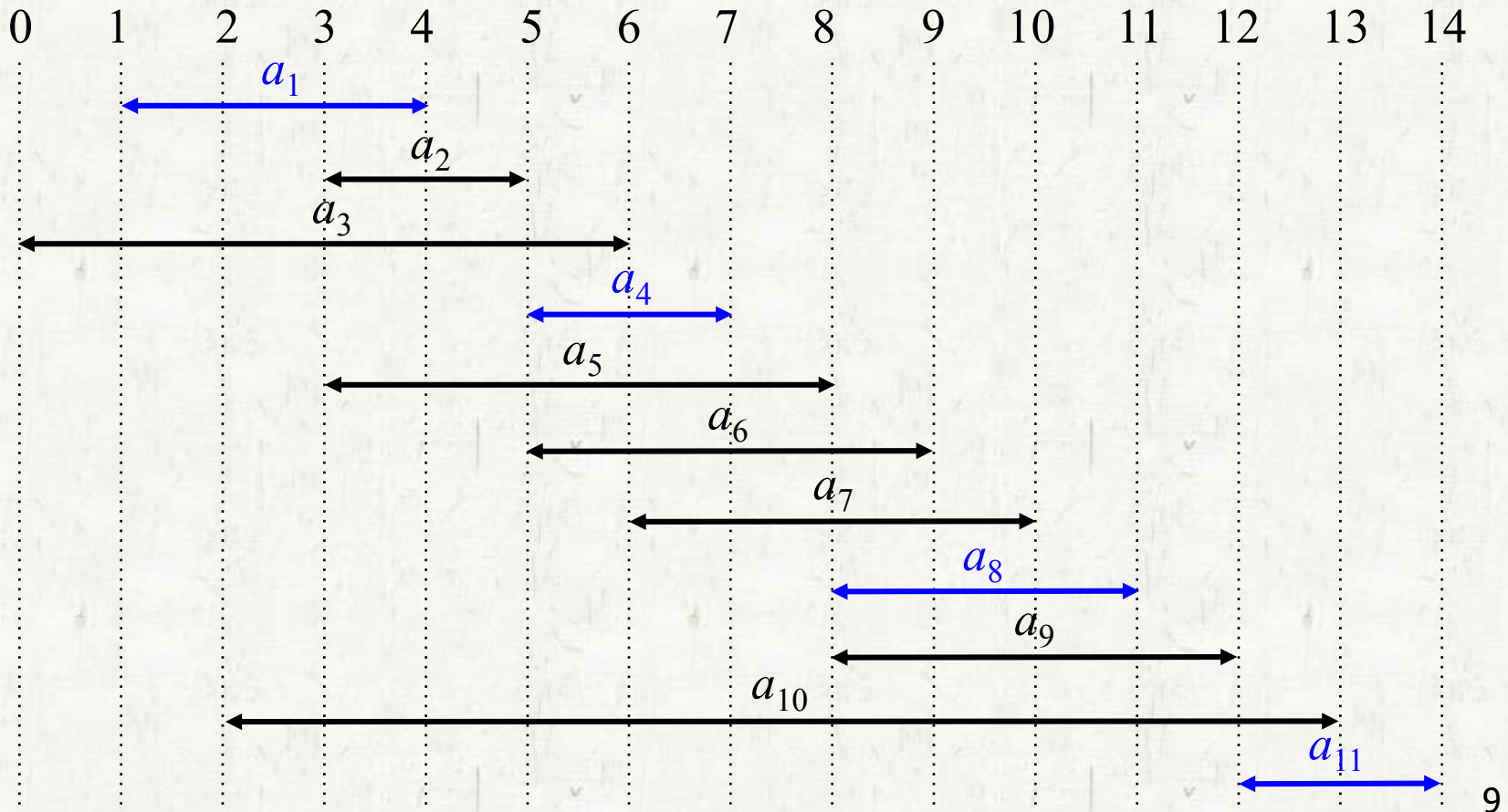# An activity selection problem

# An activity selection problem

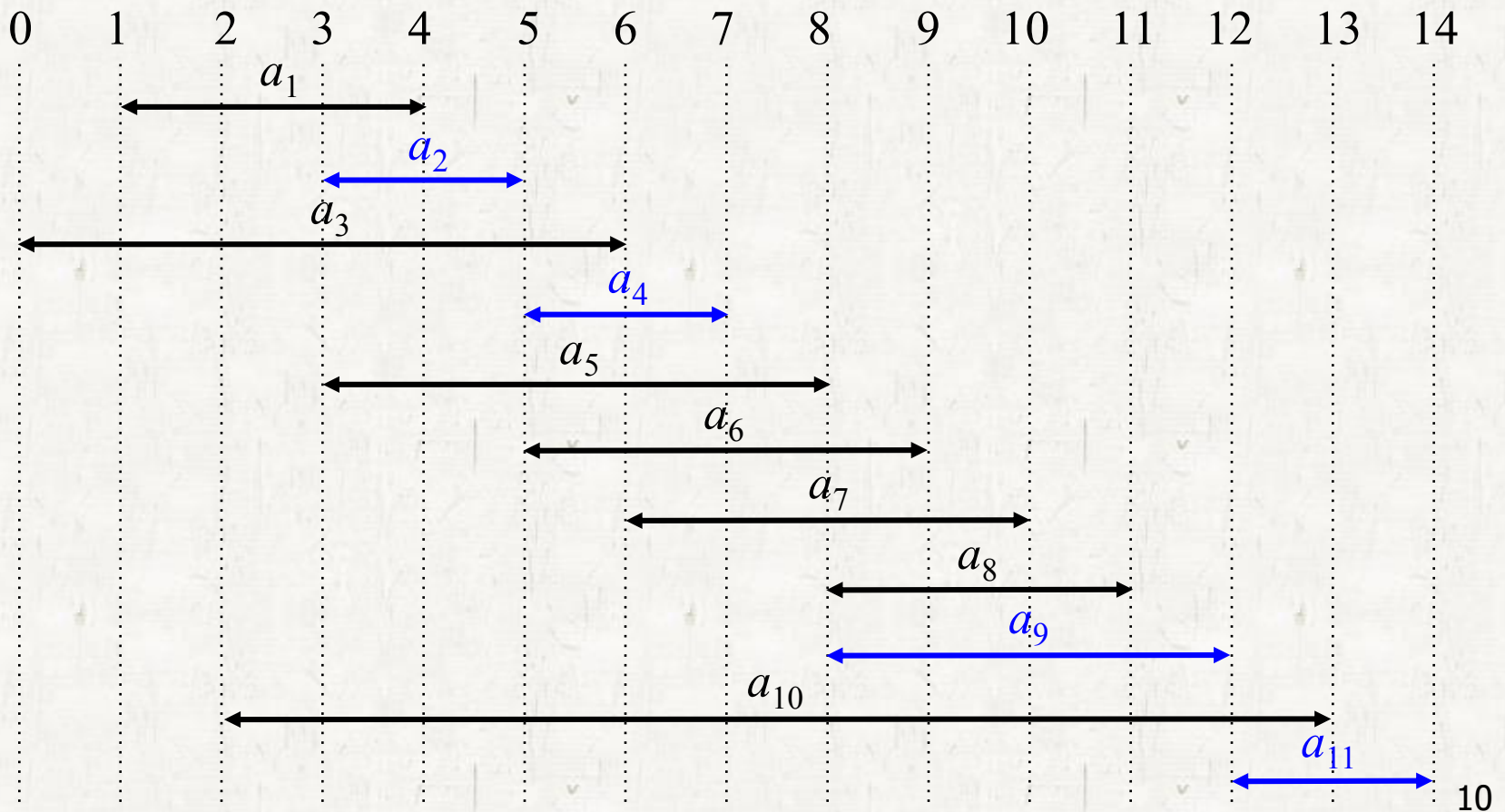- $\{a_3, a_9, a_{11}\}$: mutually compatible activities, not a largest set

# An activity selection problem

- $\{a_1, a_4, a_8, a_{11}\}$: A largest set of mutually compatible activities

# An activity selection problem

- $\{a_2, a_4, a_9, a_{11}\}$: Another largest subset

# An activity selection problem

- **Optimal substructure**
  - Assume that activities are sorted in increasing order of finish time.

$$f_0 \leq f_1 \leq f_2 \leq \cdots \leq f_n < f_{n+1}$$

| $i$ | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 | 9 | 10 | 11 |
|-----|---|---|---|---|---|---|---|---|---|----|----|
| $s_i$ | 1 | 3 | 0 | 5 | 3 | 5 | 6 | 8 | 8 | 2 | 12 |
| $f_i$ | 4 | 5 | 6 | 7 | 8 | 9 | 10 | 11 | 12 | 13 | 14 |

