

제 8장 근사 알고리즘

Goal

- Understand how to design approximation algorithms
- Understand how to evaluate approximation ratio
- Understand well-known approximation algorithms

Index

- Traveling Salesman Problem (TSP)
- Vertex Covering
- Bin Packing
- Task Scheduling
- Clustering

Strategy of Problem Solving

- For a given problem,
 - If it is easy
 - Design an algorithm solving the problem in polynomial time
 - If it is not working in time or proven as a NP-complete problem
 - Invest money for buying a supercomputer (x, not your option)
 - Build an algorithm working in polynomial time by giving up accuracy

Strategy of Problem Solving

- How to control the trade-off between accuracy and time?
 - Do not search all possible solutions
 - Prune search space by hypothesis
 - Search until you find a relatively good solution
 - Many probabilistic techniques
 - Leads to AI algorithms

Strategy of Problem Solving

- NP-완전 문제들은 실 세계의 광범위한 영역에 활용되지만, 이 문제들을 다항식 시간에 해결할 수 있는 알고리즘이 아직 발견되지 않았다.
- 또한 아직까지 그 누구도 이 문제들을 다항식 시간에 해결할 수 없다고 증명하지 못했다.
- 대부분의 학자들은 이 문제들을 해결할 다항식 시간 알고리즘이 존재하지 않을 것이라고 추측하고 있다.
- c: 여기서 다항식 시간은 deterministic polynomial time을 의미

Strategy of Problem Solving

- 이러한 NP-완전 문제들을 어떤 방식으로든지 해결하려면 다음의 3가지 중에서 1가지는 포기해야 한다.
 1. 다항식 시간에 해를 찾는 것
 2. 모든 입력에 대해 해를 찾는 것
 3. 최적해를 찾는 것
 - 8장에서는 NP-완전 문제를 해결하기 위해 3 번째 것을 포기한다. 즉, **최적해에 아주 근사한 (가까운) 해를 찾아주는 근사 알고리즘** (Approximation algorithm)을 살펴본다.
- c: 2. 와 3. 은 같은 의미. 2.가 안되면 3.이 보장이 안됨. 3.과 관련 없이 2.를 할 수 있다는 것은 이미 우리가 확실하게 적용할 수 있는 휴리스틱이 있다는 이야기.

Approximation Ratio

- 근사 알고리즘은 근사해를 찾는 대신에 다항식 시간의 복잡도를 가진다.
- 근사 알고리즘은 근사해가 얼마나 최적해에 근사한지 (즉, 최적해에 얼마나 가까운지)를 나타내는 근사 비율 (Approximation Ratio)을 알고리즘과 함께 제시하여야 한다.
- 근사 비율은 근사해의 값 (c')과 최적해의 값 (c)의 비율 (c'/c)로서, 1.0에 가까울수록 정확도가 높은 알고리즘이다.
- 근사 비율을 계산하려면 최적해를 알아야 하는 모순이 생긴다.
- 따라서 최적해를 대신할 수 있는 '간접적인' 최적해를 찾고, 이를 최적해로 삼아서 근사 비율을 계산한다.

c: 최적해가 무엇인지 모르지만 최적해의 값을 아는 문제들도 많음. 이 경우 모순이 아니고, 최적해의 비율만 계산. 예를 들면 regression 문제들.

8.1 여행자 문제

- 여행자 문제 (Traveling Salesman Problem, TSP)
 - 여행자가 임의의 한 도시에서 출발
 - 다른 모든 도시를 1번씩만 방문
 - 출발했던 도시로 도착
 - 돌아오는 여행 경로의 길이를 최소화

8.1 여행자 문제

- 제약조건에 따른 TSP의 다양한 변형 문제 존재
- 이 강의에서 다루는 문제의 제약 조건
 1. 도시 A에서 도시 B로 가는 거리는 도시 B에서 도시 A로 가는 거리와 같다. (대칭성, symmetry)
 2. 도시 A에서 도시 B로 가는 거리는 도시 A에서 다른 도시 C를 경유하여 도시 B로 가는 거리보다 짧다. (삼각 부등식 특성, triangle inequality)

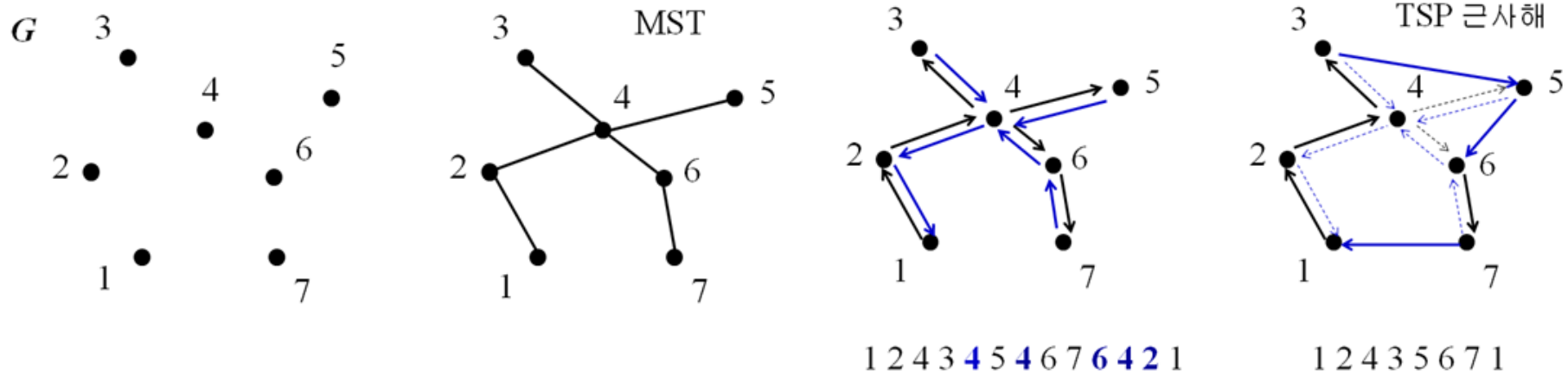
C: weight는 각 도시간의 distance를 의미 (distance의 정의 만족)

8.1 여행자 문제

- TSP는 NP-complete 문제로 알려짐
- 근사 알고리즘 고안 방식
 - 1. 유사한 특성을 가진 문제로부터의 변환
 - 2. Minimum Spanning Tree (한 번씩 노드 방문, weight sum 최소화)
 - 3. TSP 문제로 변환하기 위하여, MST를 트리가 아닌 시작점을 중심으로 한 cycle을 찾도록 문제 변환 방법 개발
 - e.g.) 시작 도시를 제외한 다른 모든 도시를 트리 선분을 따라 1번씩 방문하도록 경로 설정 후 TSP 조건에 맞게 변형.

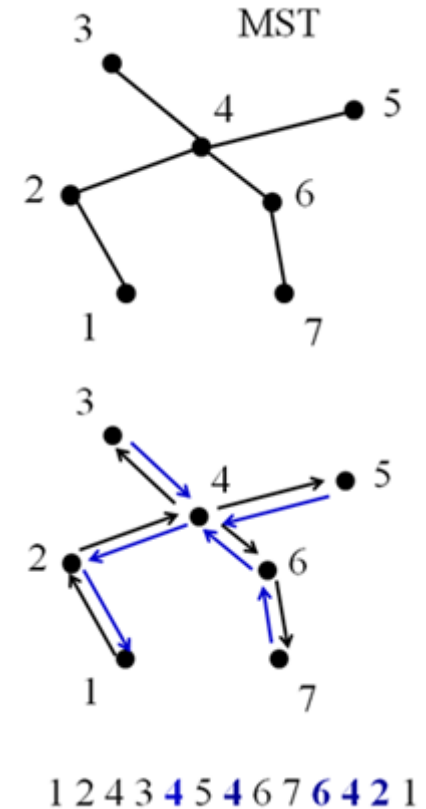
8.1 여행자 문제

MST 최적해의 TSP 근사해로 변형 과정의 예

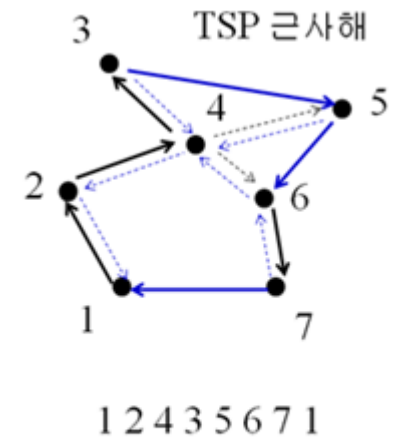
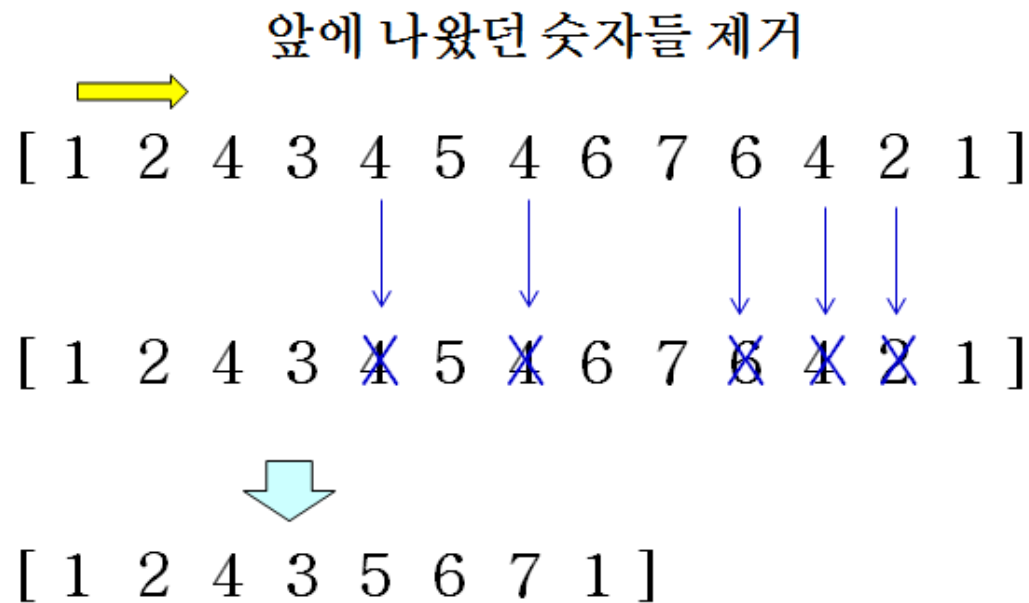


8.1 여행자 문제

- MST 최적해 탐색
- 임의의 도시를 출발점으로 선정
 - (그림에서는 도시 1)
 - TSP에서 출발점 주어지는 경우 그대로 사용
- 트리의 선분을 따라 모든 도시를 방문하도록 path 생성
- [1 2 4 3 4 5 4 6 7 6 4 2 1].
- 중복되는 도시들의 제거 (출발점 제외)



8.1 여행자 문제



- 제거 시 triangle inequality 적용 (가능하면 더 짧은 path 생성)

8.1 여행자 문제

Pseudo Code

입력: a weighted graph (weights are distance values)

출력: a simple path starting from and ending to a randomly selected node

1. Find a MST
2. Convert a MST to a TSP path
3. Remove repeatedly visited nodes in the path
4. return the simple path

8.1 여행자 문제

Explanation of the Pseudo Code

- Line 1에서는 4.2절의 크루스컬 (Kruskal)이나 프림 (Prim) 알고리즘을 사용하여 MST를 찾는다.
- Line 2에서는 하나의 도시에서 어떤 트리 선분을 선택하여 다른 도시를 방문할 때 지켜야 할 순서가 정해져 있지 않으므로, 임의의 순서로 방문해도 괜찮다.
- Line 3에서는 삼각 부등식의 원리를 적용함과 동시에 각 도시를 1번씩만 방문하기 위해서 중복 방문된 도시를 제거한다. 단, 가장 마지막 도시인 출발 도시는 중복되지만 TSP 문제의 출발 도시로 돌아와야 한다는 조건에 따라 제거하지 않는다.

