

Memory Utilization

Paul Franklin, Yangqi Su, Yaying Shi

1.7 Memory utilization

Input: A disk of size D . A set F of n files of size f_1, \dots, f_n .

Output: A subset $S \subseteq F$ of the files that can fit on the disk.

Metric: Maximize the number of file in the disk.

(Bonus question: adapt the algorithm for knapsack, why isn't it optimal?)

Presentation of the problem

- Input:

For any n files set, we have size set F for those files.

$$F = \{f_i \mid \text{where } i \text{ from } 1 \text{ to } n\} = \{f_1, f_2, f_3, \dots, f_n\}$$

You have a capacity of D which is a size limit.

- Output:

A subset of F

$$S = \{s_i \mid \text{where } i \in S\} = \{s_1, s_2, s_3, \dots, s_k\}$$

- Metric

Max(k)

Important of this problem

- Same as Knapsack problem but all the object value is same.
- Any arrangement and schedule problem with a threshold
 1. If you have an exam tomorrow, how could you select the most value part in limited time.
 2. How could you schedule your travel plan with limited money?
 3. Data download management
 4. Cache select

Properties

- Strong Dominance property

For question we have all optimal solutions $S^* = (k, s_1, \dots, s_{k-1}, s_k)$ that D -
 $\sum_{i \in k} s_i < s = \{s \in F \cap \overline{S^*}\}$

- Weak Dominance property

There is an optimal solution that all the elements is in increased order in F.
That's to say, S^* is first k elements of sorted F.

Proofs

- Strong Proofs: Assume $S^* = (k, s_1, \dots, s_k)$ is an optimal solution of problem $I = (D, f_1, \dots, f_n)$ of k files, but $D - \sum_{i \in k} s_i >= s \{s \mid s \in F \cap \overline{S^*}\}$. Then, we can construct a valid solution $S = (k+1, s_1, \dots, s_k, s_i)$ for problem I where $\{s_i \mid (s_i \in F \cap \overline{S^*}) \text{ and } \sum_{i \in k+1} s_i < D\}$ of $k+1$ files. It's conflicts with S^* is optimum for I .
- Weak Proofs: For any solution $S = (i, s_1, \dots, s_i)$, we can constructor a solution $S' = (k', s'_1, s'_2, \dots, s'_k)$ by replace each elements in the solution by one or more smaller elements in F , which can guarantee $k <= k'$. That would also be true for an optimal solution.

Greedy Algorithm

- Algorithm

```
for i in range(len(F)):  
    for j in range(len(F)-1)  
        find smallest;  
        sum += smallest;  
        if (sum > D):  
            break  
    return index
```

- Complexity:

$$O(n * n) = O(n^2)$$

Complexity and optimized

- Optimized:
 1. Sort F to F'
 2. for item in range(len(F')):
 sum += item;
 if (sum > D):
 break
 return index

- Complexity

$$O(n * \log n) + O(n) = O(n * \log n)$$

Adapted Algorithm for knapsack

- Problem

- Input:

- For any n object set, we have size set S and value set V for those objects.

- $S = \{s_i \mid \text{where } i \text{ from } 1 \text{ to } n\} = \{s_1, s_2, s_3, \dots, s_n\}$

- $V = \{v_i \mid \text{where } i \text{ from } 1 \text{ to } n\} = \{v_1, v_2, v_3, \dots, v_n\}$

- You have a capacity of C which is a size limit.

- Output:

- A subset of S

- $T = \{s_i \mid \text{where } i \in T\}$

- Metric

- $\text{Max} \sum_{i \in T} v_i$

- Algorithm

- Sort S to S'

- for item in range(len(S')):

- sum += item;

- val += V[i];

- if (sum > V):

- break

- return val

Why is not optimal in knapsack?

Memory Utilization is a special case in knapsack. When the smaller size object has same or great value than larger one, the knapsack is same as the Memory Utilization. But in knapsack problem, size is not the only threshold for question, the value is another hidden threshold. You must satisfied both of them. In Memory Utilization, when you satisfied one threshold, another is automatically satisfied.