# An activity selection problem

- **Greedy algorithm**
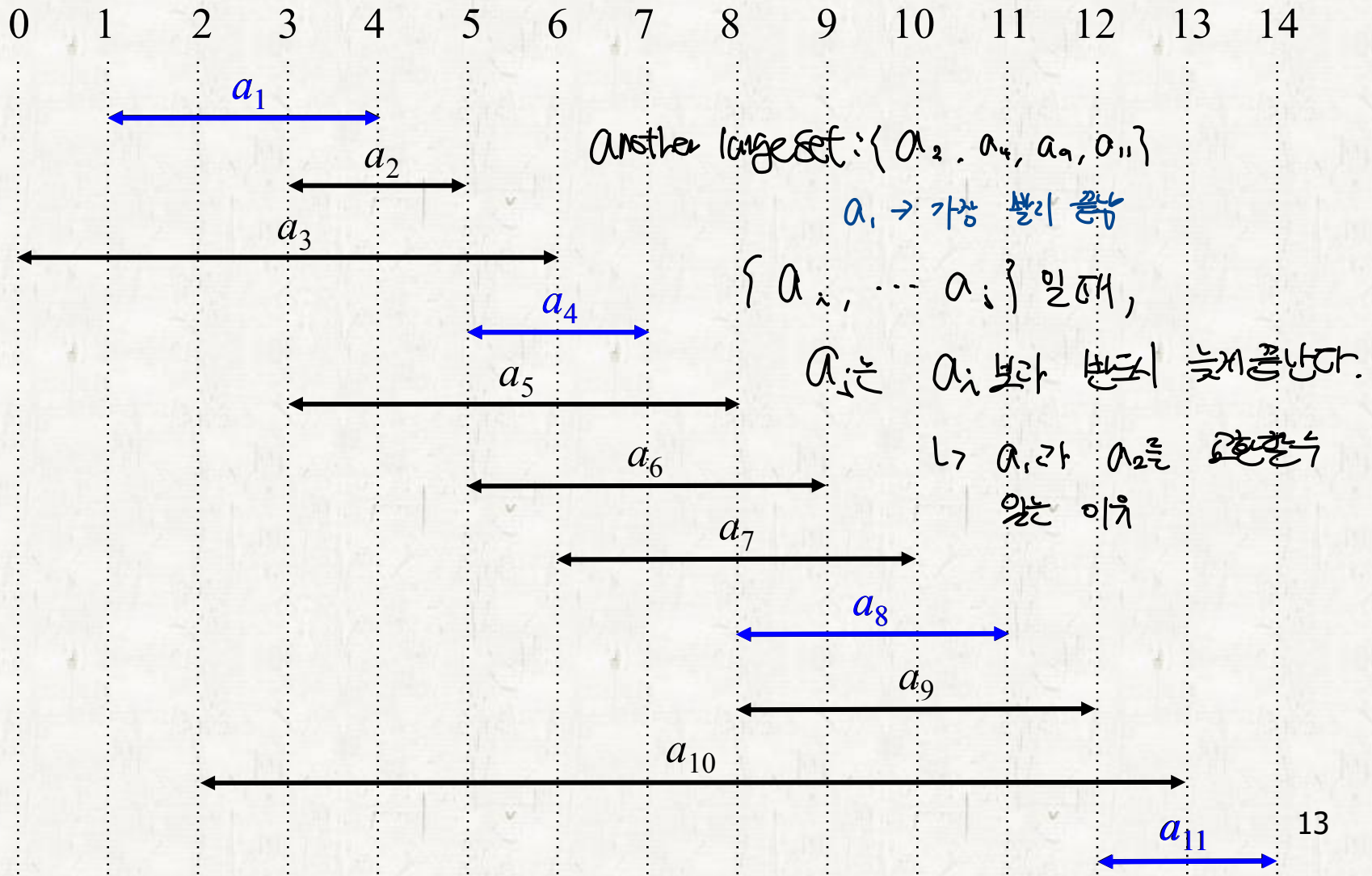  - Select the earliest finishing activity one by one.

Running Time : $\Theta(n)$

ㅅㄱH시 저정렬 안해 된상

new activity 라 latest activity를 비교

# An activity selection problem

0    1    2    3    4    5    6    7    8    9    10    11    12    13    14

$a_1$

$a_2$

$a_3$

Another largest : {$a_2$, $a_4$, $a_9$, $a_{11}$}

$a_1$ → 가장 빨리 끝남

$a_4$

$a_5$

{$a_i$, ··· $a_j$} 일때,

$a_j$는 $a_i$보다 반드시 늦게끝난다.

$a_6$

$a_7$

↳ $a_1$과 $a_2$를 교환할수 없는 이유

$a_8$

$a_9$

$a_{10}$

$a_{11}$

# Contents

- *Introduction*

- *An activity selection problem*

- **Elements of the greedy strategy**

- **Huffman codes**

14

# Elements of the greedy strategy

- **Greedy-choice property**
  - Make the choice *before* the subproblems are solved.
  - Only one subproblem is generated.

- **Dynamic programming**
  - Make a choice *after* the subproblems are solved.
  - Several subproblems may be generated.

# Elements of the greedy strategy

- **Greedy vs. Dynamic programming**
  - **0-1 knapsack**
    - A thief robbing a store finds $n$ items.
    - The $i$th item is worth $v_i$ dollars and weighs $w_i$ pounds.
    - He can carry at most $W$ pounds in his knapsack.
    - The $n$, $v_i$, $w_i$, and $W$ are integers.
    - Which items should he take?

  - **Fractional knapsack**
    - In this case, the thief can take <u>fractions of items</u>.
      Item의 나어질 수 있다.

# Elements of the greedy strategy

**Fractional knapsack**

- The greedy strategy works.
- Compute the value per pound $v_i/w_i$ for each item.
- Take as much as possible of the item with the greatest value per pound.

$$\boxed{\frac{V_1}{W_1}} \;=\; \frac{V_2}{W_2} \;\geq\; \frac{V_3}{W_3}$$

①  $\longrightarrow$

# Elements of the greedy strategy

- **0-1 knapsack**
  - The greedy strategy does not work.

fraction 0-1 knapsack
→ Greedy work!

Empty



| | | | |
|---|---|---|---|
| item 1 | item 2 | item 3 | |
| 10 | 20 30 | 30 50 | 50 |
| $60 | $100 | $120 | Knapsack |
| | | | (a) |

30 $120
+
20 $100
=$220

20/30 $80
+
20 $100
+
10 $60
= $240

30 $120
+
20 $100
+
10 $60
=$160

30 $120
+
20 $100
+
10 $60
=$180

Greedy

(b)     (c)

18

# Self-study

- **Exercise 16.2-1**

- **Exercise 16.2-2**

- **Exercise 16.2-5**

- **Exercise 16.2-7**

# Contents

- *Introduction*

- *An activity selection problem*

- *Elements of the greedy strategy*

- **Huffman codes**

# Huffman codes

- **Huffman Codes**
  - A widely used technique for compressing data.

$\rightarrow 3 \times 100,000 = 300,000$

- Consider representing 100,000 characters from {a, b, c, d, e, f}.
  - 3-bit *fixed-length code* is used in general.
  - It takes 300,000 bits in total

|  | a | b | c | d | e | f |
|---|---|---|---|---|---|---|
| Frequency (in thousands) | 45 | 13 | 12 | 16 | 9 | 5 |
| Fixed-length codeword | 000 | 001 | 010 | 011 | 100 | 101 |

# Huffman codes

- We can reduce the space if *variable-length code* is used.

|                            | a   | b   | c   | d   | e    | f    |
|----------------------------|-----|-----|-----|-----|------|------|
| Frequency (in thousands)   | 45  | 13  | 12  | 16  | 9    | 5    |
| Fixed-length codeword      | 000 | 001 | 010 | 011 | 100  | 101  |
| Variable-length codeword   | 0   | 101 | 100 | 111 | 1101 | 1100 |

- Shorter **codewords** for frequent characters.

- 224,000 bits in total
  - $(45 \cdot 1 + 13 \cdot 3 + 12 \cdot 3 + 16 \cdot 3 + 9 \cdot 4 + 5 \cdot 4) \cdot 1000$ bits

# Huffman codes

*radix Tree*



- **Encoding and decoding of variable-length code**
  - Encoding abc : $0 \cdot 101 \cdot 100$ : $abc \rightarrow 0101100$
  - Decoding 001011101 $=$ aabe!
    - $0 \cdot 0 \cdot 101 \cdot 1101$: aabe

prefix 관계 X
→ 정확 decoding 된다

|  | a | b | c | d | e | f |
|---|---|---|---|---|---|---|
| Variable-length codeword | 0 | 101 | 100 | 111 | 1101 | 1100 |

- Decoding 001 when a: 0 b: 01 c: 1

이런건, a가 b의 prefix

  - 001: aac or ab  → 둘 다 가능
  - The codeword 0 for a is a prefix of the codeword 01 for b.

23

# Huffman codes

**Prefix codes**

- No codeword is a prefix of some other codeword.

|  | a | b | c | d | e | f |
|---|---|---|---|---|---|---|
| Variable-length codeword | 0 | 101 | 100 | 111 | 1101 | 1100 |

# Huffman codes

- **Prefix codes**
  - 3-bit fixed-length code is also a prefix code.

we need
300,000 bits

→ 224,000 bits



$100 \rightarrow e$  )   $10 \rightarrow e$  )
$101 \rightarrow f$  )2   $11 \rightarrow f$  )2
고정길이 → fixed-length
0은 중요

- The left tree is a *full binary tree* while the right one is not. 0을 추가

→ full binary로
안함 (11→를 할수가
없다)

  - Every node is either leaf or has two children
  - A full binary tree for alphabet $C$ has $|C|$ leaves and $|C|-1$
    internal nodes.   {a,b,c,d,e,f}   `6
    알파벳

25

# Huffman codes

- **The cost of tree $T$**
  - $f(c)$: frequency of a character $c$
  - $d_T(c)$: length of the codeword for $c$

$$B(T) = \sum_{c \in C} f(c)d_T(c)$$

  - An optimal code is represented by a full binary tree.