8.1 여행자 문제

시간복잡도

- Line 1: MST를 찾는 데에는 크루스컬이나 프림 알고리즘의 시간 복잡도만큼 시간이 걸리고,
- Line 2에서 트리 선분을 따라서 도시 방문 순서를 찾는 데는 $O(n)$ 시간이 걸린다. 왜냐하면 트리의 선분 수가 $(n-1)$ 이기 때문이다.
- Line 3: line 2에서 찾은 도시 방문 순서를 따라가며, 단순히 중복된 도시를 제거하므로 $O(n)$ 시간이 걸린다.
- 시간복잡도는 (크루스컬이나 프림 알고리즘의 시간복잡도) + $O(n)$ + $O(n)$ 이므로 **크루스컬이나 프림 알고리즘의 시간복잡도와 같다.**

8.1 여행자 문제

근사 비율

- 근사 비율의 upperbound 제시
 - 에러가 최적해 대비 xxx% 를 넘어서지 않는다는 주장을 하기 위한 분석
 - 최적해를 모르면 최적해보다 더 낮은 가능한 값에 대해 근사 비율 설정
 - Upperbound는 정확한 최적해까지 알지 못해도 유도 가능

8.1 여행자 문제

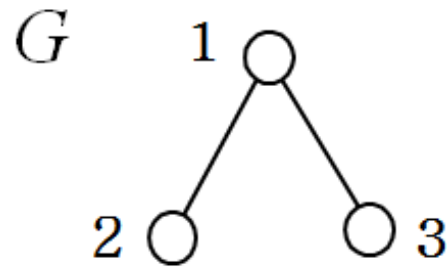
근사 비율

- 가능한 최적해 값
 - MST의 가중치 합 M
 - MST 가중치 $M < \text{TSP 최적해 } M'$ (why?)
- Approx_MST_TSP $M'' < 2M$
 - A raw TSP path uses each edge of MST two times
 - Triangle inequality makes the sum less than $2M$
- approximation ratio = $M'' / M' < 2M / M' < 2M / M = 2$

8.2 정점 커버 문제

문제 정의

- 정점 커버 (Vertex Cover) 문제는 주어진 그래프 $G=(V,E)$ 에서 각 선분의 양 끝점들 중에서 적어도 하나의 끝점을 포함하는 점들의 집합들 중에서 최소 크기의 집합을 찾는 문제이다.
- 정점 커버를 살펴보면, 그래프의 모든 선분이 정점 커버에 속한 점에 인접해있다.
- 즉, 정점 커버에 속한 점으로서 **그래프의 모든 선분을 '커버'하는 것이다.**

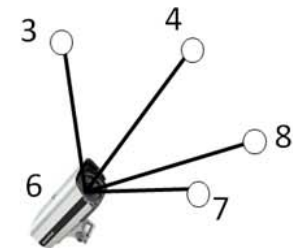
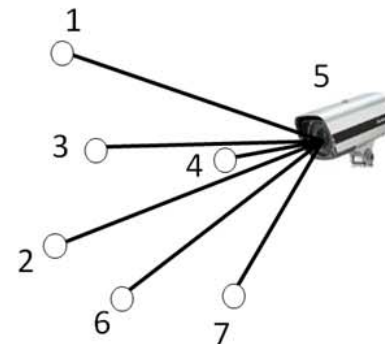
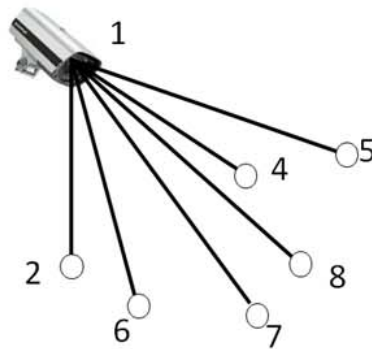
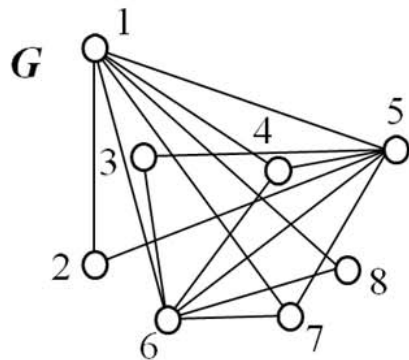


정점 커버 문제의 예

- 위의 그래프 G 에서 $\{1, 2, 3\}$, $\{1, 2\}$, $\{1, 3\}$, $\{2, 3\}$, $\{1\}$ 이 각각 정점 커버이다.
- $\{2\}$ 또는 $\{3\}$ 은 정점 커버가 아니다. 왜냐하면 $\{2\}$ 는 선분 $(1, 3)$ 을 커버하지 못하고, $\{3\}$ 은 선분 $(1, 2)$ 를 커버하지 못한다.
- 따라서 위의 그래프에 대한 정점 커버 문제의 해는 $\{1\}$ 이다.

정점 커버 문제 예 2

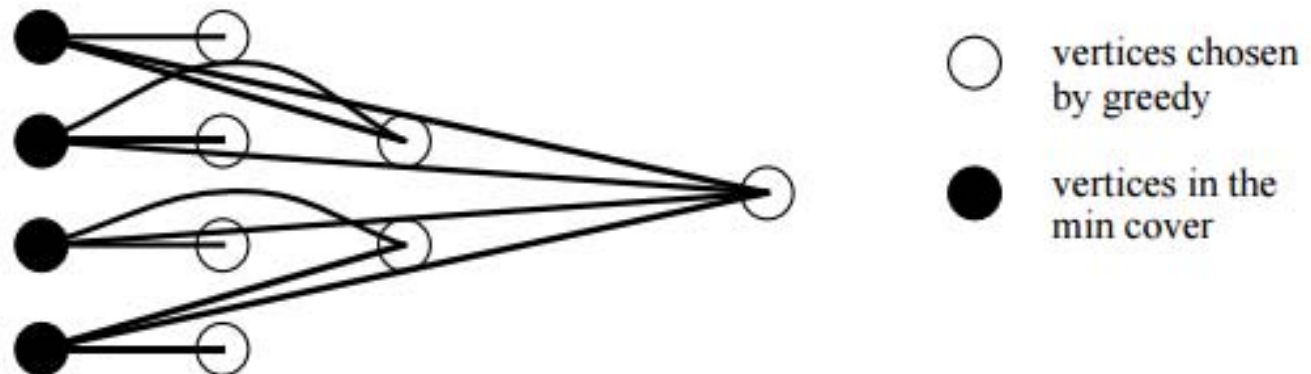
- 아래의 그래프 G 는 어느 건물의 내부도면을 나타낸다.
- 건물의 모든 복도를 감시하기 위해 가장 적은 수의 CCTV 카메라를 설치하고자 한다.
- 이를 위해서 3대의 카메라를 각각 점 1, 5, 6에 설치하면 모든 복도 (선분)을 '커버'할 수 있다.



정점 커버 문제를 위한 Greedy Algorithm의 Approximation ratio

- 그래프 G 의 정점 커버를 찾는 것은 건물 내부 경비를 위한 최소의 카메라 수와 각 카메라의 위치를 구하는 것과 동일
- How to select such vertices?
- Select the most frequently used vertex in all edges.
- Approximation ratio?

– $\log n$

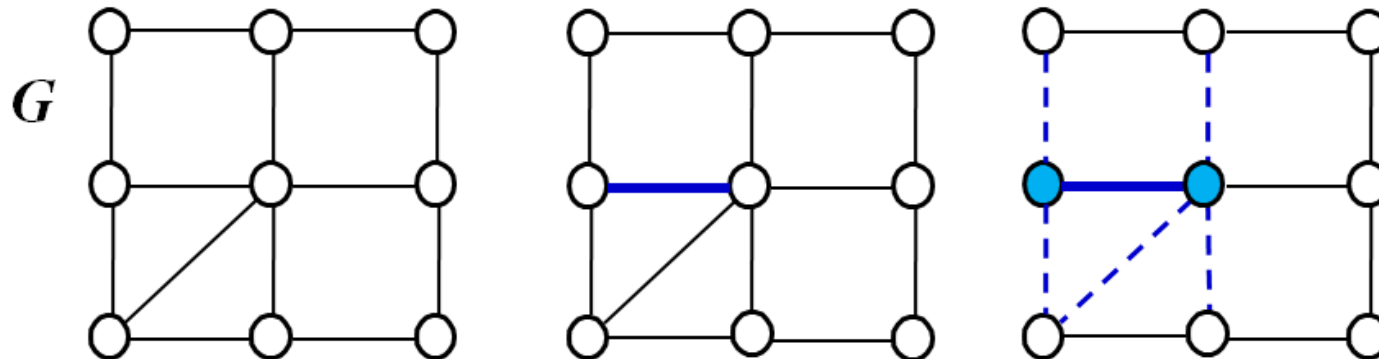


Another Algorithm

- Selecting edges instead of vertices
- Selecting an edge \rightarrow selecting two vertices
- E : selected edges
- V : vertices of e in E
- Select a disjoint edge (a set of vertices) with selected V' at each iteration

Matching의 예

- 그림에서 1개의 선분이 임의로 선택되었을 때, 선택된 선분 주변의 5개의 선분 (점선으로 표기된 선분)은 정점 커버를 위해 선택되지 않는다.
- 그 이유는 선택된 선분의 양 끝점 (파란색 점)들이 점선으로 표시된 선분을 커버하기 때문이다.



- 이러한 방식으로 선분을 선택하다가 더 이상 선분을 추가 수 없을 때 중단한다.
- 이렇게 선택된 선분의 집합을 **극대 매칭(maximal matching)**이라고 한다.
- **매칭(matching)**이란 각 선분의 양쪽 끝점들이 중복되지 않는 선분의 집합이다.
- 극대 매칭은 이미 선택된 선분에 기반을 두고 새로운 선분을 추가하려 해도 더 이상 추가할 수 없는 매칭을 말한다.

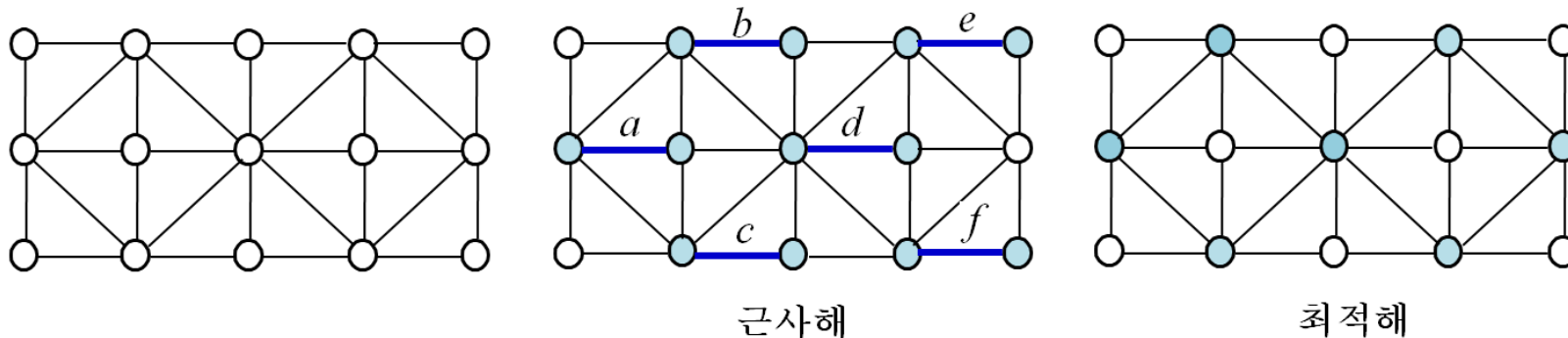
Approx_Matching_VC

입력: 그래프 $G=(V,E)$

출력: 정점 커버

1. 입력 그래프에서 극대 매칭 M 을 찾는다.
2. return 매칭 M 의 선분의 양 끝점들의 집합

- 아래의 그림에서 극대 매칭으로서 선분 a, b, c, d, e, f 가 선택되었다. 따라서 근사해는 선분 a, b, c, d, e, f 의 양 끝점들의 집합이다.
- 즉, 근사해는 총 12개의 점으로 구성된다. 반면에 오른쪽 그림은 입력 그래프의 최적해로서 7개의 점으로 구성되어 있다.



