

Lab03-Matroid & DP

CS2308-Algorithm and Complexity, Xiaofeng Gao, Spring 2023.

- * Contact TA Jiasen Li for any questions. Include both your .pdf and .tex files in the uploaded .rar or .zip file.
* Name:张祥东 521030910206 Time Spent on this Lab:5-6h

1. (Matroid) Provide an example of (S, \mathbf{C}) , which is an independent system but not a matroid. Give an instance of S such that $v(S) \neq u(S)$ (should be different from the examples in class).

Definition 1. S consists of n objects, which have weight.

Definition 2. C consists of subsets of S the weight summation of which is less than or equal to the given weight limit X .

For the subset B and the subset A of S . If $B \subset A$ and $A \in C$:

Due to $A \in C$, the summation of elements in $A \leq$ the given X . And due to $B \subset A$, obviously the summation of elements in $B \leq$ the given X . Namely

$$B \subset A \text{ and } A \in C \rightarrow B \in C$$

It's hereditary property, thus the system is a independent system.

However, it's not a matroid.

Let us give a concrete instance to illustrate this.

- $X = 10$
- $n = 10$ and $\text{weight}(a_1) = 5, \text{weight}(a_2) = 5, \text{weight}(a_3) = \dots = \text{weight}(a_{10}) = 1$

Easy to find that

$$u(S) = 2$$

$$v(S) = 8$$

We know (S, C) is a matroid iff for any $F \subseteq S, u(F) = v(F)$.

So this independent system is not a matroid.

2. (Knapsack & Greedy) Let $A = [a_1, a_2, \dots, a_n]$, $\text{weight}(a_i), \text{volume}(a_i) \in \mathbb{Z}^+, \forall 1 \leq i \leq n$. You should find a subset of m items with maximum density. Your time complexity should be at most $O(n \log n)$. Give your algorithm in the form of pseudo-code, prove the correctness and analyze its time complexity.

Algorithm 1: find the maximum desity subset of size m

Input: $A = [a_1, a_2, \dots, a_n]$, weight, volume**Output:** DensityNeeded and the solution

```
1  $maxDensity \leftarrow 0$ ;
2  $minDensity \leftarrow \infty$ ;
3 for  $i \leftarrow 1$  to  $n$  do
4    $maxDensity \leftarrow \max(maxDensity, \frac{weight(a_i)}{volume(a_i)})$ ;
5    $minDensity \leftarrow \min(minDensity, \frac{weight(a_i)}{volume(a_i)})$ ;
6  $b \leftarrow$  new array of size  $n$ ;
7  $solution \leftarrow$  new array of size  $m$  and it is not assigned now;
8 while  $abs(maxDensity - minDensity) \geq 1e - 5$  do
9    $mid \leftarrow (maxDensity + minDensity)/2$ ;
10  for  $i \leftarrow 1$  to  $n$  do
11     $b_i \leftarrow (weight(a_i) - mid * volume(a_i), i)$ ;
12  sort first element in the the array  $b$  by the relation  $\geq$ .
13   $isValid \leftarrow 0$ 
14  for  $i \leftarrow 1$  to  $m$  do
15     $isValid \leftarrow isValid + b_i.first$ ;
16  if  $isValid < 0$  then
17     $maxDensity = mid$ ;
18  else
19     $solution \leftarrow b[1, \dots, m].second$ ;
20     $minDensity = mid$ ;
21 if  $solution$  has not been assigned then
22   for  $i \leftarrow 1$  to  $n$  do
23      $b_i \leftarrow (weight(a_i) - minDensity * volume(a_i), i)$ ;
24   sort first element in the the array  $b$  by the relation  $\geq$ .
25    $solution \leftarrow b[1, \dots, m].second$ ;
26 return  $minDensity, solution$ ;
```

We take **density needed** denoting the maximum density that we need to figure out in this question.

Here comes the proof of the correctness of our algorithm:

Obviously, the maximum density that satisfy the demand in the question should satisfy

$$\min\{\frac{weight(a_i)}{volume(a_i)}\} \leq \text{density needed} \leq \max\{\frac{weight(a_i)}{volume(a_i)}\}$$

For every density that is less than the *density needed*, it can be satisfied, and for every density that is greater than the *density needed*, it cannot be satisfied. Thus we can do binary search to find the answer.