# Huffman codes

- Huffman invented a greedy algorithm that constructs an optimal prefix code called an ***Huffman code***.

|                          | a  | b  | c  | d  | e | f |
| ------------------------ | -- | -- | -- | -- | - | - |
| Frequency (in thousands) | 45 | 13 | 12 | 16 | 9 | 5 |

| f : 5 | e : 9 | c : 12 | b : 13 | d : 16 | a : 45 |
| ----- | ----- | ------ | ------ | ------ | ------ |

# Huffman codes

Choose two least
frequency character

| f : 5 | e : 9 | c : 12 | b : 13 | d : 16 | a : 45 |

최소 2개를 찾기 위해 Θ(n)

Virtual node

→ 5+9

g, 14

0        1

| f : 5 | e : 9 |

| c : 12 | b : 13 | d : 16 | a : 45 |

build sub tree θ(1)

앞장에서   한   장   넘판

5개 중 frequency가 작은 2개

g, 14    c : 12    b : 13    d : 16    a : 45

$\theta(n-1)$

f : 5    e : 9

$\theta(1)$

g, 14    h, 25    d : 16    a : 45

0    1    0    1

f : 5    e : 9    c : 12    b : 13

29

# Huffman codes

**g, 14**

d : 16

**h, 25**

a : 45

1

0          1

f : 5          e : 9

c : 12          b : 13

**i, 30**

**h, 25**

a : 45

0          1

0          1

**g, 14**

d : 16

c : 12          b : 13

0          1

f : 5          e : 9

# Huffman codes

# Huffman codes

$$\theta(n) + \theta(n-1) + \cdots + \theta(1)$$
$$= \theta(n^2)$$

```
            k,100
          0/     \1
       a : 45    j, 55
                0/    \1
             h, 25    i, 30
            0/   \1   0/   \1
        c : 12  b : 13  g, 14  d : 16
                       0/   \1
                    f : 5   e : 9
```

# Huffman codes

| f : 5 | e : 9 | c : 12 | b : 13 | d : 16 | a : 45 |

- Min Heap

  → extract two least node

  $\theta(\log n)$

  $\theta(n \log n)$

  $\theta(n)$

  ```
        f, 5
       /    \
     e, 9   c, 12
     /  \      \
  b, 13 d, 16  a, 45
  ```

# Huffman codes

| f : 5 | e : 9 | c : 12 | b : 13 | d : 16 | a : 45 |

- Min Heap

# Huffman codes

| f : 5 | e : 9 | c : 12 | b : 13 | d : 16 | a : 45 |

- Min Heap

# Huffman codes

- Min Heap

# Huffman codes

- Min Heap

e, 9

b, 13

c, 12

a, 45

d, 16

f, 5

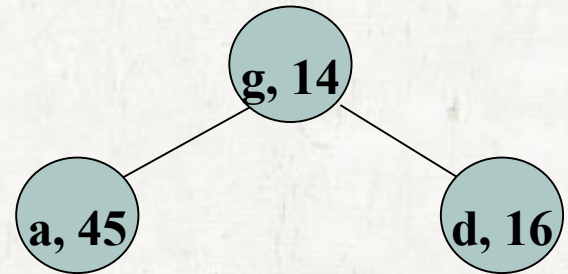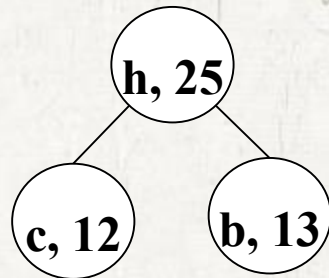# Huffman codes

- Min Heap
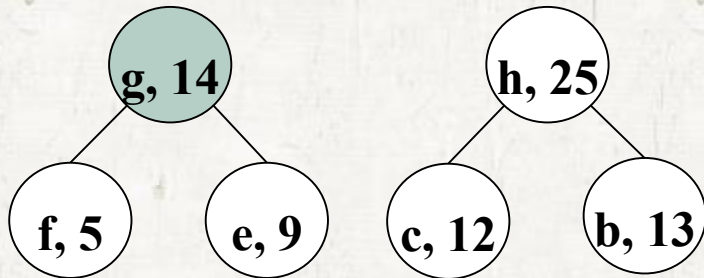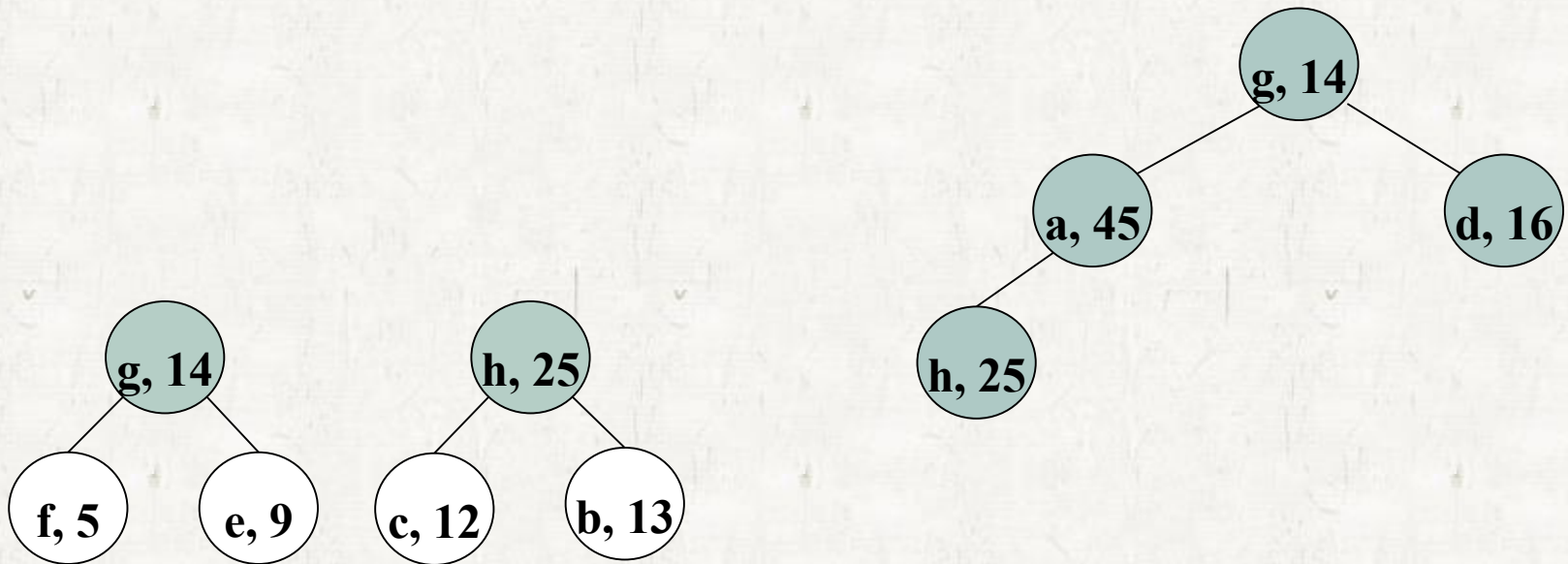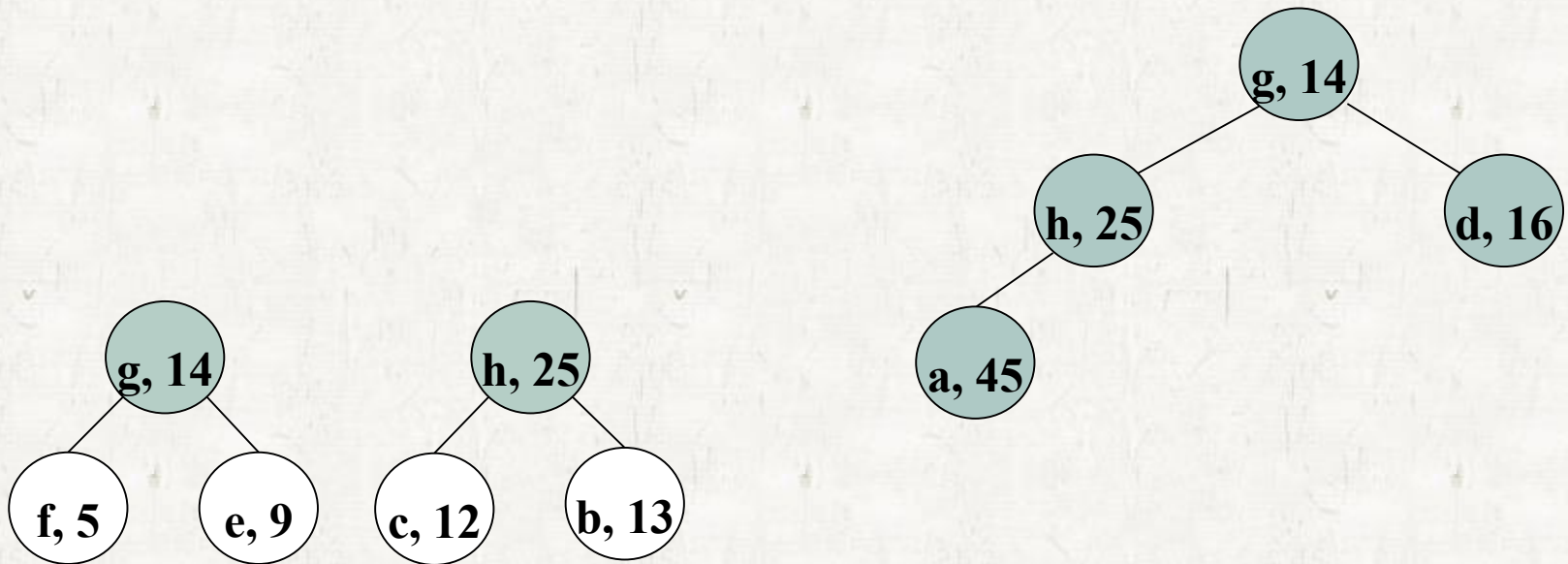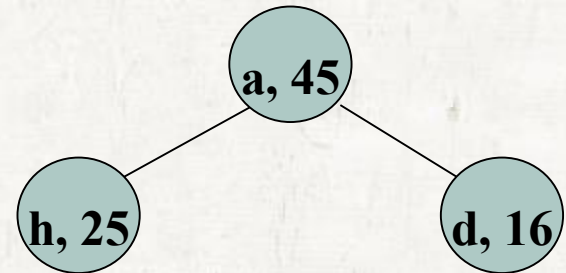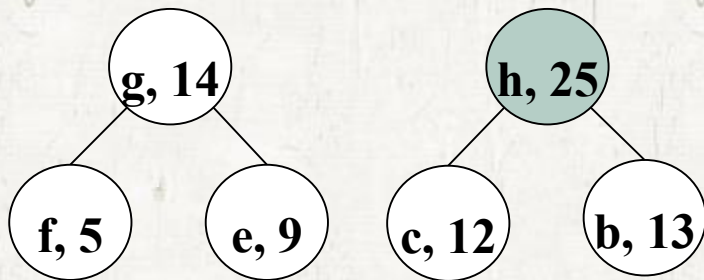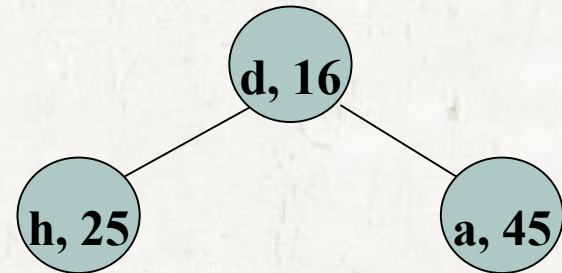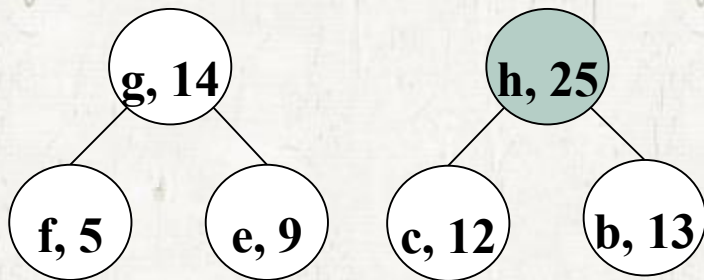
# Huffman codes