시간복잡도

- Approx_Matching_VC 알고리즘의 시간복잡도는 주어진 그래프에서 극대 매칭을 찾는 과정의 시간복잡도와 같다.
- 극대 매칭을 찾기 위해 하나의 선분을 선택할 때, 양 끝점들이 이미 선택된 선분의 양 끝점과 동일한지를 검사해야 하므로, 이는 $O(n)$ 시간이 걸린다.
- 그런데 입력 그래프의 선분 수가 m 이면, 각 선분에 대해서 $O(n)$ 시간이 걸리므로, Approx_Matching_VC 알고리즘의 시간복잡도는 $O(n) \times m = O(nm)$ 이다.

근사 비율

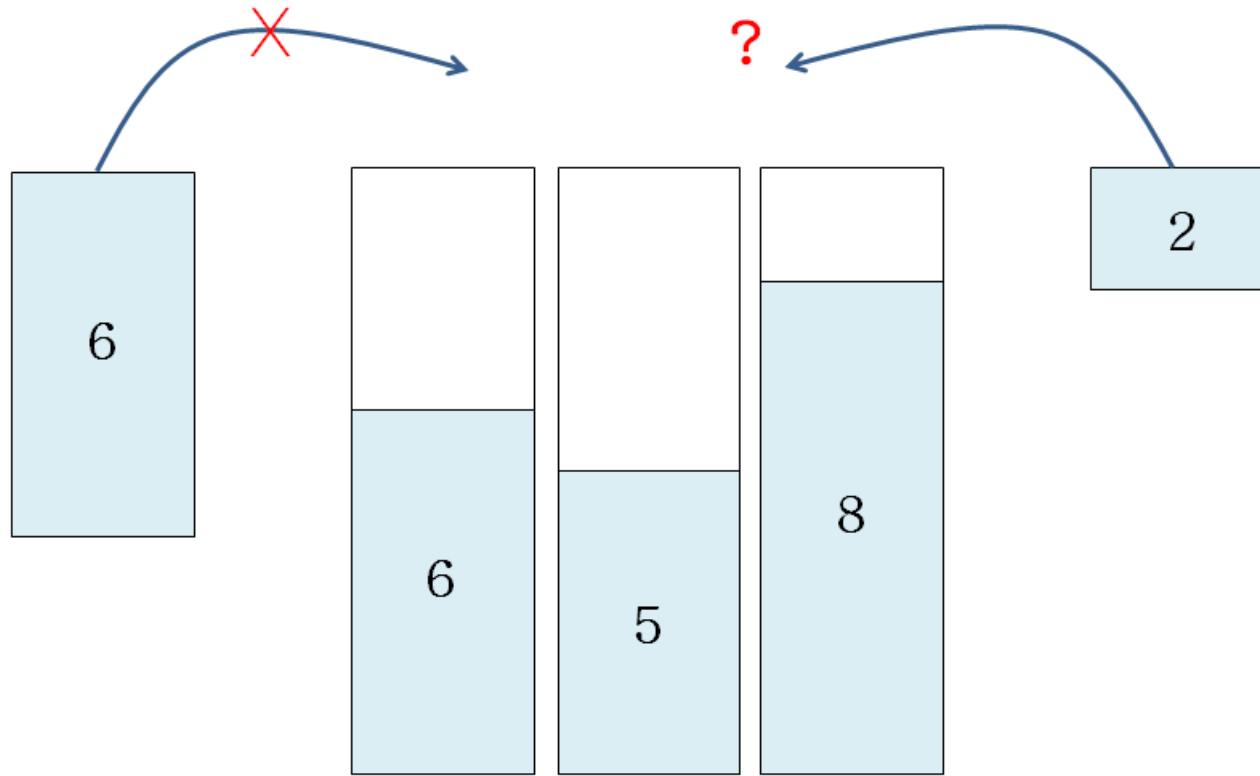
- M : the number of edges in the maximal matching S
- M' : the number of vertices of true optimal T
- $M \leq M'$ (why?)
 - Vertices of edges in S are all unique
 - Vertices in T should be included at least one time in all edges of S

근사 비율

- $M' \leq 2M$
 - Vertices of all edges of $S \rightarrow$ a vertex cover (not optimal)
 - Because M' is the optimal M' should be less than $2M$ (by definition)
- $M \leq M' \leq 2M$
- Our estimation is $M'' = 2M$!!
- Approximation ratio $R : 1/2 \leq M/M' \leq 1$,
- $1 \leq M''/M' \leq 2$

8.3 통 채우기 문제

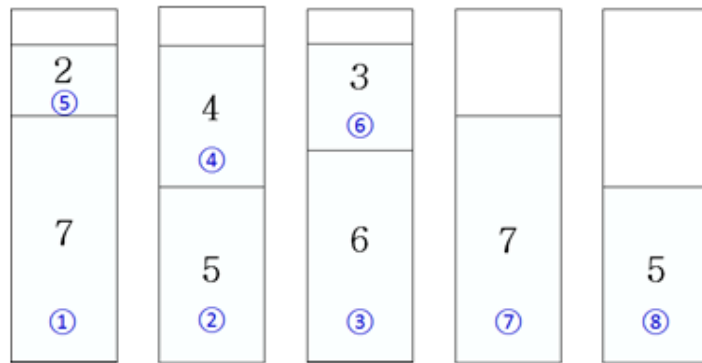
- 통 채우기(Bin Packing) 문제는 n 개의 물건이 주어지고, 통(bin)의 용량이 c 일 때, 모든 물건을 가장 적은 수의 통에 채우는 문제이다. 단, 각 물건의 크기는 c 보다 크지 않다.
- 물건을 통에 넣으려고 할 때에는 그 통에 물건이 들어갈 여유가 있어야만 한다.
- 예를 들어, 통의 크기가 10이고, 현재 3개의 통에 각각 6, 5, 8 만큼씩 차있다면, 크기가 6인 물건은 기존의 3개의 통에 넣을 수 없고 새 통에 넣어야 한다. 그런데 만일 새 물건의 크기가 2라면, 어느 통에 새 물건을 넣어야 할까?



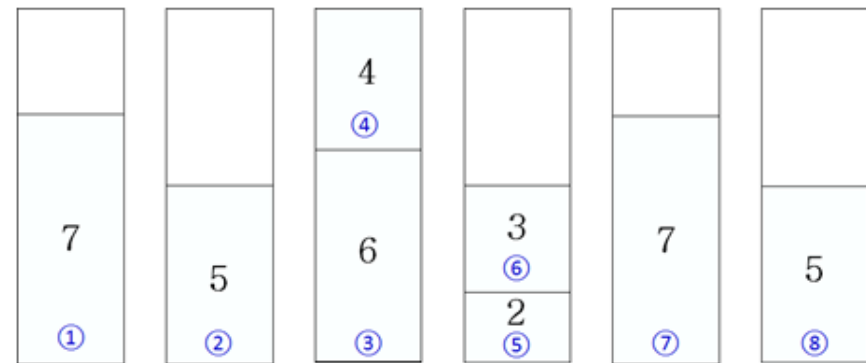
- 위의 질문에 대한 간단한 답은 **그리디 방법**으로 넣을 통을 정하는 것이다.

- 그리디 방법은 '무엇에 욕심을 낼 것인가'에 따라서 아래와 같이 4 종류로 분류할 수 있다.
1. **최초 적합 (first fit)**: 첫 번째 통부터 차례로 살펴보며, 가장 먼저 여유가 있는 통에 새 물건을 넣는다.
 2. **다음 적합 (next fit)**: 직전에 물건을 넣은 통에 여유가 있으면 새 물건을 넣는다.
 3. **최선 적합 (best fit)**: 기존의 통 중에서 새 물건이 들어가면 남는 부분이 가장 작은 통에 새 물건을 넣는다.
 4. **최악 적합 (worst fit)**: 기존의 통 중에서 새 물건이 들어가면 남는 부분이 가장 큰 통에 새 물건을 넣는다.

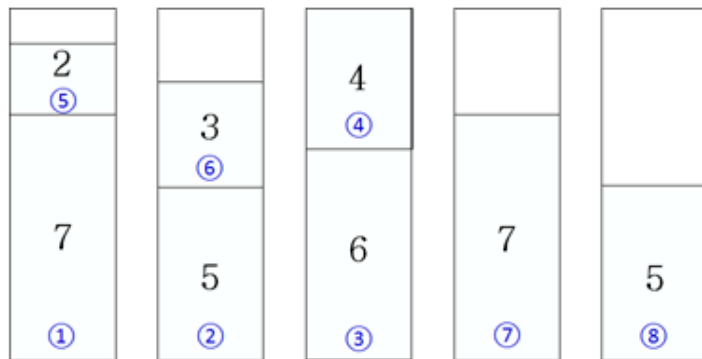
- 통의 용량 $c=10$ 이고, 물건의 크기가 각각 **[7, 5, 6, 4, 2, 3, 7, 5]**일 때, 최초 적합, 다음 적합, 최선 적합, 최악 적합을 각각 적용한 결과는 다음과 같다.



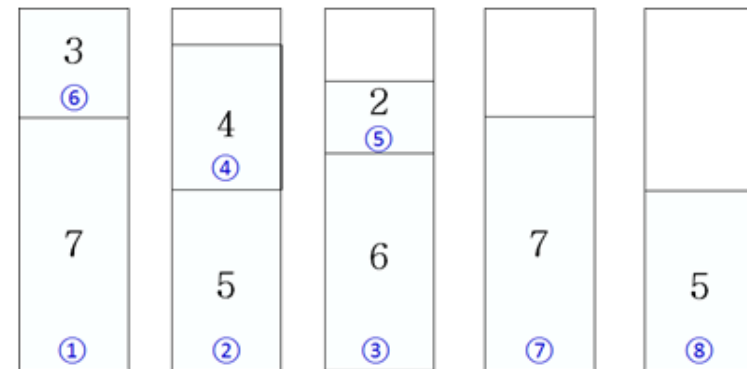
처음 적합



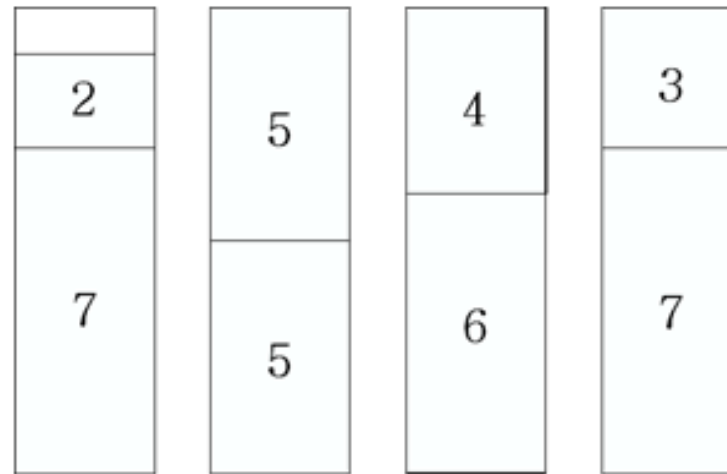
다음 적합



최선 적합



최악 적합



최적해

- 최후 적합 (Last Fit)
- 감소순 최초 적합 (First Fit Decrease)
- 감소순 최선 적합 (Best Fit Decrease)

Approx_BinPacking

입력: n 개의 물건의 각각의 크기

출력: 모든 물건을 넣는데 사용된 통의 수

1. $B = 0$ // 사용된 통의 수
2. for $i = 1$ to n
3. if (물건 i 를 넣을 여유가 있는 기존의 통이 있으면)
4. 그리디 방법에 따라 정해진 통에 물건 i 를 넣는다.
5. else
6. 새 통에 물건 i 를 넣는다.
7. $B = B + 1$ // 통의 수를 1 증가시킨다.
8. return B

시간복잡도

- 다음 적합을 제외한 다른 3가지 방법은 새 물건을 넣을 때마다 기존의 통을 살펴보아야 한다.
- 따라서 통의 수가 n 을 넘지 않으므로 수행시간은 $O(n^2)$ 이다.
- 다음 적합은 새 물건에 대해 직전에 사용된 통만을 살펴보면 되므로 수행시간은 $O(n)$ 이다.

