



南方科技大学
SOUTHERN UNIVERSITY OF SCIENCE AND TECHNOLOGY

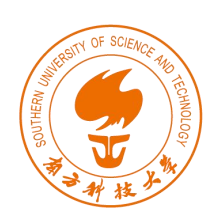
Algorithm Design and Analysis (H)

CS216

Prof. Shiqi Yu (于仕琪)

yusq@sustech.edu.cn

<http://faculty.sustech.edu.cn/yusq/>



南方科技大学
SOUTHERN UNIVERSITY OF SCIENCE AND TECHNOLOGY

Dynamic Programming



南方科技大学
SOUTHERN UNIVERSITY OF SCIENCE AND TECHNOLOGY

5. Sequence Alignment



Correct Typos Automatically

发件人: Shiqi Yu(于仕琪) – yusq@sustech.edu.cn

邮件大小: 24 KB

Dear [redacted],

Glad to knoww



Prof. Shiqi Yu

Department of
Southern Univ
Shenzhen, Ch
<https://faculty.sustech.edu.cn/yusq/>

know
knows
known

忽略拼写
学习拼写

查询“knoww”
翻译“knoww”
用 Google 搜索

neering,
logy,



The course Agorithm Design is the best one!↵

Algorithm

Algorithms

全部忽略

添加到词典

添加到自动更正





It can do more ...

- Gmail

« ▾ iapr@papercept.net, yusq@sustech.edu.cn

Thank you for your email



« ▾ iapr@papercept.net, yusq@sustech.edu.cn

Thank you for your kind response



↶ ↷ Sans Serif ▾ ↕ ▴ ▾ **B** *I* U A ▾ ≡ ▾

Send ▾

A



↶ ↷ Sans Serif ▾ ↕ ▴ ▾ **B** *I* U A ▾ ≡ ▾ 1.
2.
3.

Send ▾

A





String Similarity

- How similar are two strings?

- **ocurrence**
- **occurrence**

o	c	u	r	r	a	n	c	e	-
---	---	---	---	---	---	---	---	---	---

o	c	c	u	r	r	e	n	c	e
---	---	---	---	---	---	---	---	---	---

6 mismatches, 1 gap

o	c	-	u	r	r	a	n	c	e
---	---	---	---	---	---	---	---	---	---

o	c	c	u	r	r	e	n	c	e
---	---	---	---	---	---	---	---	---	---

1 mismatch, 1 gap

o	c	-	u	r	r	-	a	n	c	e
---	---	---	---	---	---	---	---	---	---	---

o	c	c	u	r	r	e	-	n	c	e
---	---	---	---	---	---	---	---	---	---	---

0 mismatches, 3 gaps



Edit Distance

- Applications.

- Basis for Unix diff.
- Speech recognition.
- Computational biology.

- Edit distance. [Levenshtein 1966, Needleman-Wunsch 1970]

- Gap penalty δ ; mismatch penalty α_{pq} .
- Cost = sum of gap and mismatch penalties.

C T G A C C T A C C T

- C T G A C C T A C C T

C C T G A C T A C A T

C C T G A C - T A C A T

$$\alpha_{TC} + \alpha_{GT} + \alpha_{AG} + 2\alpha_{CA}$$

$$2\delta + \alpha_{CA}$$



Sequence Alignment

- Goal: Given two strings $X = x_1 x_2 \dots x_m$ and $Y = y_1 y_2 \dots y_n$ find alignment of minimum cost.
- Def. An **alignment** M is a set of ordered pairs x_i-y_j such that each item occurs in at most one pair and no crossings.
- Def. The pair x_i-y_j and $x_{i'}-y_{j'}$ **cross** if $i < i'$, but $j > j'$.
- Ex: CTACCG **vs.** TACATG.
Sol: $M = x_2-y_1, x_3-y_2, x_4-y_3, x_5-y_4, x_6-y_6$.

x_1	x_2	x_3	x_4	x_5		x_6
C	T	A	C	C	-	G
-	T	A	C	A	T	G
	y_1	y_2	y_3	y_4	y_5	y_6



Sequence Alignment: Problem Structure

•Def. $\text{OPT}(i, j)$ = min cost of aligning strings $x_1 x_2 \dots x_i$ and $y_1 y_2 \dots y_j$.