- Min Heap

# Huffman codes

- Min Heap

# Huffman codes

- Min Heap

# Huffman codes

- Min Heap

: 2개를 min heap로 추출하고
2개를 결합한 것을 삽입하여 } $\theta(\log n)$



(g, 14)를 분리
$\theta(\log n)$

# Huffman codes

- Min Heap

# Huffman codes

- Min Heap

g, 14

b, 13          d, 16

a, 45

g, 14

f, 5     e, 9     c, 12

# Huffman codes

- Min Heap

**b, 13**

**g, 14**

**d, 16**

**a, 45**

**g, 14**

**f, 5**

**e, 9**

**c, 12**

# Huffman codes

- Min Heap



g, 14

d, 16

a, 45

g, 14

f, 5    e, 9    c, 12    b, 13

# Huffman codes

- Min Heap

# Huffman codes

- Min Heap

g, 14
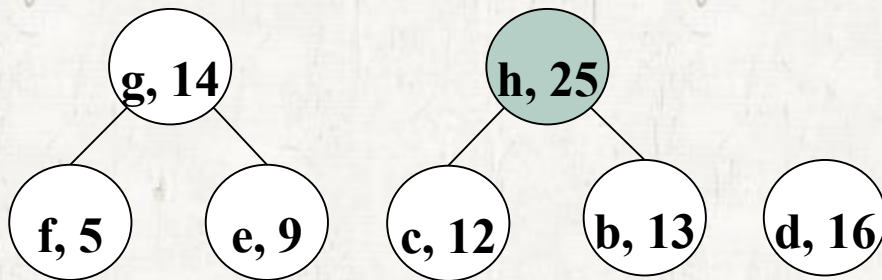
a, 45          d, 16

g, 14

f, 5     e, 9     c, 12     b, 13

# Huffman codes

- Min Heap

g, 14

a, 45          d, 16

g, 14          h, 25

f, 5     e, 9     c, 12     b, 13

# Huffman codes

- Min Heap

g, 14

a, 45

d, 16

h, 25

g, 14

f, 5   e, 9

h, 25

c, 12   b, 13

# Huffman codes

- Min Heap



g, 14
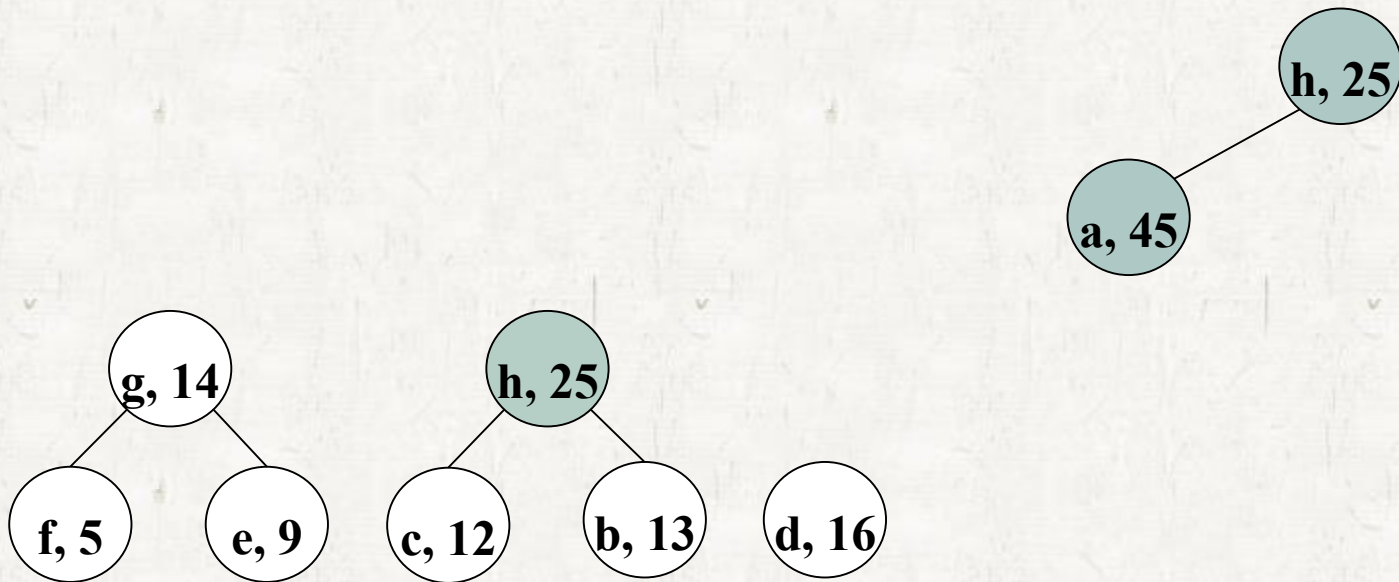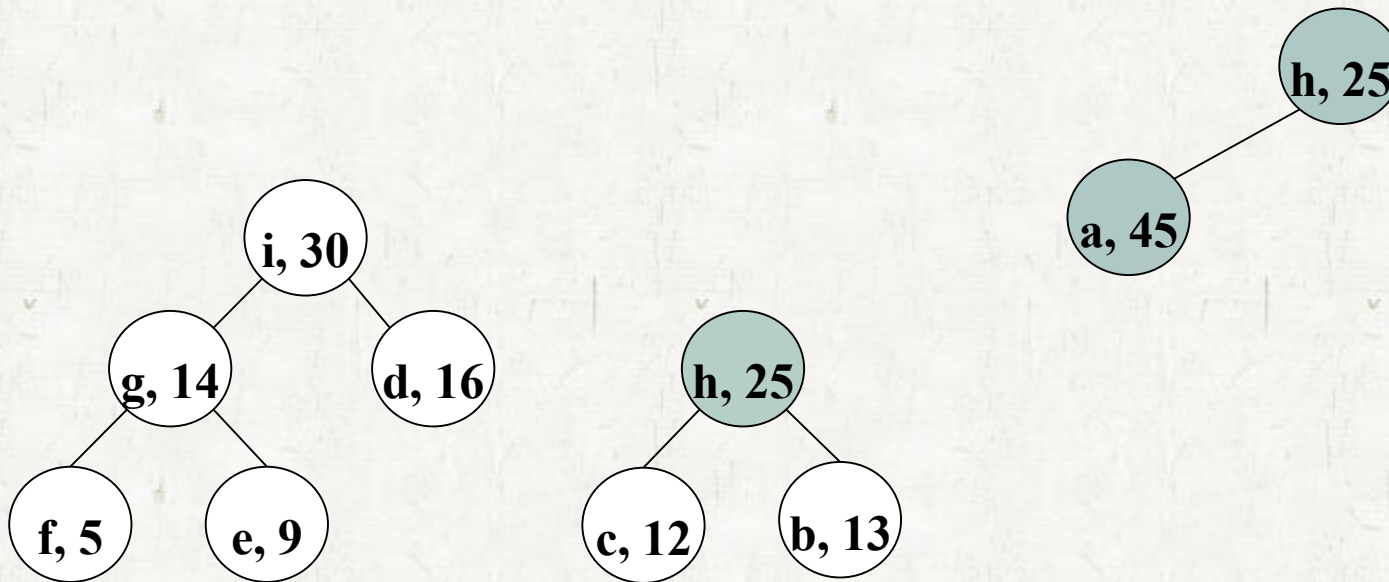
h, 25      d, 16