근사 비율

- Lowerbound of the used portion of all bins in the optimal solution
- In an ideal case, the least number of used bin is
 - $L = \text{Total volume of products} / \text{bin capacity}$
 - $L < M'$
- M : the number of used bins in a greedy algorithm searching all used bins to put a new product

근사 비율

(continued)

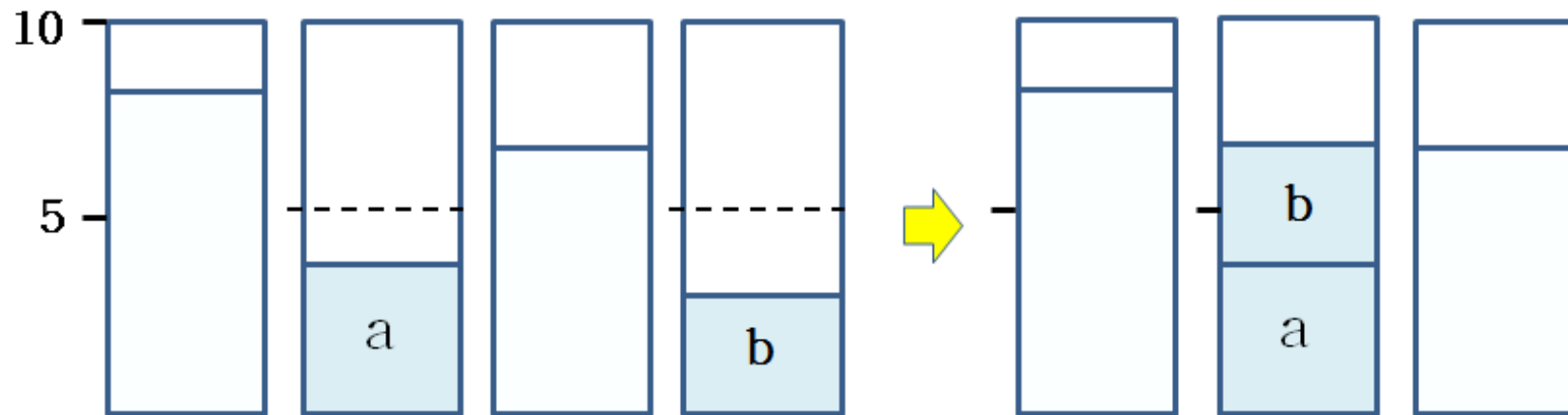
- Candidate solutions of the algorithms have only one or zero bin whose capacity is used less than half.
 - If there are two bins using less than half capacity in a solution, the bins can be merged.
 - -> the solution is not a solution generated from the algorithms.
 - Zero bin case : mean of used portion of all bins $> C/2$ (modified)
 - One bin case : mean of used portion of (all bins – 1) $> C/2$

근사 비율

(continued)

- $\text{bin capacity} / 2 < \text{total volume} / M < \text{total volume} / (M-1)$
- $(M-1) / 2 < \text{total volume} / \text{bin capacity} = L < M'$
- $(M-1) < 2M'$
- $M < 2M' + 1$
- $M \leq 2M'$
- $M/M' \leq 2$

$a, b \leq \text{통의 용량의 } \frac{1}{2}$



- 최적해에서 사용된 통의 수를 OPT 라고 하면,
 $OPT \geq (\text{모든 물건의 크기의 합})/c$ 이다.
- 단, c 는 통의 크기이다.

- 따라서 각 방법이 사용한 통의 수가 OPT' 라면, OPT' 의 통 중에서 1개의 통이 $1/2$ 이하로 차있으므로, 각 방법이 $(OPT'-1)$ 개의 통에 각각 $1/2$ 넘게 물건을 채울 때 그 물건의 크기의 합은 $((OPT'-1) \times C/2)$ 보다는 크다.
- 그러므로 다음과 같은 부등식이 성립한다.

$$(\text{모든 물건의 크기의 합}) > (OPT'-1) \times C/2$$

$$\Rightarrow (\text{모든 물건의 크기의 합})/C > (OPT'-1)/2$$

$$\Rightarrow OPT > (OPT'-1)/2, \quad \text{OPT} \geq (\text{모든 물건의 크기의 합})/C \text{이므로}$$

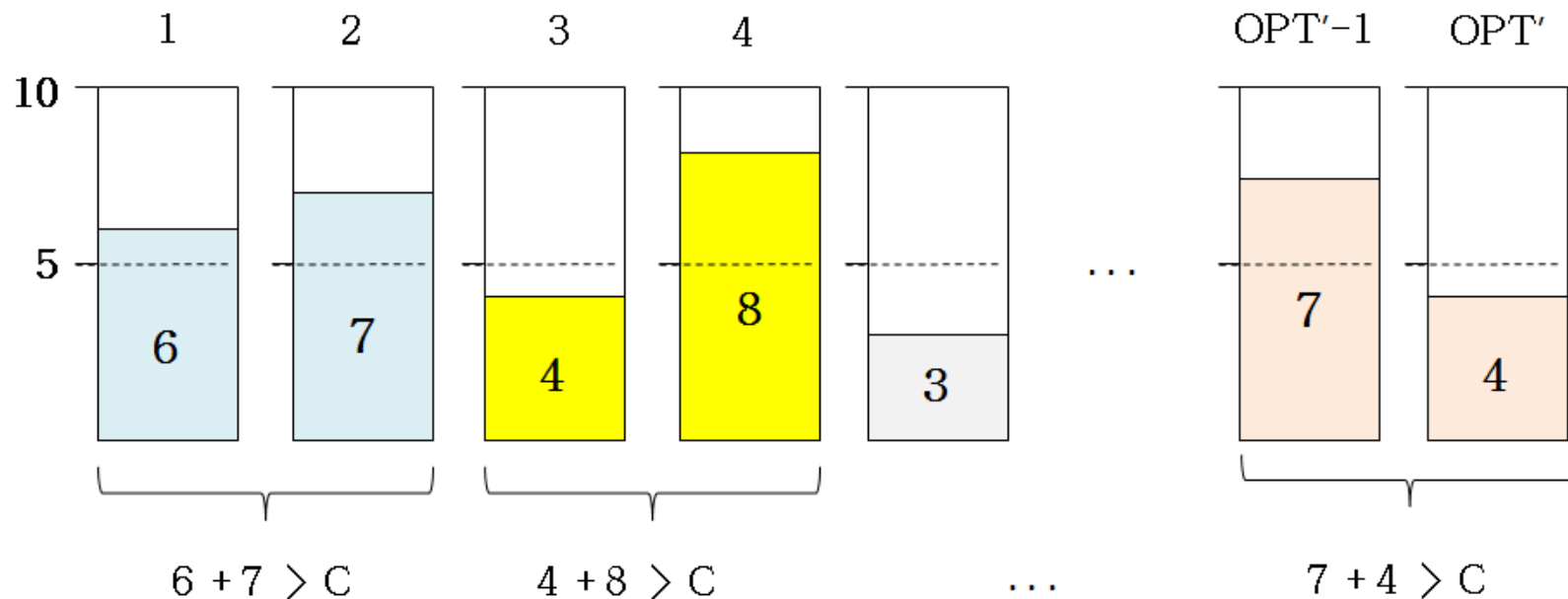
$$\Rightarrow 2OPT > OPT'-1$$

$$\Rightarrow 2OPT + 1 > OPT'$$

$$\Rightarrow 2OPT \geq OPT'$$

- 따라서 3가지 방법의 근사 비율은 2이다.

- 다음 적합은 직전에 사용된 통에 들어있는 물건의 크기의 합과 새 물건의 크기의 합이 통의 용량보다 클 때에만, 새 통에 새 물건을 넣는다.
- 다음 적합이 사용한 통의 수가 OPT' 라면, 이웃한 2개의 통을 다음 그림과 같이 나타낼 수 있다.



$$(\text{모든 물건의 크기의 합}) > \text{OPT}'/2 \times C$$

$$\Rightarrow (\text{모든 물건의 크기의 합})/C > \text{OPT}'/2$$

$$\Rightarrow \text{OPT} > \text{OPT}'/2, \text{ OPT} \geq (\text{모든 물건의 크기의 합})/C \text{ 이므로}$$

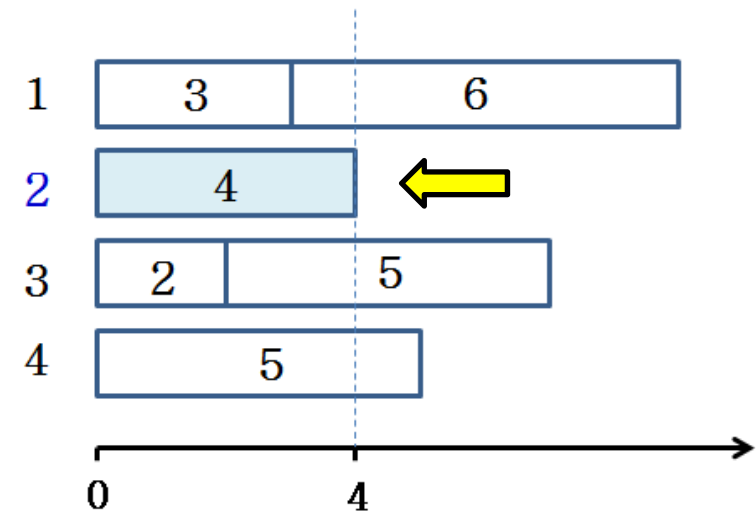
$$\Rightarrow 2\text{OPT} > \text{OPT}'$$

- 따라서 다음 적합의 근사 비율은 2이다.

8.4 작업 스케줄링 문제

- 작업 스케줄링 (Job Scheduling) 문제는 n 개의 작업, 각 작업의 수행 시간 $t_i, i = 1, 2, 3, \dots, n$, 그리고 m 개의 동일한 기계가 주어질 때, 모든 작업이 가장 빨리 종료되도록 작업을 기계에 배정하는 문제이다.
- 단, 한 작업은 배정된 기계에서 연속적으로 수행되어야 한다.
- 또한 기계는 1번에 하나의 작업만을 수행한다.

- 작업을 어느 기계에 배정하여야 모든 작업이 가장 빨리 종료될까?
- 이에 대한 간단한 답은 그리디 방법으로 작업을 배정하는 것이다.
- 즉, 현재까지 배정된 작업에 대해서 가장 빨리 끝나는 기계에 새 작업을 배정하는 것이다.
- 옆의 예제에서는 2 번째 기계가 가장 빨리 작업을 마치므로, 새 작업을 2 번째 기계에 배정한다.