Only if we find the maximum density that can be satisfied by the array a , the density needed must be that number. So to prove our algorithm's correctness, we've already known through previous state that our algorithm can give the right answer if the operation in the algorithm is right.

So next, we prove the correctness of our operations.

For mid , we want to test

mid is a proper answer

$$\begin{aligned}
&\iff \text{there exist a sequence } a_{i_1}, a_{i_2}, \dots, a_{i_m} \text{ that } \frac{\sum_{j=1}^m \text{weight}(a_{i_j})}{\sum_{j=1}^m \text{volume}(a_{i_j})} \geq mid \\
&\iff \text{there exist a sequence } a_{i_1}, a_{i_2}, \dots, a_{i_m} \text{ that } \sum_{j=1}^m \text{weight}(a_{i_j}) \geq mid * \sum_{j=1}^m \text{volume}(a_{i_j}) \\
&\iff \text{there exist a sequence } a_{i_1}, a_{i_2}, \dots, a_{i_m} \text{ that } \sum_{j=1}^m (\text{weight}(a_{i_j}) - mid * \text{volume}(a_{i_j})) \geq 0
\end{aligned} \tag{1}$$

So if we could find out whether there exist a sequence of $a_{i_1}, a_{i_2}, \dots, a_{i_m}$ that satisfies $\sum_{j=1}^m (\text{weight}(a_{i_j}) - mid * \text{volume}(a_{i_j})) \geq 0$, the legality of the mid can be proved.

So we get new array b for every mid

$$b[\cdot] \leftarrow \text{weight}(a[\cdot]) - mid * \text{volume}(a[\cdot])$$

we try to find the biggest summation subset of b with size m using greedy. By sort the array b with \geq and calculate the first m elements' summation. If the summation is smaller than 0, the mid is illegal, or it's legal.

As for why the biggest summation subset can be found with greedy, we'd like to prove it by contradiction.

First, we get sequences of size m . The optimal sequence taken as a^{opt} . The sequence obtained by greedy taken as a^{alg} . Sort the a^{opt} and a^{alg} by the relation geq .

We assume the $a^{opt} \neq a^{alg}$

We assume that $a_i^{opt} = a_i^{alg}$ until i , namely

$$a_i^{opt} \neq a_i^{alg}$$

And from our selection method, the inequality $a_i^{alg} > a_i^{opt}$ must exist. So if we replace the a_i^{opt} with a_i^{alg} , the optimal solution will be better. It's impossible for an optimal solution to get better.

So the i doesn't exist, the assumption is wrong, namely $a^{opt} = a^{alg}$.

Above all, the proof of my algorithm's correctness is completed.

As for the time complexity. By limit the precision into $1e-5$, we can consider the complexity for binary search is $O(k * 1e5) = O(1)$ times.

* k is the $\text{maxDensity} - \text{minDensity}$ in my pseudo-code and it can be constrained by a constant here, namely $k = O(1)$.

In the While part, the for in it is $O(n + m)$ and the sort in it is $O(n \log n)$.

So the time complexity is

$$O(k * 1e5 * (n + m + n \log n)) = O(n \log n)$$

It's

$$O(n \log n)$$

.

3. (DP Application) Assuming you are a programmer for Word, you are now interested in studying how to make document formatting neatly. We model the problem as follows: Denote the article as a sequence $A = \langle s_1, s_2, \dots, s_n \rangle$ of strings, each with a given length satisfying $\text{len}(s_i) \leq k, \forall 1 \leq i \leq n$. Your job is typesetting these strings into a Word file and deciding how many words are in each line. There should be a space between every two adjacent words in a line. Commonly, a line looks good with m characters, including spaces. During the composition, if the number of characters in a line is m' , there is a penalty of $(m - m')^2$. Design an algorithm to minimize the sum of penalties with time complexity at most $O(n \log n)$. Give your algorithm in the form of pseudo-code and analyze its time complexity.