a, 45

g, 14          h, 25

f, 5    e, 9    c, 12    b, 13

# Huffman codes

- Min Heap
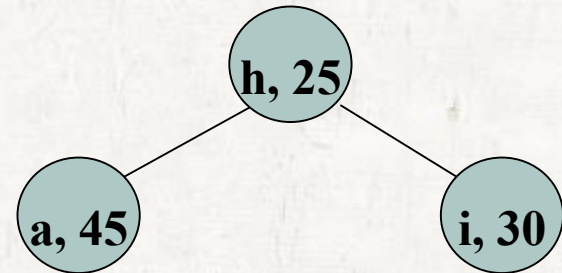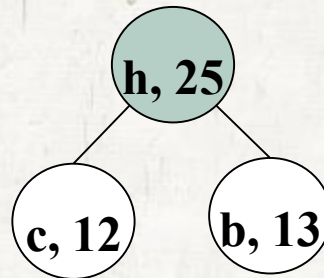
# Huffman codes

- Min Heap

d, 16

h, 25    a, 45

g, 14

h, 25

f, 5    e, 9    c, 12    b, 13

# Huffman codes

- Min Heap

h, 25          a, 45

g, 14          h, 25

f, 5    e, 9    c, 12    b, 13    d, 16

# Huffman codes

- Min Heap

a, 45

h, 25

g, 14

h, 25

f, 5

e, 9

c, 12

b, 13

d, 16

# Huffman codes

- Min Heap

h, 25

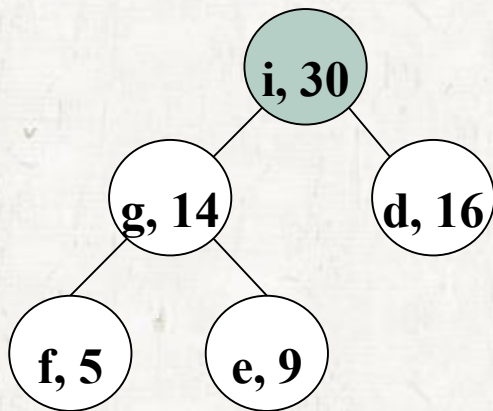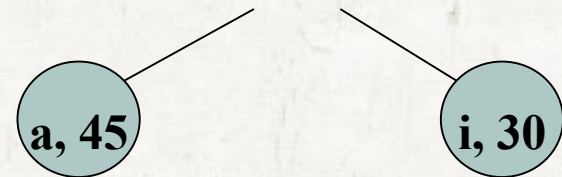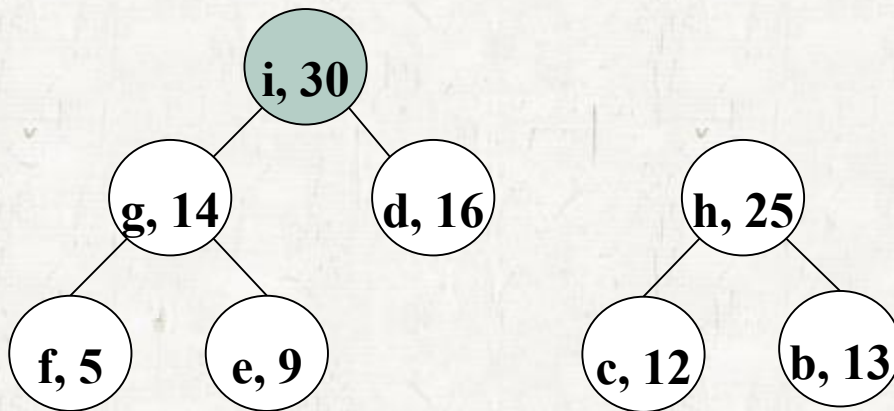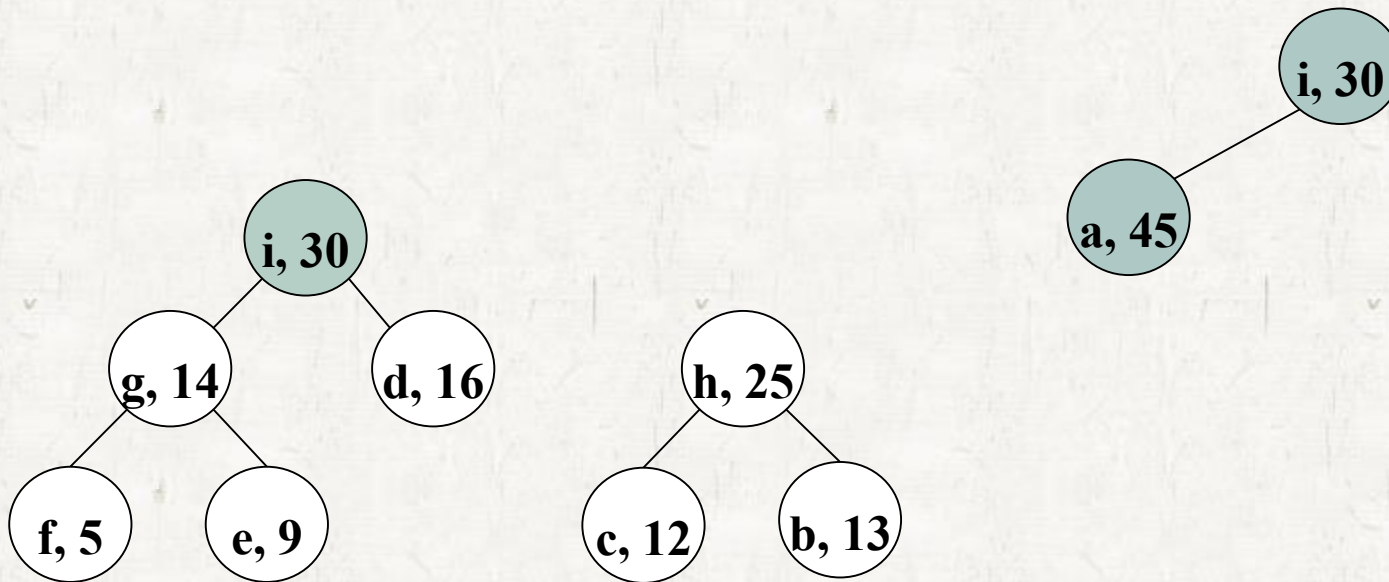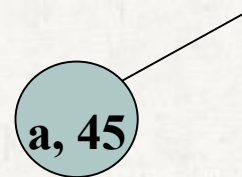a, 45

g, 14

f, 5    e, 9

h, 25

c, 12    b, 13    d, 16

# Huffman codes

- Min Heap

# Huffman codes

- Min Heap

# Huffman codes

- Min Heap

# Huffman codes

- Min Heap

# Huffman codes

- Min Heap

# Huffman codes

- Min Heap

**a, 45**

**i, 30**

**g, 14**     **d, 16**     **h, 25**

**f, 5**     **e, 9**     **c, 12**     **b, 13**

# Huffman codes

- Min Heap

# Huffman codes

- Min Heap



A min heap diagram showing two trees:

First tree rooted at j, 55:
- j, 55
  - h, 25
    - c, 12
    - b, 13
  - i, 30
    - g, 14
      - f, 5
      - e, 9
    - d, 16