Approx_JobScheduling

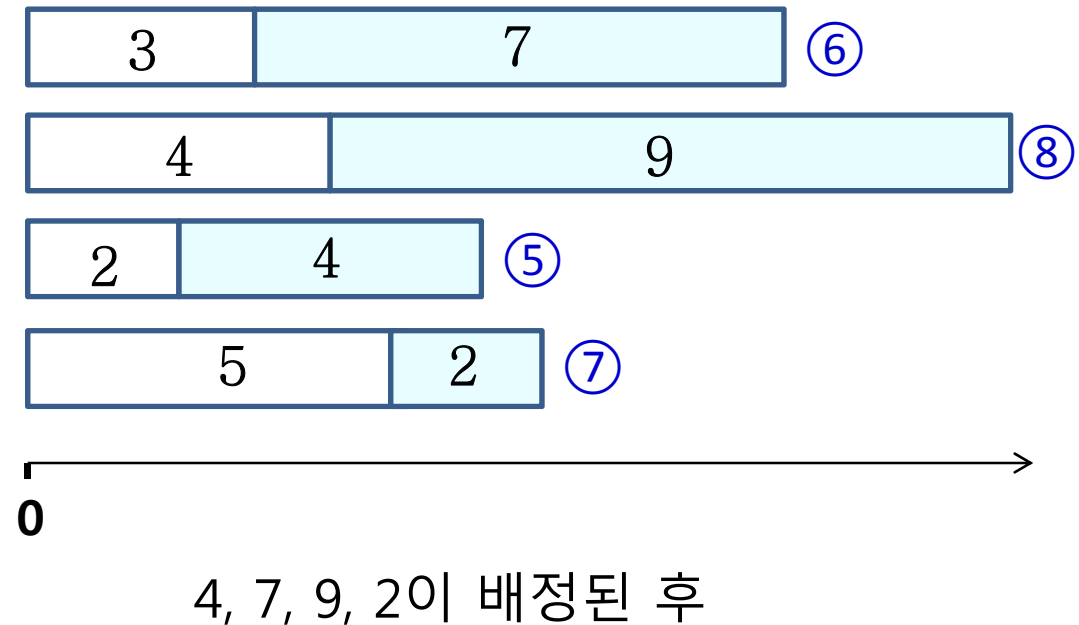
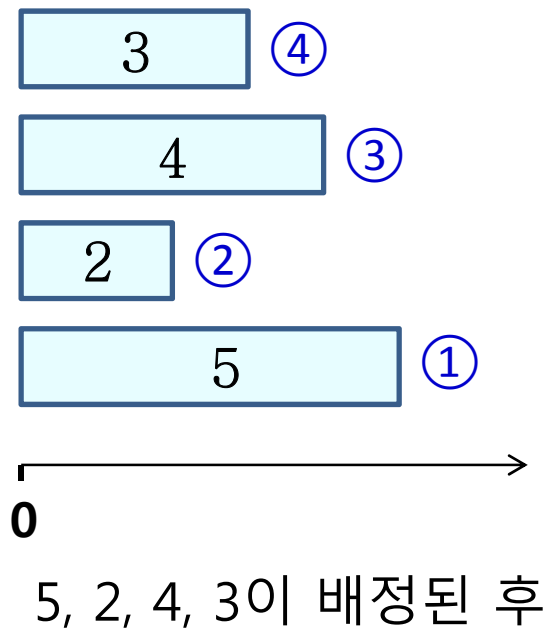
입력: n 개의 작업, 각 작업 수행 시간 $t_i, i = 1, 2, \dots, n$, 기계 $M_j, j = 1, 2, \dots, m$

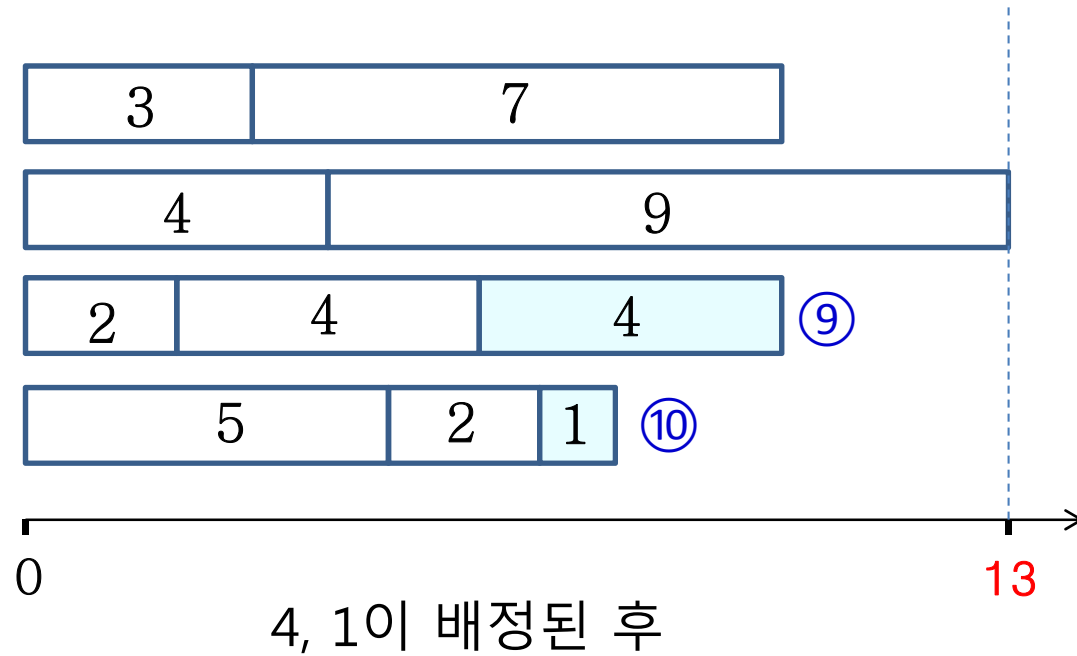
출력: 모든 작업이 종료된 시간

1. for $j = 1$ to m
2. $L[j] = 0$ // $L[j]$ =기계 M_j 에 배정된 마지막 작업의 종료 시간
3. for $i = 1$ to n
4. $\text{min} = 1$
5. for $j = 2$ to m { // 가장 일찍 끝나는 기계를 찾는다.
6. if ($L[j] < L[\text{min}]$)
7. $\text{min} = j$ }
8. 작업 i 를 기계 M_{min} 에 배정한다.
9. $L[\text{min}] = L[\text{min}] + t_i$
10. return 가장 늦은 작업 종료 시간

- Line 1~2: 각 기계에 배정된 마지막 작업의 종료 시간 $L[j]$ 를 0으로 초기화시킨다. 왜냐하면 초기엔 어떤 작업도 기계에 배정되지 않은 상태이기 때문이다. 단, 기계 번호는 1부터 시작된다.
- Line 3~9의 for-루프에서는 n 개의 작업을 1개씩 가장 일찍 끝나는 기계에 배정한다.
- Line 4: 가장 일찍 끝나는 기계의 번호인 min 을 1 (즉, 기계 M_1)로 초기화시킨다.
- Line 5~7의 for-루프에서는 각 기계의 마지막 작업의 종료 시간을 검사하여, min 을 찾는다.
- Line 8~9: 작업 i 를 기계 M_{min} 에 배정하고, $L[min]$ 을 작업 i 의 수행 시간을 더하여 갱신한다.
- Line 10: 배열 L 에서 가장 큰 값을 찾아서 리턴한다.

- 작업의 수행시간이 각각 5, 2, 4, 3, 4, 7, 9, 2, 4, 1이고, 4개의 기계가 있을 때, Approx_JobScheduling 알고리즘을 수행시킨 결과는 다음과 같다.





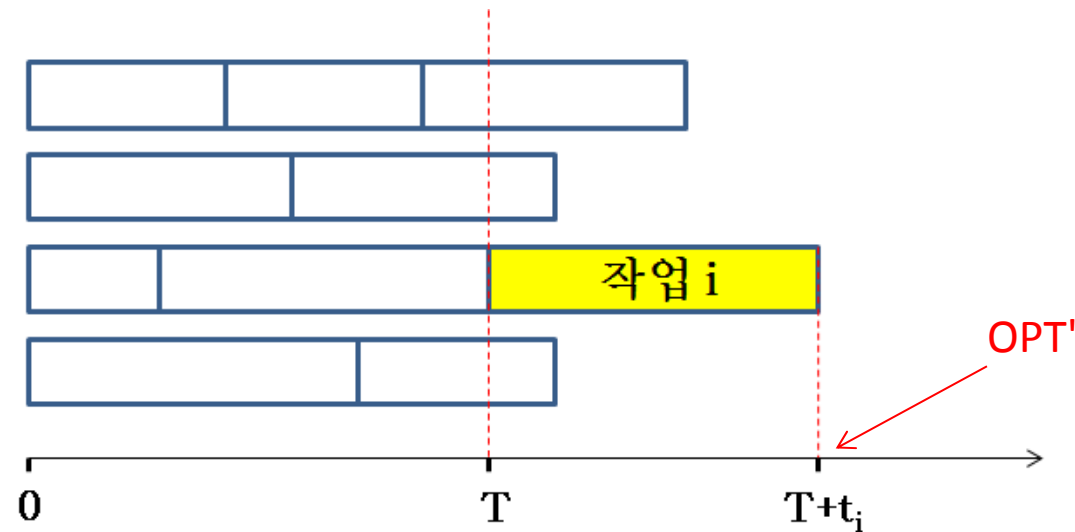
- Approx_JobScheduling 알고리즘은 가장 늦게 끝나는 작업의 종료 시간인 **13**을 리턴한다.

시간복잡도

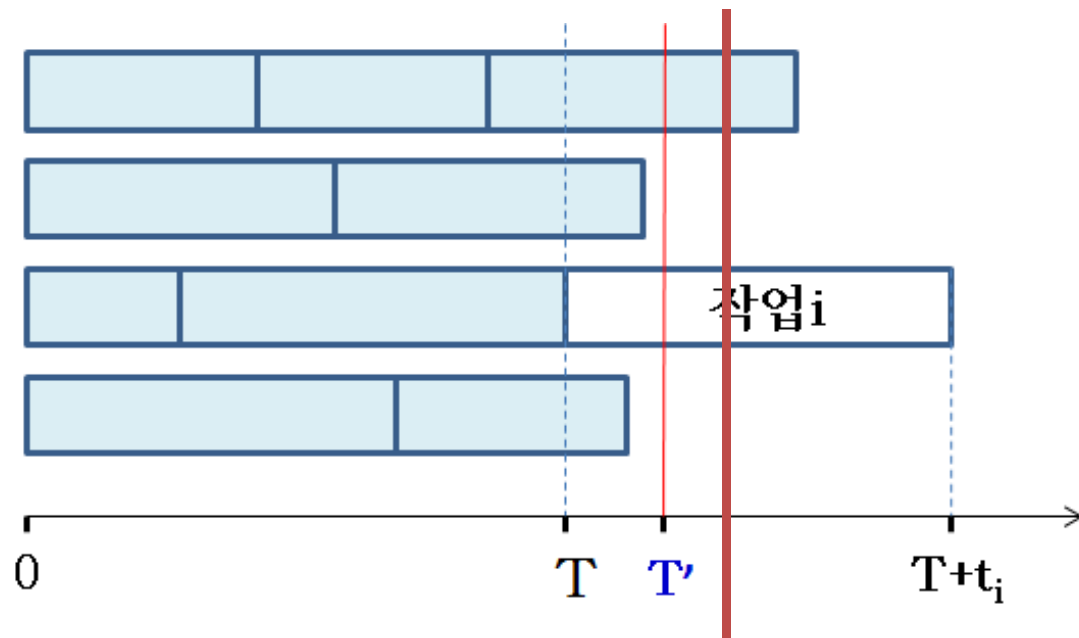
- Approx_JobScheduling 알고리즘은 n 개의 작업을 하나씩 가장 빨리 끝나는 기계에 배정한다.
- 이러한 기계를 찾기 위해 알고리즘의 line 5~7의 for-루프가 $(m-1)$ 번 수행된다.
- 즉, 모든 기계의 마지막 작업 종료 시간인 $L[j]$ 를 살펴보아야 하므로 $O(m)$ 시간이 걸린다.
- 따라서 시간복잡도는 n 개의 작업을 배정해야하고, line 10에서 배열 L 을 탐색해야하므로 $n \times O(m) + O(m) = O(nm)$ 이다.

근사 비율

- Approx_JobScheduling 알고리즘의 근사해를 OPT' 라 하고, 최적해를 OPT 라고 할 때, $OPT' \leq 2 \times OPT$ 이다.
- 즉, 근사해는 최적해의 2배를 넘지 않는다. 이를 아래의 그림을 통해서 이해해보자. 단, t_i 는 작업 i 의 수행 시간이다.