- Case 1: OPT matches x_i - y_j .
 - ✓ pay mismatch for x_i - y_j + min cost of aligning two strings $x_1 x_2 \dots x_{i-1}$ and $y_1 y_2 \dots y_{j-1}$
- Case 2a: OPT leaves x_i unmatched.
 - ✓ pay gap for x_i and min cost of aligning $x_1 x_2 \dots x_{i-1}$ and $y_1 y_2 \dots y_j$
- Case 2b: OPT leaves y_j unmatched.
 - ✓ pay gap for y_j and min cost of aligning $x_1 x_2 \dots x_i$ and $y_1 y_2 \dots y_{j-1}$

$$\text{OPT}(i, j) = \begin{cases} j\delta, & \text{if } i = 0 \\ \min \begin{cases} \alpha_{x_i y_j} + \text{OPT}(i-1, j-1) \\ \delta + \text{OPT}(i-1, j) \\ \delta + \text{OPT}(i, j-1) \end{cases}, & \text{otherwise} \\ i\delta, & \text{if } j = 0 \end{cases}$$



Sequence Alignment: Algorithm

```
Sequence-Alignment(m, n,  $x_1x_2\dots x_m$ ,  $y_1y_2\dots y_n$ ,  $\delta$ ,  $\alpha$ ) {  
  for i = 0 to m  
    M[i, 0] =  $i\delta$   
  for j = 0 to n  
    M[0, j] =  $j\delta$   
  
  for i = 1 to m  
    for j = 1 to n  
      M[i, j] = min( $\alpha[x_i, y_j] + M[i-1, j-1]$ ,  
                    $\delta + M[i-1, j]$ ,  
                    $\delta + M[i, j-1]$ )  
  
  return M[m, n]  
}
```

$$\text{OPT}(i, j) = \begin{cases} j\delta, & \text{if } i = 0 \\ \min \begin{cases} \alpha_{x_i y_j} + \text{OPT}(i-1, j-1) \\ \delta + \text{OPT}(i-1, j) \\ \delta + \text{OPT}(i, j-1) \end{cases}, & \text{otherwise} \\ i\delta, & \text{if } j = 0 \end{cases}$$

- Analysis. $\Theta(mn)$ time and space.
- English words or sentences: $m, n \leq 10$.
- Computational biology: $m = n = 100,000$. 10 billions ops OK, but 10GB array?



南方科技大学
SOUTHERN UNIVERSITY OF SCIENCE AND TECHNOLOGY

6. Shortest Paths

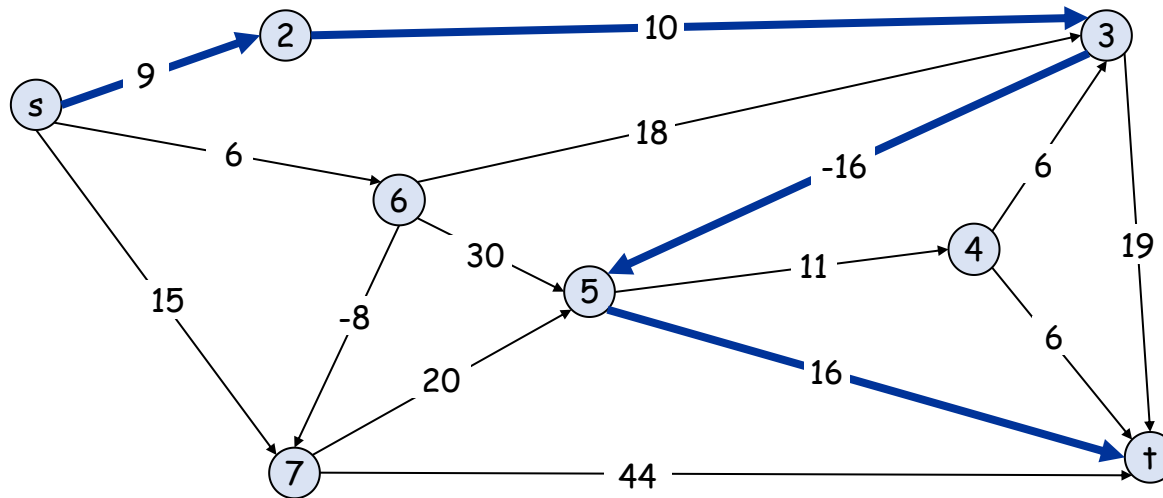


Shortest Paths

- Shortest path problem. Given a directed graph $G = (V, E)$, with edge weights c_{vw} , find shortest path from node s to node t .