Algorithm 2: DP application

Input: a sequence $A = \langle s_1, s_2, \dots, s_n \rangle$ of strings and the given m

Output: Solution to determine after which string the line change should be applied to get the minimize penalty and the minimum penalty

```
1  $pre \leftarrow \text{size } n + 1$  and index from 0 to  $n$  with initial value 0;
2 for  $i \leftarrow 1$  to  $n$  do
3    $pre_i \leftarrow pre_{i-1} + 1$ ;
4  $from \leftarrow \text{size } n$  and with initial value -1;
5  $dp \leftarrow \text{size } n + 1$  and index from 0 to  $n$  with initial value 0;
6  $dp_1 \leftarrow dp_0 + (m - (pre_1 - 1))^2$ 
7  $convexPack$  is a convex pack and maintain this  $convexPack$  in  $O(\log n)$  for every added
   point operation using a balance tree.
8  $convexPack \leftarrow \text{point}(pre_0, dp_0 + (pre_0 - 1)^2 + 2mpre_0 + 2m + m^2)$ ;
9  $convexPack \leftarrow \text{point}(pre_1, dp_1 + (pre_1 - 1)^2 + 2mpre_1 + 2m + m^2)$ ;
10  $from[1] \leftarrow 0$ ;
11  $pointStack$  is a point sequence maintaining counterclockwise point sequence in the
    $convexPack$ 's edge and renewed danamically with the  $convexpack$  which only contains
   the points laying at the down convex pack.
12 for  $i \leftarrow 2$  to  $n$  do
13    $k \leftarrow 2pre_i$ ;
14    $l \leftarrow 2$ ;  $r \leftarrow \text{sizeof}(pointStack)$ ;
15   while  $l < r$  do
16      $mid \leftarrow (l + r) / 2$ ;
17      $cur = (pointStack_{mid}.second - pointStack_{mid-1}.second) / (pointStack_{mid}.first -$ 
        $pointStack_{mid-1}.first)$ ;
18     if  $cur < k$  then
19        $l \leftarrow mid + 1$ ;
20     else
21        $r \leftarrow mid$ ;
22    $ans1 \leftarrow -k * pointStack_l.first + pointStack_l.second$ ;
23    $ans2 \leftarrow -k * pointStack_{l-1}.first + pointStack_{l-1}.second$ ;
24   if  $ans1 < ans2$  then
25      $dp_i = ans1 + pre_i^2 - (2m + 2)pre_i$ ;
26      $from[i] \leftarrow l$ ;
27   else
28      $dp_i = ans2 + pre_i^2 - (2m + 2)pre_i$ ;
29      $from[i] \leftarrow l - 1$ ;
30    $convexPack \leftarrow \text{add point}(pre_i, dp_i + (pre_i - 1)^2 + 2mpre_i + 2m + m^2)$ ;
31  $solution \leftarrow \text{empty}$ ;
32  $cur \leftarrow from[n]$ ;
33 while  $from[cur] \neq 0$  do
34    $solution.push\_front(cur)$ ;
35    $cur \leftarrow from[cur]$ ;
36 return  $dp_n, solution$ ;
```

Definition 3. $pre \rightarrow$ an array of size $n + 1$. pre_i denotes the characters before i_{th} string (including i_{th} string) and add $i - 1$ characters which are spaces among i strings.

Definition 4. $dp \rightarrow$ an array of size $n + 1$. dp_i denotes the minimize penalty when take the i_{th} string as the end string.

We consider this question in one demesion perspective, here comes the OPT function

$$dp_i = \begin{cases} \min_{j=0}^{i-1} \{dp_j + (m - (pre_i - pre_j - 1))^2\} & i \geq 1 \\ 0 & i = 0 \end{cases}$$

Using this OPT function, we can solve this question in $O(n^2)$.

To accerlerate this process, we can reconsider the following equation

$$dp_i = \min_{j=0}^{i-1} \{dp_j + (m - (pre_i - pre_j - 1))^2\}$$

$$dp_i \geq dp_j + (m - (pre_i - pre_j - 1))^2$$

$$dp_i + 2mpre_i - pre_i^2 - 2pre_i \geq dp_j + (pre_j - 1)^2 + 2mpre_j + m^2 + 2m - 2pre_i pre_j$$

Record the above inequality as the following form:

$$tar(i) \geq h(j) + f(i)g(j)$$

$$tar(i) = dp_i + 2mpre_i - pre_i^2 - 2pre_i$$

$$h(j) = dp_j + (pre_j - 1)^2 + 2mpre_j + m^2 + 2m$$

$$f(i) = -2pre_i$$

$$g(j) = pre_j$$

Then $tar(i)$ equals to the $\max_{j=0}^{i-1} \{h(j) + f(i)g(j)\}$, we take $k = h(j) + f(i)g(j)$, namely

$$h(j) = -f(i)g(j) + k$$

After find the $tar(i)$, through

$$dp_i = tar(i) - 2mpre_i + pre_i^2 + 2pre_i$$

we can easily get the dp_i .