Second tree:
- a, 45
  - j, 55

# Huffman codes

- Min Heap

**j, 55**

**h, 25**  **i, 30**

**c, 12**  **b, 13**  **g, 14**  **d, 16**

**f, 5**  **e, 9**  **a, 45**

**j, 55**

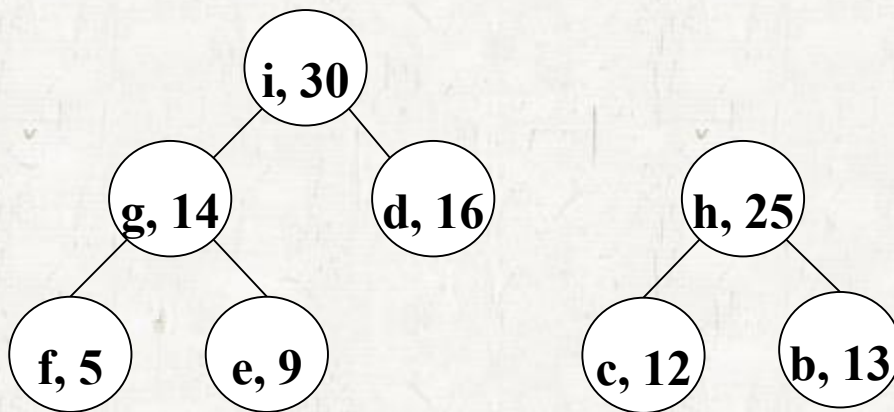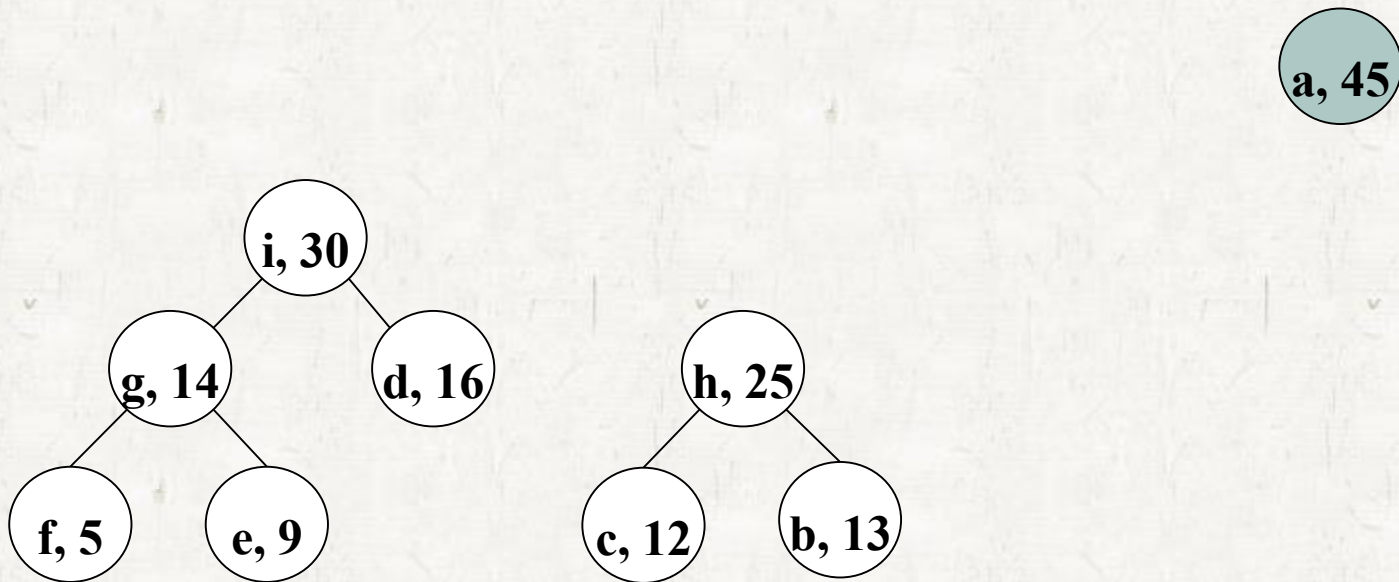# Huffman codes