↖ allow negative weights

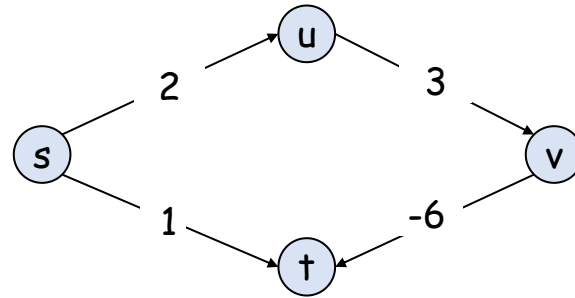
- Ex. Nodes represent agents in a financial setting and c_{vw} is cost of transaction in which we buy from agent v and sell immediately to w .



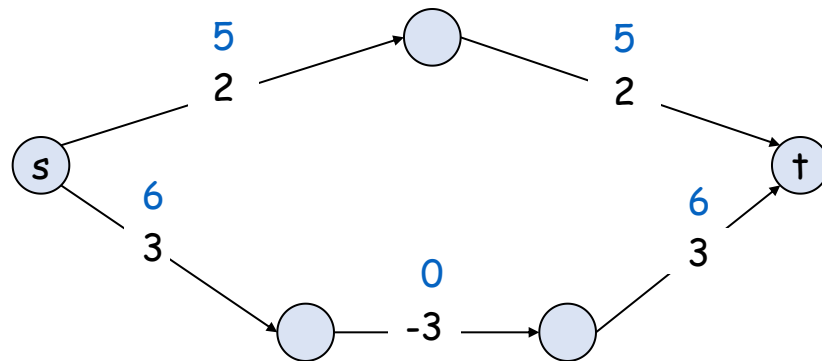


Shortest Paths: Failed Attempts

- Dijkstra. Can fail if negative edge costs.



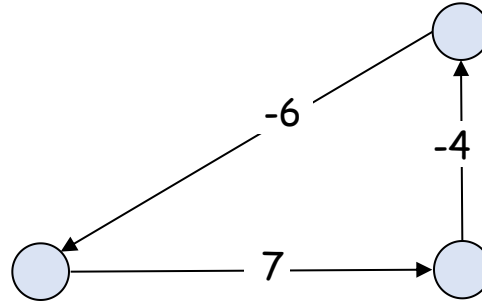
- Re-weighting. Adding a constant to every edge weight can fail.



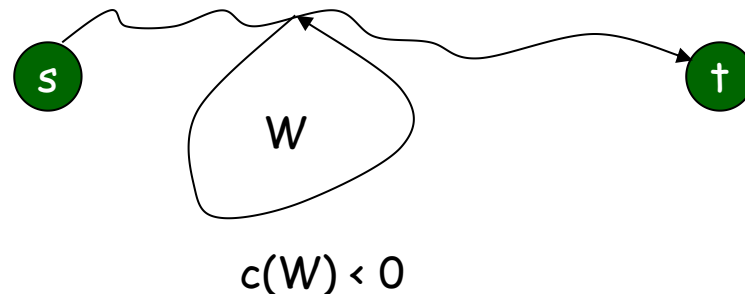


Shortest Paths: Negative Cost Cycles

- Negative cost cycle.



- Observation. If some path from s to t contains a negative cost cycle, there does not exist a shortest s - t path; otherwise, there exists one that is simple.





Shortest Paths: Dynamic Programming

- Def. $OPT(i, v)$ = length of shortest v - t path P using at most i edges.
 - Case 1: P uses at most $i-1$ edges.
 - ✓ $OPT(i, v) = OPT(i-1, v)$
 - Case 2: P uses exactly i edges.
 - ✓ if (v, w) is first edge, then OPT uses (v, w) , and then selects best w - t path using at most $i-1$ edges

$$OPT(i, v) = \min(OPT(i-1, v), \min_{w \in V} (OPT(i-1, w) + C_{vw}))$$

- Remark. By previous observation, if no negative cycles, then $OPT(n-1, v)$ = length of shortest v - t path.