Presentation of the problem

Input

D: size

F = (1, ..., n) : 'set of files'

f_i ≥ i ∈ F: file size of file i

This definition
changes to vector
notation.

Output

S ⊆ F: solution; set of files on disk

Metric

max(|S|): maximize the solution set cardinality

Without a Dominance Property

An optimal solution to the problem or any subproblem of the problem has the properties:

$\max(|S|)$: maximize the solution set cardinality

and

combination satisfying the cost constraint

WHY: Just find one (with as little effort as possible)

Note: It is impossible to compute what the value equivalent to the metric is until all combinations satisfying the cost constraint have been computed. Alternatively, the metric is a function of the set of ALL combinations satisfying the cost constraint. Also, the metric is a function of information that has not been obtained at any nonfinal, local decision.

(Weak) Dominance Property

Dominance Rule: eliminating uninteresting solutions, or selecting interesting solutions.

$$S' = \{(1, 1, 1, 1), (1, 1, 1, 2)\}$$

IF dominance rule selects minimum size,
then $S'' = \{(1, 1, 1, 1)\}$.

• IF, $\forall s \in S$, $s \in S''$, THEN "Strong!"

Formula

Parallel Computing

$$MU(D, F) \exists i, [f_i] \in F = S \exists \forall f \subseteq F, |f| = \max(|f|) \wedge \sum_{i \in I(f)} f_i \leq D$$

$$\rightarrow f \in S' \wedge \forall s \in S', s = s^*$$

Formula (Dominance Property)

Dominance Rule (1), weak

$$MU(D, F) \exists i, [f_i] \in F = S \exists f \subseteq F, |f| = \max_{f_i \in F} (|f_i|) \wedge \sum f_i \leq D$$

$$\rightarrow f \in S' \wedge f \in S, \sum_{i \in I \setminus f} f_i = \min_{i \in I \setminus f} (\sum f_i) \rightarrow f \in S'' \wedge S' \subseteq S,$$

$$S = S^*$$

Alternatively, parallel computing distribution is a function of the first element. However, if our intent is not to implement a dynamic programming algorithm, there is a contradiction.

Formulas

We can write a greedy algorithm in its familiar form from the weak dominance rule, but it is not the same problem as defined at the beginning. When we change the metric from $\max(|S|)$ to $\min(\sum_{i \in S} f_i)$, the goal function no longer requires information from later steps and we substitute the recursive S^* for S . The locally-optimal decision becomes $\min_{i \in E \setminus S^*} (f_i)$ and the base case is $\forall f_i \in E \setminus S^*, f_i > D^*$.

"Globally
Optimal" Metric

Proof(s)

we can prove that the dominance rule generates a subset of the solution set satisfying the original metric; that it is a valid solution to the problem and neglect that there may be other valid solutions or that the solution set is a subset of the original solution set.

Proof(s) - Optimal Substructure

Let $T_k \subseteq S$.

Either assume all elements have distinct f_i , or
that an arbitrary element or an element minimizing
computational time complexity is selected from
the extremes.

$$\max_{i \in T_k} (f_i) = \min_{i \in F-T_{k-1}} (f_i)$$

Proof(s) - Optimal Substructure

If $T_k = (t_1, \dots, t_k)$ is an optimal solution
of problem $I_k = (D^*, I_1, \dots, I_n, T_{k-1})$ of value
 $(|T_k| = T^*, \sum_{k \in T_k} f_k = f^*)$, then $T_{k-1} = (t_1, \dots, t_{k-1})$ is an
optimal solution of problem $I_{k-1} = (D^* - T^*, I_1, \dots, I_n, T_{k-2})$
of value $(|T_{k-1}| = |T_k| - 1 = T^* - 1, \sum_{k \in T_{k-1}} f_k = f^* - \max_{i \in I \in T_k} (f_i))$.