- Min Heap

# Huffman codes
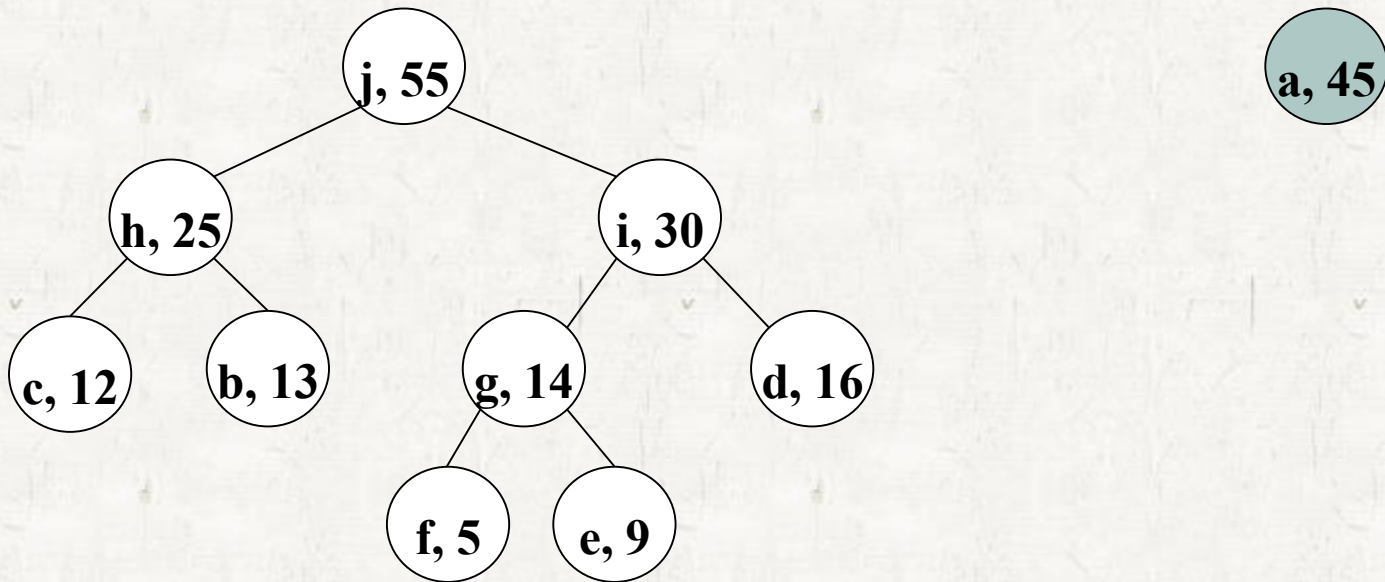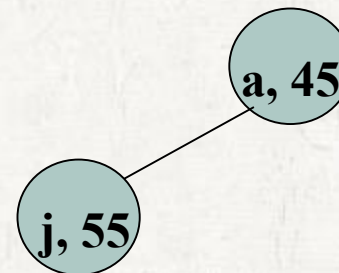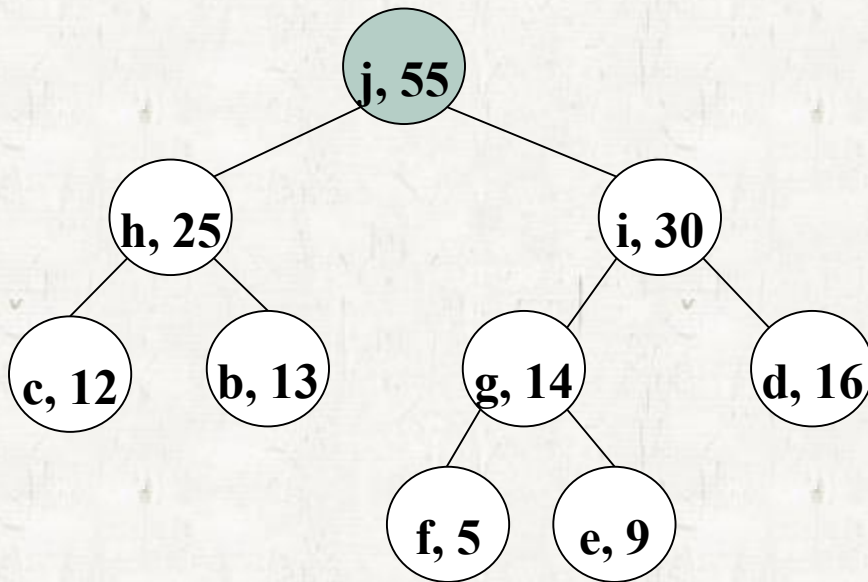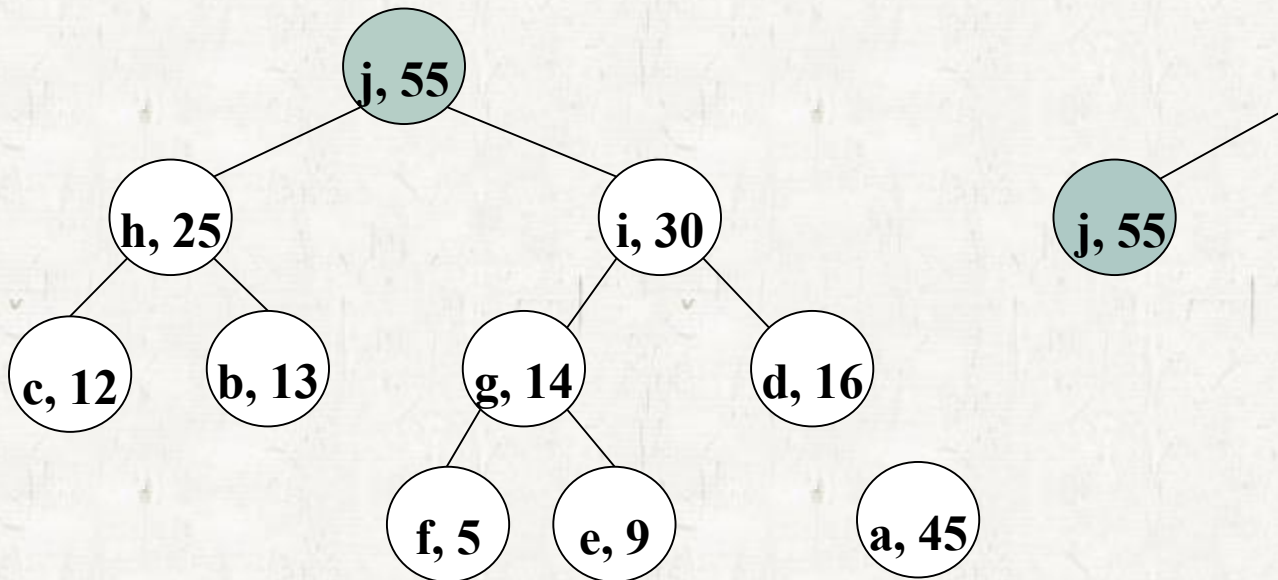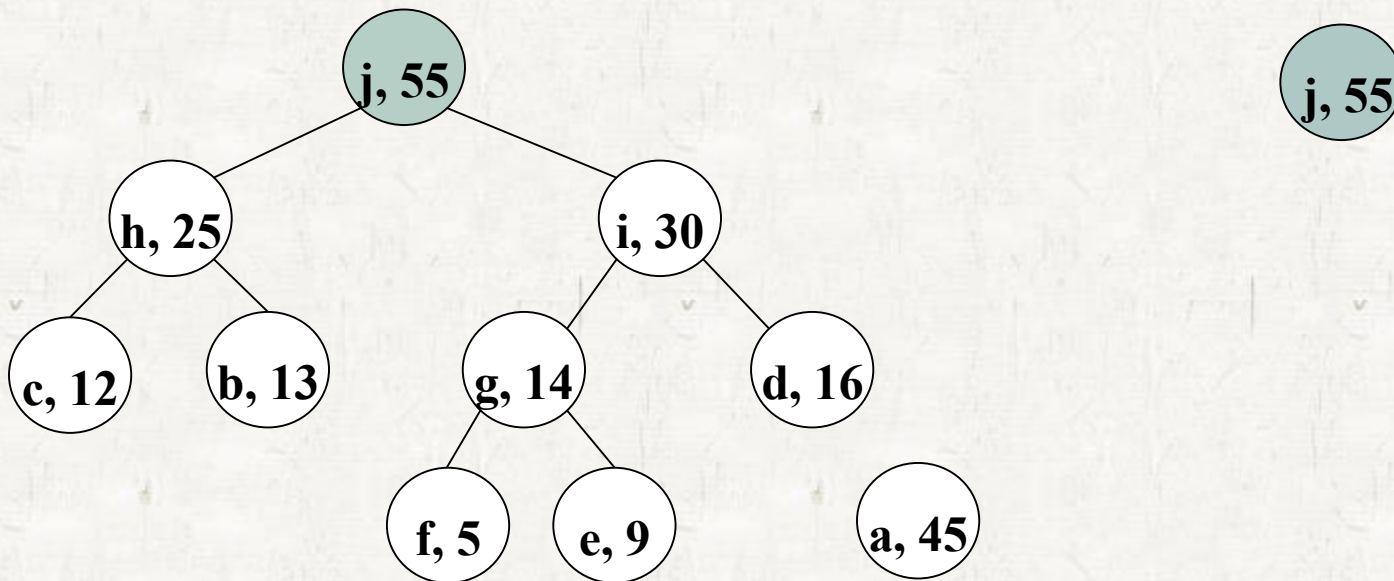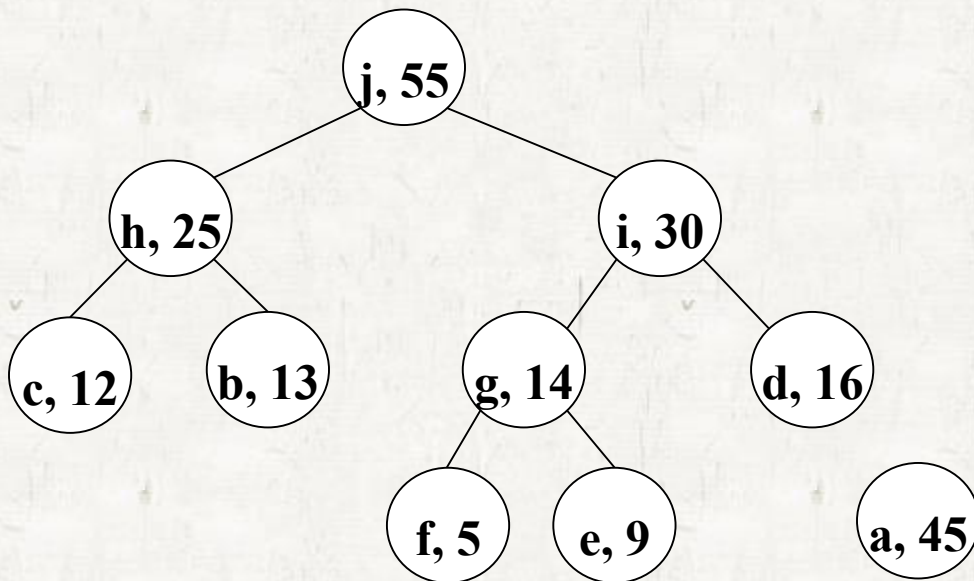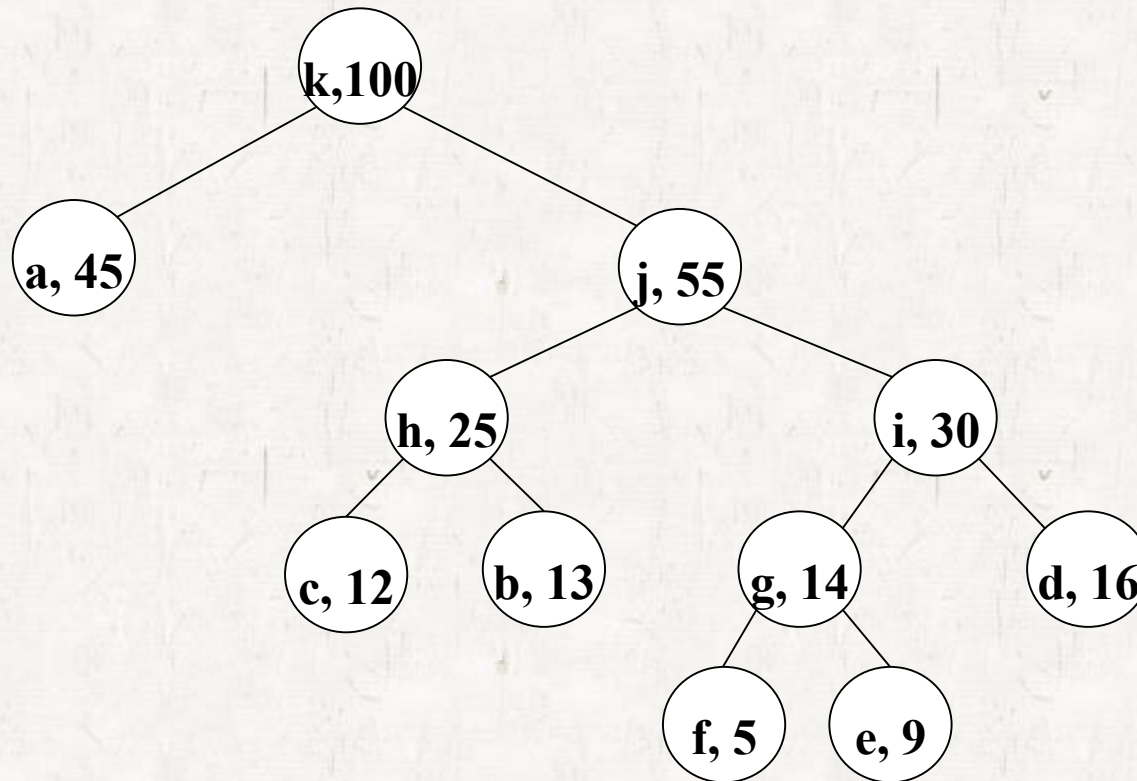
- Min Heap

# Huffman codes

- Min Heap

# Huffman codes

- Min Heap

# Huffman codes

- **Running time**: O($n$lg$n$)
  - Build min heap: $O(n)$
  - Merge: $n$-1 times
    - Each merge requires two minimum selection: $O(\lg n)$
      + one     insertion

# Huffman codes

- **Correctness**
  - *Lemma 16.2*
    - Let $C$ be an alphabet in which each character $c \in C$ has frequency $f[c]$.
    - Let $x$ and $y$ be two characters in $C$ having the lowest frequencies.
    - Then there exists an optimal prefix code for $C$ in which the *codewords for x and y have the same length and differ only in the last bit*.