- 위의 그림은 Approx_JobScheduling 알고리즘으로 작업을 배정하였고, 가장 마지막으로 배정된 작업 i 가 T 부터 수행되며, 모든 작업이 $T+t_i$ 에 종료된 것을 보이고 있다. 그러므로 $OPT' = T+t_i$ 이다.



A lowerbound of optimal solutions (larger than T')

- 앞 그림에서 T' 는 작업 i 를 제외한 모든 작업의 수행 시간의 합을 기계의 수 m 으로 나눈 값이다.
- 즉, T' 는 작업 i 를 제외한 **평균 종료 시간**이다.
- 그러면 $T \leq T'$ 이 된다.
- 왜냐하면 작업 i 가 배정된 (가장 늦게 끝나는) 기계를 제외한 모든 기계에 배정된 작업은 적어도 T 이후에 종료되기 때문이다.

- 다음은 T와 T'의 관계인, $T \leq T'$ 를 가지고 $OPT' \leq 2 \times OPT$ 를 증명한다.

$$OPT' = t_i + T \leq t_i + T' \text{ — ①}$$

$$= t_i + \frac{(\sum_{j=1}^n t_j) - t_i}{m}, \text{ 왜냐하면 } T' = \frac{(\sum_{j=1}^n t_j) - t_i}{m}$$

$$= t_i + \frac{(\sum_{j=1}^n t_j)}{m} - \frac{t_i}{m}$$

$$= \frac{1}{m} \sum_{j=1}^n t_j + \left(1 - \frac{1}{m}\right) t_i$$

$$\leq OPT + \left(1 - \frac{1}{m}\right) OPT \text{ — ②} \quad \begin{array}{l} \text{By definition of OPT} \\ \text{(the time to finish all tasks)} \end{array}$$

$$= \left(2 - \frac{1}{m}\right) OPT$$

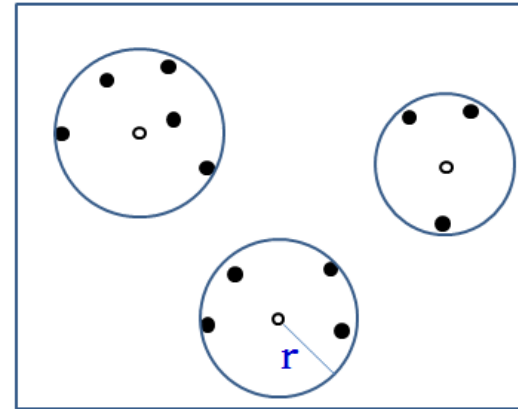
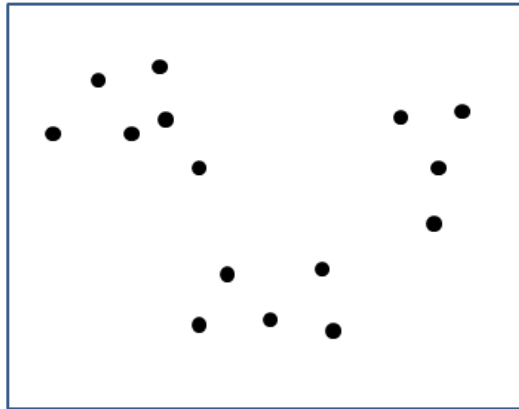
$$\leq 2 OPT \leq 2 M'$$

OPT term here is
a lowerbound of the real OPT
($OPT \leq M'$)

- 첫 번째 부등식은 위의 그림에서 살펴본 $T \leq T'$ ①을 이용한 것이다.
- 식 ②로의 변환은 최적해 OPT는 모든 작업의 수행 시간의 합을 기계의 수로 나눈 값 (평균 종료 시간)보다 같거나 크고 또한 하나의 작업 수행 시간과 같거나 크다는 것을 부등식에 반영한 것이다.

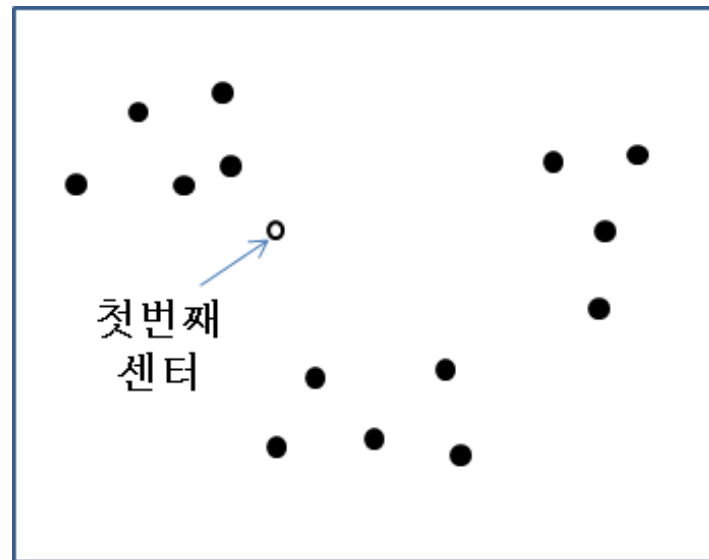
8.5 클러스터링 문제

- n 개의 점이 2차원 평면에 주어질 때, 이 점들 간의 거리를 고려하여 k 개의 그룹으로 나누고자 한다. 아래의 그림은 점들을 3개의 그룹으로 나눈 것을 보여준다.

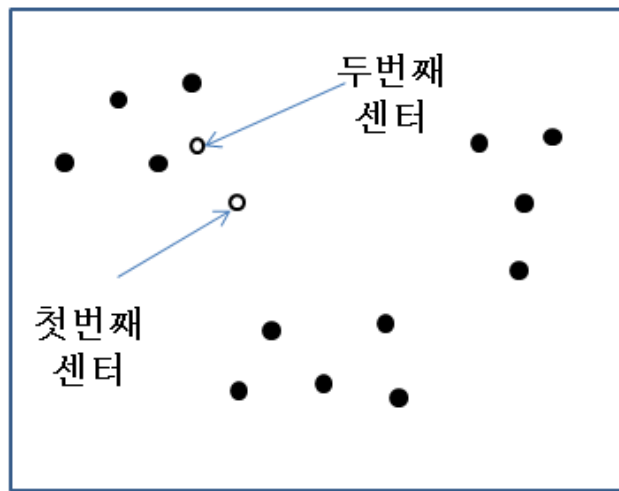


- 클러스터링 (Clustering) 문제는 입력으로 주어진 n 개의 점을 k 개의 그룹으로 나누고 각 그룹의 중심이 되는 k 개의 점을 선택하는 문제이다.
- 단, 가장 큰 반경을 가진 그룹의 직경이 최소가 되도록 k 개의 점이 선택되어야 한다.

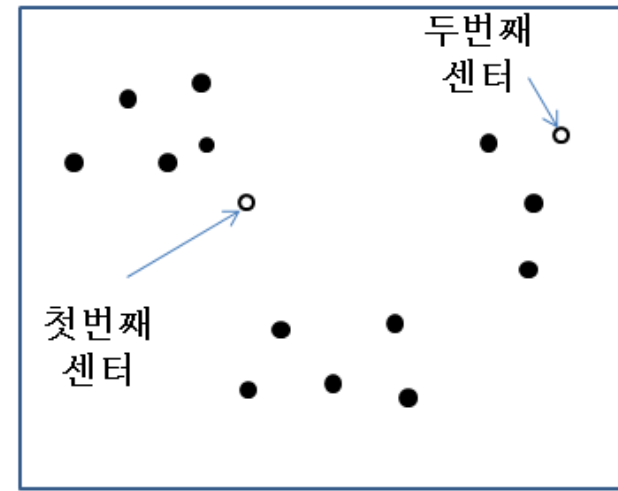
- n 개의 점 중에서 k 개의 센터를 선택하여야 하는데, 이를 한꺼번에 선택하는 것보다는 1개씩 선택하는 것이 쉽다.
- 아래의 그림 같이 1 번째 센터가 랜덤하게 정해졌다고 가정해 보자.



- 어느 점이 두 번째 센터가 되면 좋을까?
- 첫 번째 센터에서 가장 가까운 점과 가장 먼 점 중에서 어느 점이 좋을까?

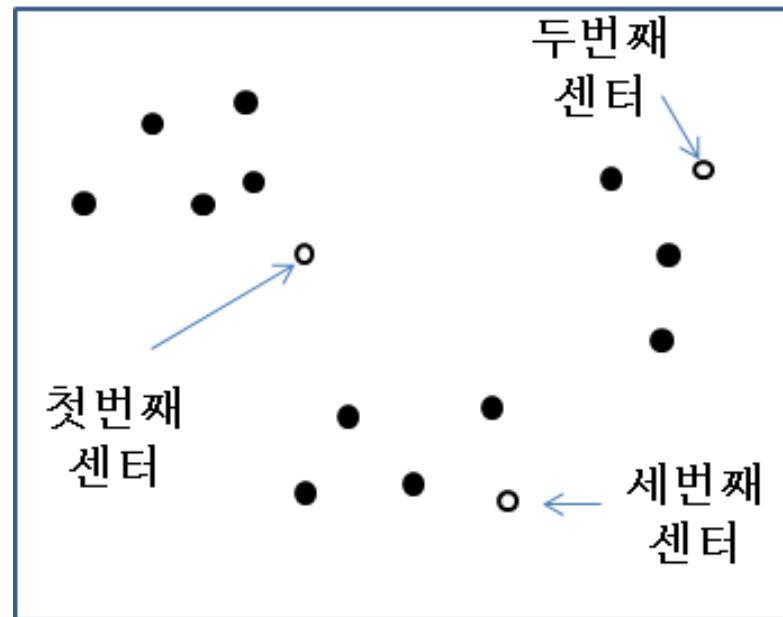


첫번째 센터에서 가장 가까운 점



첫번째 센터에서 가장 먼 점

- 위의 그림에서 보면 2개의 센터가 서로 가까이 있는 것보다 멀리 떨어져 있는 것이 좋다.
- 그 다음 세 번째 센터는 첫 번째와 두 번째 센터에서 가장 멀리 떨어진 점을 선택한다.



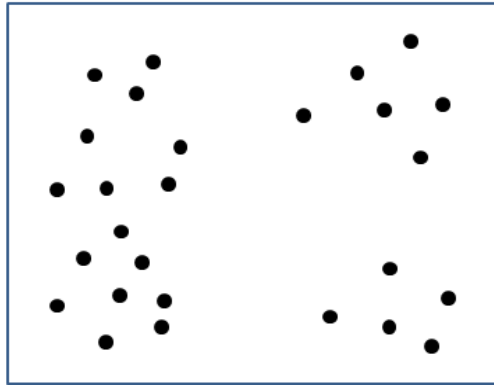
Approx_k_Clusters

입력: 2차 평면상의 n 개의 점 $x_i, i=0, 1, \dots, n-1$, 그룹의 수 $k > 1$

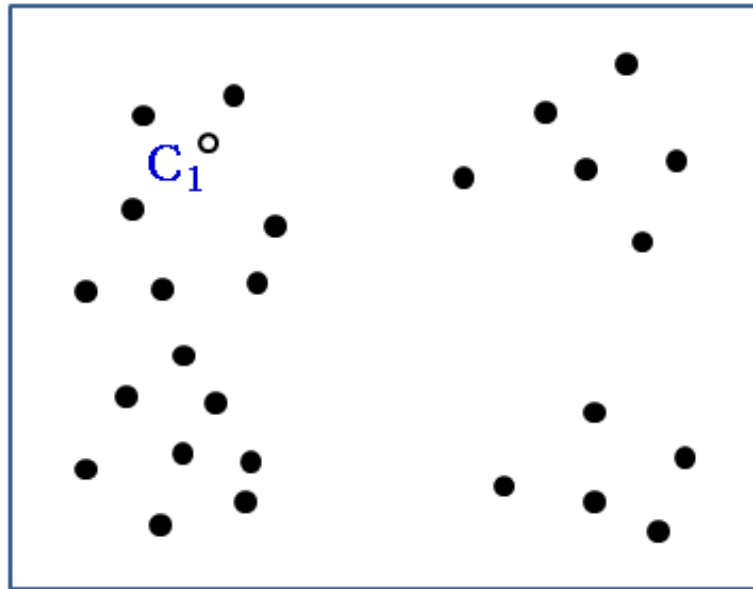
출력: k 개의 점의 그룹 및 각 그룹의 센터

1. $C[1] = r$, 단, x_r 은 n 개의 점 중에서 랜덤하게 선택된 점이다.
2. for $j = 2$ to k {
3. for $i = 0$ to $n-1$ {
4. if ($x_i \neq \text{센터}$)
5. x_i 와 각 센터까지의 거리를 계산하여, x_i 와 가장 가까운 센터까지의 거리를 $D[i]$ 에 저장한다. }
6. $C[j] = i$, 단, i 는 배열 D 의 가장 큰 원소의 인덱스이고, x_i 는 센터가 아니다. }
7. 센터가 아닌 각 점 x_i 로부터 위에서 찾은 k 개의 센터까지 거리를 각각 계산하고 그 중에 가장 짧은 거리의 센터를 찾는다. 이때 점 x_i 는 가장 가까운 센터의 그룹에 속하게 된다.
8. return 배열 c 와 센터가 아닌 각 점 x_i 가 속한 그룹의 센터