So we only need to record the former point $(g(j), h(j))$ and using the slope of value $-f(i)$ to find the point whose intercept is the minimum. For i , we've maintained the previous $i - 1$ point and find a convex pack for this points, to find the best point we only need to binary the points in the convex pack and using the $-f(i)$ as the search basis. Thus we can update dp_i in $O(\log n)$. **It's noteworthy that the $-f(i) = 2pre_i$ which is positive, so in our pseudo-code, we only need to main the points laying at the down part of convex pack.**

Above all, we maintained our convex pack in $O(n \log n)$ and do $n - 1$ times binary search with $O(\log n)$ each, so the time complexity is

$$O(n \log n + (n - 1) \log n)$$

It's

$$O(n \log n)$$

4. **(DP & Binary Search)** Let $A = [a_1, a_2, \dots, a_n]$, $a_i \in \mathbb{Z}$. A sub-sequence of A can be described as $A' = [a_{i_1}, a_{i_2}, \dots, a_{i_k}]$, where $1 \leq i_1 < i_2 < \dots < i_k \leq n$. An increasing sub-sequence of A satisfies $a_{i_1} < a_{i_2} < \dots < a_{i_k}$. Find the longest increasing sub-sequence in worst case $O(n \log n)$ time complexity. Give your algorithm in the form of pseudo-code and analyze its time complexity.

Algorithm 3: LIS problem

Input: $A[a_1, a_2, \dots, a_n]$
Output: The longest increasing sub-sequence in $A[\cdot]$

```

1   $dp[1, 2, \dots, n] \leftarrow \{\infty, -1\}$  foreach in  $dp[\cdot]$ ;
2   $pre[1, 2, \dots, n] \leftarrow -1$  foreach in  $pre[\cdot]$ ;
3  for  $i \leftarrow 1$  to  $n$  do
4       $l \leftarrow 1$ ;  $r \leftarrow n$ ;
5      while  $l < r$  do
6           $m \leftarrow \frac{l+r}{2}$ ;
7          if  $dp_m.first < a_i$  then
8               $l \leftarrow m + 1$ ;
9          else
10              $r \leftarrow m$ 
11      $dp_i \leftarrow \{a_i, i\}$ ;
12      $pre_i \leftarrow dp_{l-1}.second$ ;
13  $ans \leftarrow 1$ ;
14 while  $dp_{ans}.first \neq \infty$  do
15      $ans \leftarrow ans + 1$ ;
16  $ans \leftarrow ans - 1$ ;
17  $solution \leftarrow$  empty deque;
18  $cur \leftarrow dp_{ans}.second$ ;
19 while  $cur \neq -1$  do
20      $solution.push\_front(a_{cur})$ ;
21      $cur \leftarrow pre_{cur}$ ;
22 return  $solution$ ;
```

Here comes the OPT relationship

$$OPT[i] = \begin{cases} \max\{OPT[j] + 1\} & i \neq 1 \text{ or } 0 \\ 1 & i = 1 \\ 0 & i = 0 \end{cases}$$

***j in 1 to i - 1 for every j that $a[j] < a[i]$. If j doesn't exist, $j = 0$.**

If we do this by enumerating, there will be $O(n)$ iterations and in each iteration, the enumerating times will be $O(n)$. If so, the time complexity will be $O(n^2)$.

But by using binary search to accelerate the times of enumerate operation, we get $O(\log n)$ operation in each iteration.

As for why the binary search can be used, we know that the optimal only have something to do with the sequence that is longest and its last element is smaller than the current one. So it can be used.

The overall the time complexity is

$$O(n \log n)$$