# Huffman codes

- *Proof*
  - **Idea**: take an arbitrary optimal prefix code tree $T$ and modify it and to <mark>make a tree representing another optimal prefix code such that the characters $x$ and $y$ appear as sibling leaves of maximum depth in the new tree.</mark>



Some length, 미리 보면안 다음

# Huffman codes

- **The cost of tree *T***
  - *f*(*c*): frequency of a character *c*
  - $d_T$(*c*): length of the codeword for *c*

$$B(T) = \sum_{c \in C} f(c)d_T(c)$$

$f_{(x)}, f_{(y)} \leq f_{(a)}, f_{(b)}$

→ 가장 아래 nodes가 작음

$B(T) = \cdots 2f_{(x)} + 1 \cdot f_{(y)} + 3f_{(a)} + 3f_{(b)}$

$B(T'') = \cdots 2f_{(a)} + 1 \cdot f_{(b)} + 3f_{(x)} + 3f_{(y)}$

$2(f_{(x)} + f_{(a)} + f_{(y)} - f_{(b)}) + 3(f_{(a)} - f_{(x)}) + 3(f_{(b)} - f_{(y)})$

$= f_{(a)} - f_{(x)} + 2(f_{(b)} - f_{(y)}) \geq 0$
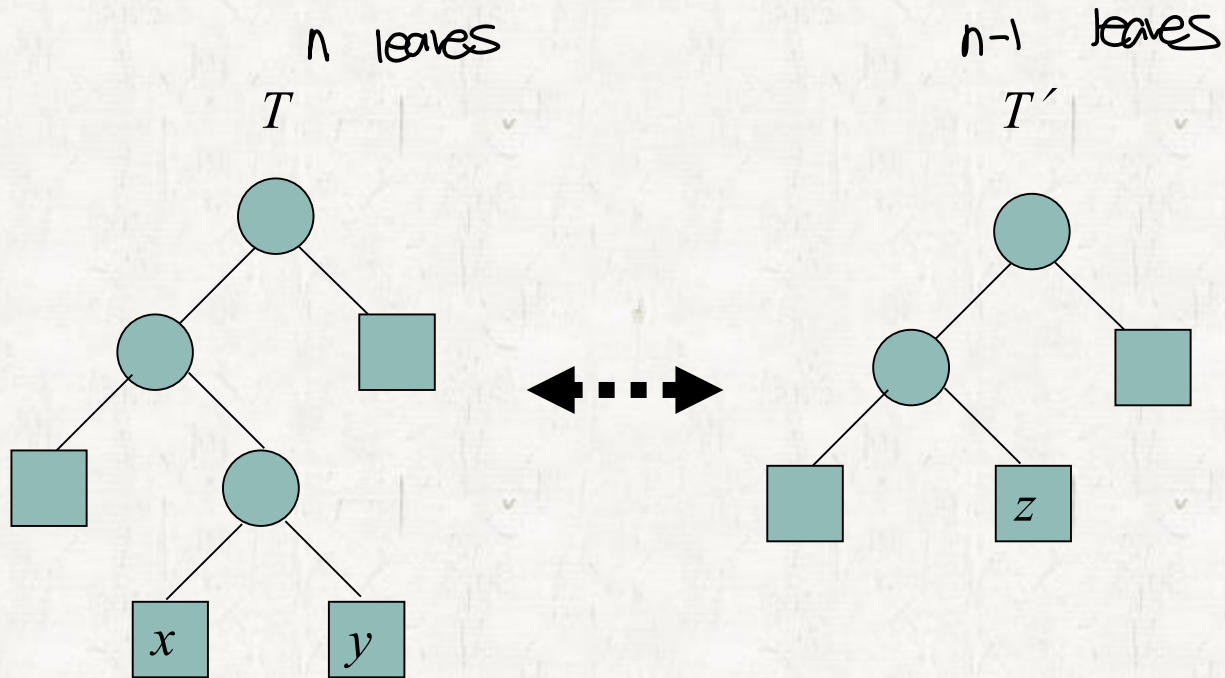
⟹ $B(T) \geq B(T'')$  T가 Optimal → T'도 Optimal tree 이다.

73

- ## *Lemma 16.3*

  - Let $x$ and $y$ be two characters in a given alphabet $C$ with minimum frequency.

  - Let $C'$ be the alphabet $C$ with characters $x, y$ removed and character $z$ added, so that $C' = C - \{x, y\} \cup \{z\}$; define $f$ for $C'$ as for $C$, except that $f[z] = f[x] + f[y]$.  것에는 virtual nodes  X,Y 없음

  - Let $T'$ be any tree representing an optimal prefix code for the alphabet $C'$.

  - Then the optimal prefix code tree $T$ for $C$ can be obtained from $T'$ by replacing the leaf node for $z$ with an internal node having $x$ and $y$ as children.

# Huffman codes



n leaves

n-1 leaves

$T$

$T'$

$x$   $y$

$z$

# Huffman codes

- *Proof*
  - Show $B(T) = B(T') + f[x] + f[y]$
    - For each $c \in C - \{x, y\}$, we have $d_T(c) = d_{T'}(c)$, and hence $f[c]d_T(c) = f[c]d_{T'}(c)$.
    - Since $d_T(x) = d_T(y) = d'(z) + 1$, we have
      $$f[x]d_T(x) + f[y]d_T(y) = (f[x] + f[y])(d_{T'}(z) + 1)$$
      $$= f[z]d_{T'}(z) + (f[x] + f[y])$$
    - From which we conclude that $B(T) = B(T') + f[x] + f[y]$ or, equivalently $B(T') = B(T) - f[x] - f[y]$.

# Huffman codes

- ○ **_Proof_**
  - ● Suppose $T$ does not represent an optimal prefix code for $C$.
  - ● There exists $T''$ such that $B(T'') < B(T)$.
  - ● By Lemma 16.2, there exists $T''$ having $x$ and $y$ as siblings.
  - ● Let $T'''$ be the tree $T''$ with the common parent of $x$ and $y$ replaced by a leaf $z$ with frequency $f[z] = f[x] + f[y]$.
  - ● Then, $B(T''') = B(T'') - f[x] - f[y]$
    $$< B(T) - f[x] - f[y]$$
    $$= B(T')$$
    - ➔ Contradiction
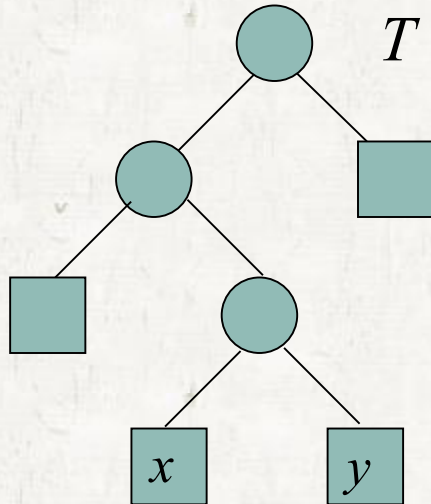  - ● $T$ must represent an optimal prefix code for the alphabet $C$.

# Huffman codes

$$B(T) = B(T') - d_T(z)f(z) + d_T(x)f(x) + d_T(y)f(y)$$

$$f(z) = f(x) + f(y)$$
$$d_T(z) = d_T(x) - 1$$
$$d_T(y) = d_T(x)$$

$$\Rightarrow B(T) = B(T') + f(x) + f(y)$$

$T$    $T'$

$T''$ n leaves

$T'''$

$T'$ optimal $\rightarrow$ T is optimal
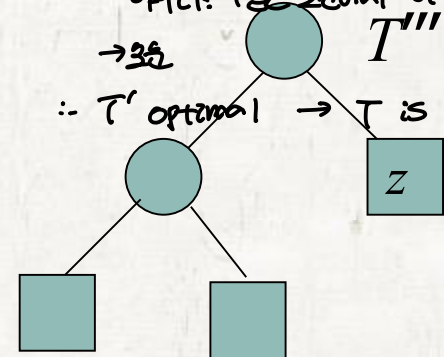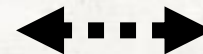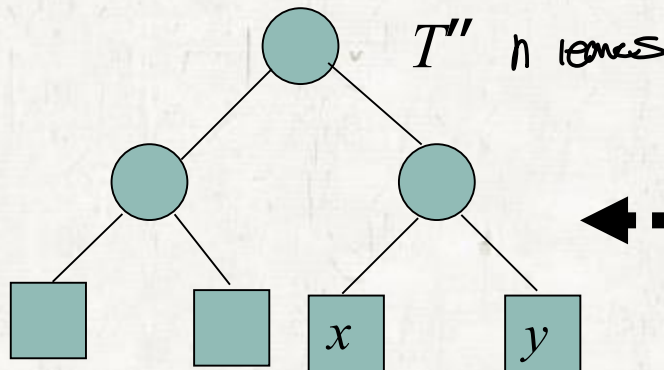모순형 )  가정: $T'$ $\rightarrow$ T is not optimal
$B(T) > B(T'')$ 라고 하자
$B(T'') + B(T''') + f(x) + f(y)$
$B(T') > B(T''') \rightarrow B(T')$ 는 optimal이
아니다. (좀 크다에서 더 나은 Tree존재)
$\rightarrow$ 모순
∴ $T'$ optimal $\rightarrow$ T is optimal

78

# Self-study

- **Exercise 16.3-3 (16.3-2 in the 2$^{nd}$ ed.)**

  - Fibonacci number definition is in p. 59 (p. 56 in the 2$^{nd}$ ed.)

- **Exercise 16.3-7 (16.3-6 in the 2$^{nd}$ ed.)**