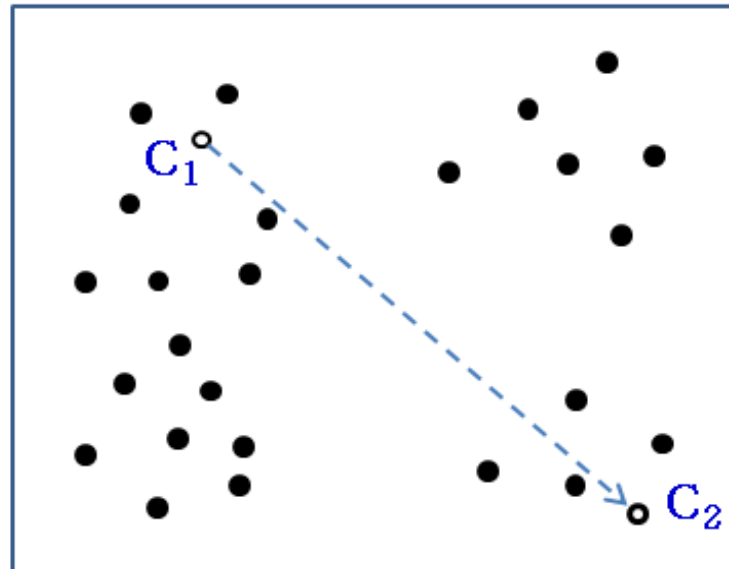
- 다음 그림에 주어진 점들을 Approx_k_Clusters 알고리즘을 수행하여 4개의 센터를 찾고, 나머지 점들을 4개의 센터를 기준으로 4개의 그룹으로 나누어 보자.



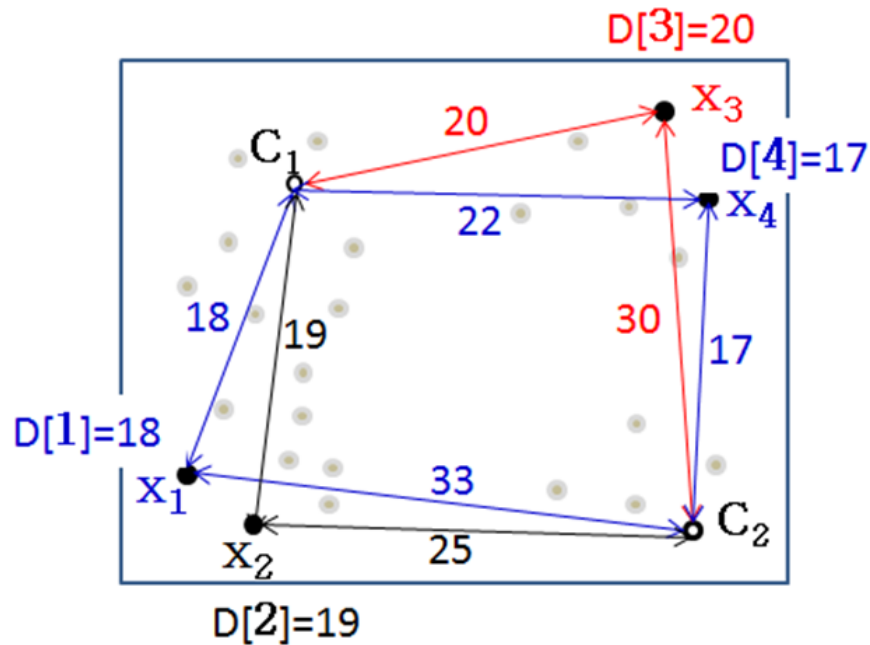
- Line 1 에서 아래의 그림처럼 임의의 점 하나를 첫 번째 센터 C_1 로 정한다.



- Line 3~5에서 c_1 이 아닌 각 점 x_i 에서 c_1 까지의 거리 $D[i]$ 를 계산한다.
- Line 6에서는 c_1 로부터 거리가 가장 먼 점을 다음 센터 c_2 로 정한다.



- Line 2: $j=3$ 이 된다. 즉, 3 번째 센터를 찾는다.
- Line 3~5: C_1 과 C_2 를 제외한 각 점 x_i 에서 각각 C_1 과 C_2 까지의 거리를 계산하여 그 중에서 작은 값을 $D[i]$ 로 정한다.
- Line 6: 배열 D 에서 가장 큰 값을 가진 원소의 인덱스가 i 라고 하면, 점 x_i 가 다음 센터 C_3 이 된다. 단, x_i 는 이전에 센터가 아닌 점이다.



$D[1]=18, \min\{\text{dist}(x_1, C_1), \text{dist}(x_1, C_2)\}$
 $=\min\{18, 33\}$ 이므로

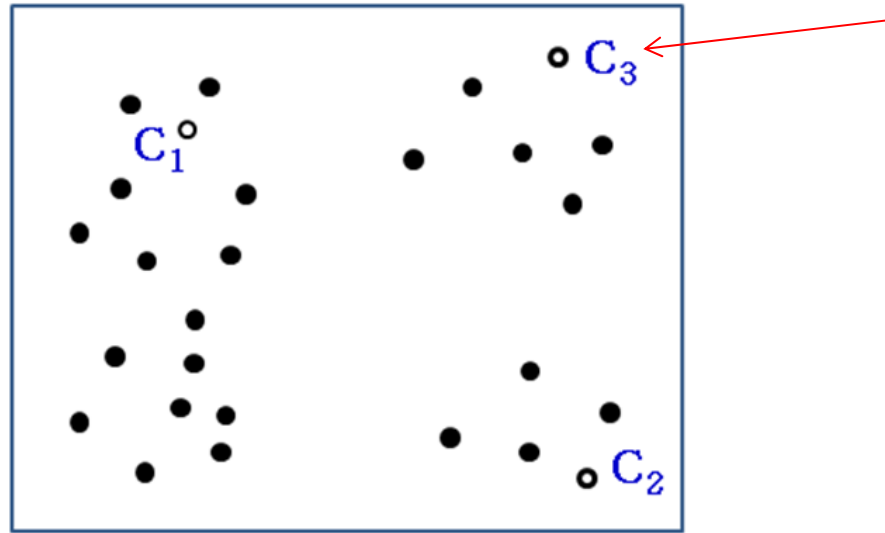
$D[2]=19, \min\{\text{dist}(x_2, C_1), \text{dist}(x_2, C_2)\}$
 $=\min\{19, 25\}$ 이므로

$D[3]=20, \min\{\text{dist}(x_3, C_1), \text{dist}(x_3, C_2)\}$
 $=\min\{20, 30\}$ 이므로

$D[4]=17, \min\{\text{dist}(x_4, C_1), \text{dist}(x_4, C_2)\}$
 $=\min\{22, 17\}$ 이므로

이외의 다른 모든 점 x_i 에 대해서 $D[i]$ 는
 17보다 작다고 가정하자.

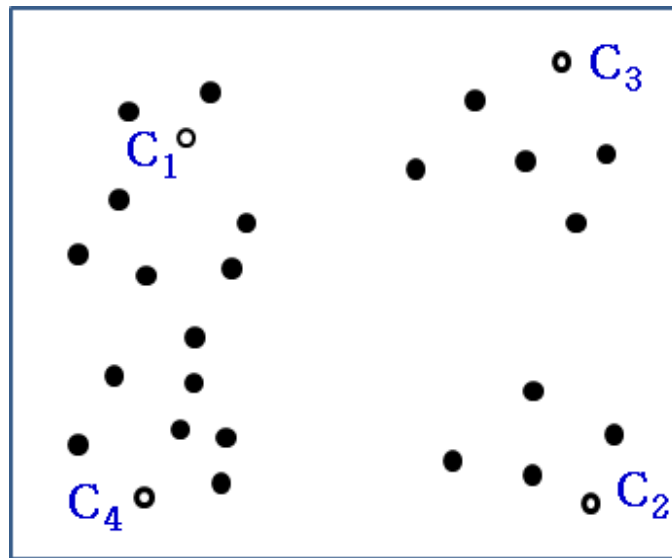
- 그림은 4개의 점 x_1, x_2, x_3, x_4 에 대해 각각 $D[x_1], D[x_2], D[x_3], D[x_4]$ 가 계산된 것을 보여주고 있다. 단, $\text{dist}(x, C)$ 는 점 x 와 센터 C 사이의 거리이다.



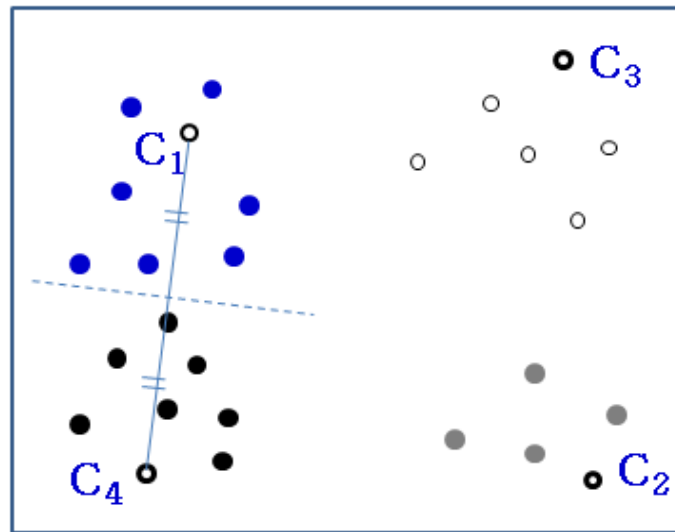
C_1 과 C_2 로부터 가장 먼 점

- 센터가 아닌 점 중에서 $D[3]$ 이 가장 큰 값이므로, 점 x_3 이 세 번째 센터인 C_3 으로 정해진 것을 그림이 보여주고 있다.

- Line 2: $j=4$ 가 된다. 즉, 4 번째 센터를 찾는다.
- Line 3~5: C_1, C_2, C_3 을 제외한 각 점 x_i 에서 각각 C_1, C_2, C_3 까지의 거리를 계산하여 그 중에서 작은 값을 $D[i]$ 로 정한다.
- Line 6: 배열 D 에서 가장 큰 값을 가진 원소의 인덱스가 i 라고 하면, 점 x_i 가 다음 센터 C_4 로 된다. 단, x_i 는 이전에 센터가 아닌 점이다.



- Line 7: 센터가 아닌 각 점 x_i 로부터 위에서 찾은 4개의 센터까지 거리를 각각 계산하고 그 중에 가장 짧은 거리의 센터를 찾는다.
- 이때 점 x_i 는 가장 가까운 센터의 그룹에 속하게 된다.
- 그리고 각 점에서 가까운 센터를 찾으면 아래와 같이 4개의 그룹으로 나누어진다.