This could be further edited to say that the set of all vectors in any instance of input is represented by the variable F^* , where $F^*_{\leq k} = F - T_{k-1}$.

This recursion achieves the dominance rule, while also computing the metric, and this is how I've chosen to define what is optimal with respect to the problem. I go on to prove that the solutions generated by the dominance rule are always a subset of the solutions generated by the metric.

Proof(s) – Optimal Substructure

Let $F_k^* = F - T_{k-1}$.

Let $T_k = (t_1, \dots, t_k)$ be an optimal solution of problem $I_k = (D^*, F_{k-1}^*)$ with value $(|T_k| = T^*, \sum_{k \in T_k} f_k = f^*)$. This is true because my definition of optimal is that $\sum_{k \in T_k} f_k = f^* - \max_i (f_i)$. It is then valid that $T_{k-1} = (t_1, \dots, t_{k-1})$ is an optimal solution of problem $I_{k-1} = (D^* - T^*, F_{k-1}^*)$ of value $(|T_{k-1}| = T^* - 1, \sum_{k \in T_{k-1}} f_k = f^* - \max_i (f_i))$. However, I am going to assume that this conclusion is not true and through some pair of contradicting relations determine that it must strictly be sound.

If the conclusion is false, then $\sum_{k \in T_{k-1}} f_k \neq f^* - \max_i (f_i)$. by the assumption. However, it is known that $f^* = \sum_{k \in T_{k-1}} f_k + \max_i (f_i) = \sum_{k \in T_{k-1}} f_k + \min_i (f_i)$. However, the pair of relations $\{\sum_{k \in T_{k-1}} f_k = f^* - \max_i (f_i), \sum_{k \in T_{k-1}} f_k \neq f^* - \max_i (f_i)\}$ cannot simultaneously exist; they are contradicting. By contradiction, the argument is now both valid and has true premises; it is sound.

Proof(s) – Dominance Property is Weak

The dominance property generalizes to the metric; the solutions generated by the dominance property are a subset of the solutions generated by the metric.

Note: This is not a proof that the dominance rule is globally optimal. It is a proof that the dominance rule generates a subset of the solutions that satisfy the original metric. Now, it is unnecessary to compute all possible solutions to determine what the value of metric is.

Proof(s) – Dominance Property is Weak

$$\forall I \in S'', |I| = \max(|S|) \in \sum_{j \in S \setminus I} f_j \leq D$$

Equivalence (1)

The solutions without the dominance rule satisfy the metric and the capacity constraint.

$$\forall I \in S', |I| = \max(|S|) \in \sum_{j \in S \setminus I} f_j \leq D$$

Equivalence (2)

The solutions with the dominance rule satisfy the metric and the capacity constraint.

$$f^* = \sum_{k \in T_{k-1}} f_k + \min_{i \in I \setminus F - T_{k-1}} (f_i)$$

Equivalence (3)

Refer to the proof of optimal substructure.

$$f^* = \sum_{k \in T_k} f_k$$

Equivalence (4)

Refer to the proof of optimal substructure.

$$T_k - T_{k-1} = \{ I_k \ni \min_{i \in I \setminus F - T_{k-1}} (f_i) \}$$

Equivalence (5)

By Algebra on {3, 4}.

$$T_0 = \{\}$$

Equivalence (6)

Definition of base case.

$$\forall I \in S', I = T_{|I|}$$

Equivalence (7)

From equivalence rules {5, 6}, derived from the definition of the dominance rule.

Proof(s) – Dominance Property is Weak

$$S' \subseteq S''$$

Subset (8)
Implication of the conjunction of relations {1, 2}.

$$S' \subset S''$$

Proper Subset (9)
Implication of the conjunction of relations {7, 8}.

$$S' \subset S'' \rightarrow \text{"weak"}$$

Definition (10)
From relations {9}, Proof by definition

Actually, the proof depends on an unspecified constraint on the input set F where the solution set without the dominance property for any capacity must have multiple solutions. Otherwise, there is a unique solution, the sets are equal and the dominance property function is strong.