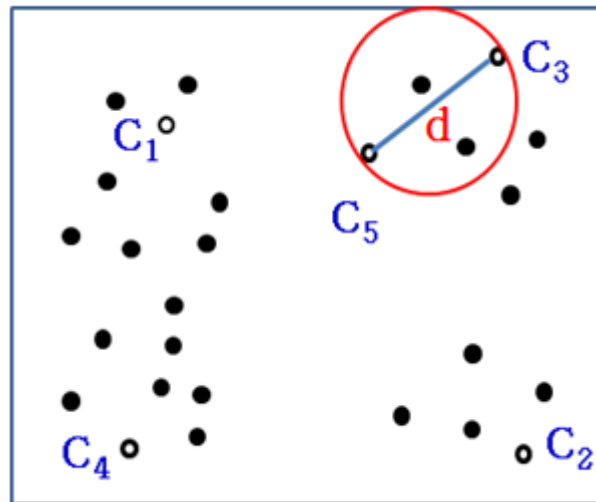
- Line 8에서는 4개의 센터와 센터가 아닌 각 점이 속한 그룹의 센터를 리턴한다.

Correct Proof

- Quality of a solution: maximum distance of a point to its assigned cluster center (let's call this as max-min distance)
- M' : optimal max-min distance
- M : lowerbound of M'
- M'' : estimate
- $M' \leq M''$ (by definition)
- $M'' < 2M'$ (By the meaning of the quality)
 - M' means the

Correct Proof

- $M' \leq M''$ (by definition)
- $M'' < 2M'$ (why??)
- Ma



시간복잡도

- Line 1: 임의의 점을 선택하므로 $O(1)$ 시간이 걸린다.
- Line 3~5: 각 점에서 각 센터까지의 거리를 계산하므로 $O(kn)$ 시간이 걸리며,
- Line 6에서는 그 중에서 최대값을 찾으므로 $O(n)$ 시간이 걸린다.
- 그런데 line 2의 for-루프는 $(k-1)$ 회 반복되므로, line 6까지의 수행 시간은 $O(1) + (k-1) \times (O(kn) + O(n))$ 이다.
- Line 7: 센터가 아닌 각 점으로부터 k 개의 센터까지의 거리를 각각 계산하면서 최솟값을 찾는 것이므로 $O(kn)$ 시간이 소요된다.
- 시간복잡도는 $O(1) + (k-1) \times (O(kn) + O(n)) + O(kn) = O(k^2n)$ 이다.

Correct Proof

- $M' \leq M''$ (by definition)
 - $M'' < 2M'$ (why??)
1. Optimal cluster set C is given (quality is M')
 2. Run the approximation algorithm
 - We pick a data point as the first cluster.
 - It should be in a cluster of C denoted by c_0

Theorem 1.

If each selected cluster is in separate clusters in C , then the next selected cluster should be in another unused cluster in C

Correct Proof

- If theorem 1 is correct, any optimal cluster has a selected cluster by the approximation algorithm.
- If the max-min distance of the optimal is d , then selected clusters can cover all data points with max-min distance lower than d . (the end of the proof)

Proof of Theorem 1?

If each selected cluster is in separate clusters in C , then the next selected cluster should be in another unused cluster in C .

Correct Proof

Definition. Occupied cluster: A cluster of optimal solution where a selected cluster center is in distance d with the center of the optimal cluster.

If a newly selected cluster is occupied again by a already occupied cluster, then all the remaining data points should be in the occupied clusters.

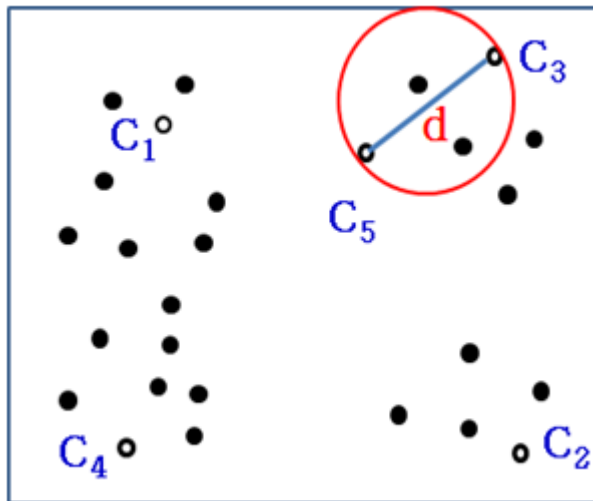
- Because their max-min distances to previously selected clusters are less than d ,
- which means that they are occupied by a cluster in the optimal solution.
- Then, removing the newly selected cluster maintains the max-min distance d of the optimal clusters.
- We can always find $d' < d$ by moving this additional cluster to around d . (next slide)
- Then, the new solution has lower d' than optimal solution , which is contradictory to the definition of the optimal solution.

Correct Proof

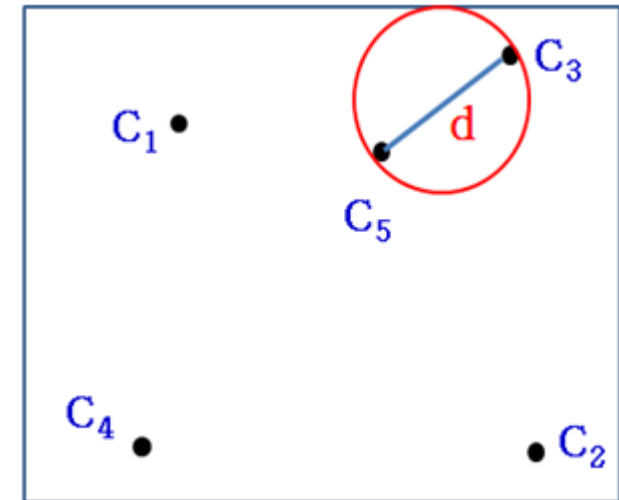
- Proof for the argument : We can always find $d' < d$ by moving this additional cluster to around d . (next slide)
- If max-min distance of a solution is d , at least one data point has distance larger than or equal to any cluster centers.
- Then, we can assign a new additional cluster to this data point, which makes its max-min distance to 0.
- Then the second largest max-min distance of a solution d_2 becomes the new max-min distance, which should be lower than or equal to d .

근사 비율 – (skip, read the English proof)

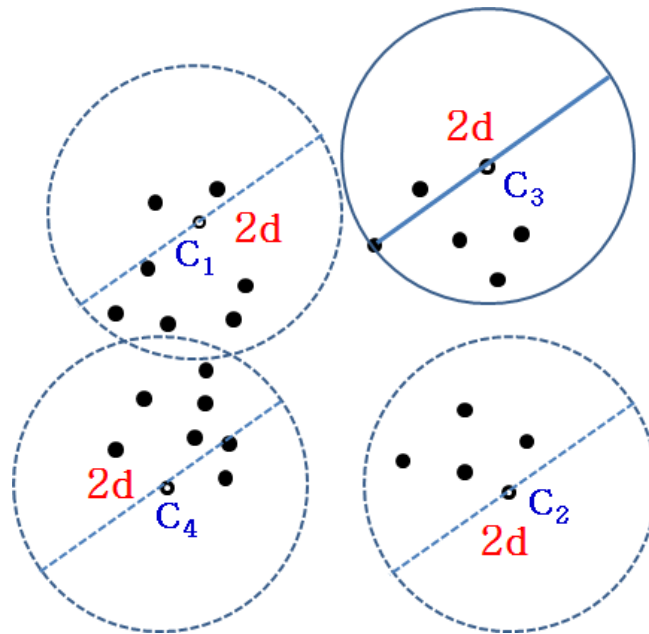
- 먼저 최적해가 만든 그룹 중에서 가장 큰 직경을 OPT라고 하자.
- 그리고 OPT의 하한을 간접적으로 찾기 위해서, Approx_k_Clusters 알고리즘이 k 개의 센터를 모두 찾고 나서 $(k+1)$ 번째 센터를 찾은 상황을 생각해보자.
- 아래의 그림은 $k=4$ 이라서 4개의 센터를 찾은 후, 1개의 센터 C_5 를 추가로 찾은 것을 나타낸다.



- c_5 에서 가장 가까운 센터인 c_3 까지의 거리를 d 라고 하자.
- 센터 배치 문제의 최적해를 계산하는 어떤 알고리즘이라도 위의 5개의 센터 점을 ($k=4$ 이니까) 4개의 그룹으로 분할해야 한다.
- 따라서 5개의 센터 중에서 2개는 하나의 그룹에 속해야만 한다.
- 그림에서는 c_3 과 c_5 가 하나의 그룹에 속한 것을 보이고 있다.
- 최적해의 가장 큰 그룹의 직경인 OPT 는 d 보다 작을 수는 없다. 즉, $OPT \geq d$ 이다.



- Approx_k_Clusters 알고리즘이 계산한 근사해의 가장 큰 그룹의 직경 OPT'는 d 와 어떤 관계인지 살펴보자.
- 다음그림을 살펴보면, 가상의 다음 센터 c_5 와 c_3 사이의 거리가 d 이므로, 센터가 아닌 어떤 점이라도 자신으로부터 가장 가까운 센터까지의 거리가 d 보다 크지 않다.
- 따라서 각 그룹의 센터를 중심으로 반경 d 이내에 그룹에 속하는 점들이 위치한다. 따라서 $OPT' \leq 2d$ 이다.



- 즉, $OPT \geq d$ 이고, $OPT' \leq 2d$ 이므로,
 $2OPT \geq 2d \geq OPT'$ 이다.
- 따라서 Approx_k_Clusters 알고리즘의 근사 비율은 2를 넘지 않는다.

응 용

- 클러스터링 알고리즘은 대단히 많은 분야에서 활용된다.
이는 일반적으로 어떠한 데이터라도 유사한 특성 (유사도)을 가진 부분적인 데이터로 분할하여 분석할 때에 사용될 수 있기 때문이다.
- 추천 시스템
- 데이터 마이닝 (Data Mining)
- VLSI 설계
- 병렬 처리 (Parallel Processing)
- 웹 탐색 (Web Searching)
- 데이터베이스
- 소프트웨어 공학 (Software Engineering)
- 컴퓨터 그래픽스 (Computer Graphics)

응 용

- 패턴 인식 (Pattern Recognition)
- 유전자 분석 (Gene Analysis)
- 소셜 네트워크 (Social Network) 분석
- 도시 계획, 사회학, 심리학, 의학, 금융, 통계, 유통 등 많은 분야에서 데이터의 그룹화는 매우 중요한 문제이다.

요약

- 근사 알고리즘은 최적해의 값에 가까운 해인 근사해를 찾는 대신에 다항식 시간의 복잡도를 가진다.
- 근사 비율 (approximation ratio)은 근사해가 얼마나 최적해에 가까운지를 나타내는 근사해의 값과 최적해의 값의 비율로서, 1.0에 가까울수록 실용성이 높은 알고리즘이다.
- 여행자 문제를 위한 근사 알고리즘은 최소 신장 트리의 모든 점을 연결하는 특성과 최소 가중치의 특성을 이용한다. 근사비율은 2이다.
- 정점 커버 문제를 위한 근사 알고리즘은 그래프에서 극대 매칭을 이용하여 근사해를 찾는다. 근사비율은 2이다.

- 통 채우기 문제는 최초 적합 (first fit), 다음 적합 (next fit), 최선 적합 (best fit), 최악 적합 (worst fit)과 같은 그리디 알고리즘으로 근사해를 찾는다. 근사비율은 각각 2이다.
- 작업 스케줄링 문제는 가장 빨리 끝나는 기계에 새 작업을 배정하는 그리디 알고리즘으로 근사해를 찾는다. 근사비율은 2이다.
- 클러스터링 문제는 현재까지 정해진 센터에서 가장 멀리 떨어진 점을 다음 센터로 정하는 그리디 알고리즘으로 근사해를 찾는다. 근사비율은 2이다.