Proof(s) – Dominance Property Implies Extreme

The last thing that could be proven is that if (7) is true, then (2) is also true.

Proof(s) – Dominance Property Implies Extreme

$$\forall I \in S', I \subseteq F$$

$$\forall I \in S', \sum_{j \in S' \setminus I} f_j \leq D$$

$$\forall I \in S', |I| \leq |F|$$

$$\begin{aligned} \forall J \in F - T_{k-1}, f_j \in J \geq f_i, \forall i \in I \\ \forall I \in T_{k-1} \end{aligned}$$

$$12 \wedge 13 \wedge 14 \rightarrow 2$$

Subset (11)

Defining constraints on the solution set

Inequality (12)

Defining constraints on the solution set

Inequality (13)

From relations {11}, by definition of cardinality

Inequality (14)

From relations {3, 4}, by Algebra and definition of minimum

Contradiction (15)

From relations {12, 13, 14}, by contradiction

$$12 \wedge 13 \wedge 14 \rightarrow \neg 2$$

Choosing a larger file size, results in less capacity consumption so that more files can be consumed.

Proof(s) – Dominance Property Implies Extreme

$$\forall I \in S^1, \sum_{j \in J \subseteq I} f_{j_k} = f_{j_1} + f_{j_2} + \dots + f_{j_{n-1}} + f_{j_n} \leq D$$

$$f_{j_1} + f_{j_2} + \dots + f_{j_{n-1}} \leq D - f_{j_n}$$

Let $\exists n^* \ni f_{j_{n^*}} < f_{j_n}$. Then, $D - f_{j_{n^*}} > D - f_{j_n}$.

That is, the file size is inversely proportional to the capacity consumption.

Proof(s) – Dominance Property Implies Extreme

Let big file not fit
 $f_{j,n} > D - f_{j,n-1} - \dots - f_{j,2} - f_{j,1}$

and small file(s) fit
 $f_{j,n^*} \leq D - f_{j,n^*-1} - \dots - f_{j,2} - f_{j,1}$

but that less small files fit than big files
 $n^* < n$

it follows that at least as many small files fit
 $n^* = 1 + \dots + (n-1) \geq n$

contradicting less small files fit
 $n^* < n$

then at least as many small files fit
 $n^* \geq n$

Formulas

$$MU(D^*, S^*) = \min_{f_i \in F - S^*} (D^* - f_i, S^* \cup \{I\}) + f_i$$

or +1
or set union

$$MU(D^*, S^*) \ni \forall f_i \in F - S^*, f_i > D^* = \emptyset$$

where the first of any set of duplicate f_i is taken and the initialization variables are (D, \emptyset) where $\emptyset = \{\}$ is the empty set.

Disk Space; File Size
Set Cardinality
Vector Set
Base Case

Complexity

Complexity without Dominance Property
 $(O(n) \times O(n) \times O(1)) = O(n^2)$

Assume that the order of elements in F are sorted according to f_i , increasing.

Then, complexity with dominance property is:

$$O(O(1) \times O(n) \times O(1)) = O(n)$$

Even a comparison sort with asymptotically optimal time complexity would be advantageous.

However, if the elements in F are randomly ordered with respect to f_i , then the complexity with dominance property is:
 $O(O(1) \times O(n) \times O(n)) = O(n^2)$

Sources

1: Antoine Jouglet, Jacques Carlier. Dominance rules in combinatorial optimization problems. European Journal of Operational Research, Elsevier, 2011, 212, pp.433-444. ff10.1016/j.ejor.2010.11.008ff. fffhal00625751

<https://hal.archives-ouvertes.fr/hal-00625751/document>

2: Advanced Algorithms: Greedy Algorithms, Spring 2020, Erik Saule

Please cite if you use any part of our solution in your work. This includes, but is not limited to, consumptions, adaptations, transformations, etc., ...