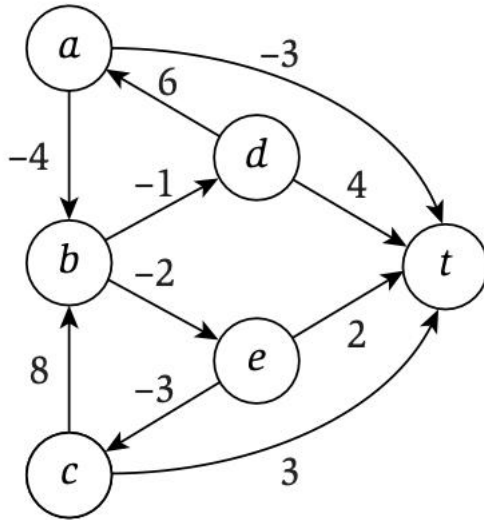
Shortest Paths: Implementation

```
Shortest-Path( $G, t$ ) {  
    foreach node  $v \in V$   
         $M[0, v] \leftarrow \infty$   
     $M[0, t] \leftarrow 0$   
  
    for  $i = 1$  to  $n-1$   
        foreach node  $v \in V$   
             $M[i, v] \leftarrow M[i-1, v]$   
            foreach edge  $(v, w) \in E$   
                 $M[i, v] \leftarrow \min \{ M[i, v], M[i-1, w] + c_{vw} \}$   
}
```

- Analysis. $\Theta(mn)$ time, $\Theta(n^2)$ space.
- Finding the shortest paths. Maintain a "successor" for each table entry.



Shortest Paths: Example



(a)

	0	1	2	3	4	5
t	0	0	0	0	0	0
a	∞	-3	-3	-4	-6	-6
b	∞	∞	0	-2	-2	-2
c	∞	3	3	3	3	3
d	∞	4	3	3	2	0
e	∞	2	0	0	0	0

(b)

$$OPT(i, v) = \min(OPT(i-1, v), \min_{w \in V} (OPT(i-1, w) + C_{vw}))$$



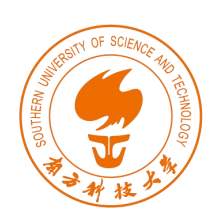
Shortest Paths: Practical Improvements

- Practical improvements.
 - Maintain only one array $M[v]$ = shortest v-t path that we have found so far.
 - No need to check edges of the form (v, w) unless $M[w]$ changed in previous iteration.
- Theorem. Throughout the algorithm, $M[v]$ is length of some v-t path, and after i rounds of updates, the value $M[v]$ is no larger than the length of shortest v-t path using $\leq i$ edges.
- Overall impact.
 - Memory: $O(m + n)$.
 - Running time: $O(mn)$ worst case, but substantially faster in practice.



Bellman-Ford: Efficient Implementation

```
Push-Based-Shortest-Path( $G, s, t$ ) {  
    foreach node  $v \in V$  {  
         $M[v] \leftarrow \infty$   
         $\text{successor}[v] \leftarrow \phi$   
    }  
  
     $M[t] = 0$   
    for  $i = 1$  to  $n-1$  {  
        foreach node  $w \in V$  {  
            if ( $M[w]$  has been updated in previous iteration) {  
                foreach node  $v$  such that  $(v, w) \in E$  {  
                    if ( $M[v] > M[w] + c_{vw}$ ) {  
                         $M[v] \leftarrow M[w] + c_{vw}$   
                         $\text{successor}[v] \leftarrow w$   
                    }  
                }  
            }  
        }  
        If no  $M[w]$  value changed in iteration  $i$ , stop.  
    }  
}
```



南方科技大学
SOUTHERN UNIVERSITY OF SCIENCE AND TECHNOLOGY

7. Distance Vector Protocol



Distance Vector Protocol

- Communication network.
 - Node \approx router.
 - Edge \approx direct communication link. \leftarrow naturally nonnegative, but Bellman-Ford used anyway!
 - Cost of edge \approx delay on link.
- Dijkstra's algorithm. Requires global information of network.
- Bellman-Ford. Uses only local knowledge of neighboring nodes.
- Synchronization. We don't expect routers to run in lockstep. The order in which each `foreach` loop executes is not important. Moreover, algorithm still converges even if updates are asynchronous.



Push-based shortest path

Push-Based-Shortest-Path(G, s, t)

n = number of nodes in G

Array $M[V]$

Initialize $M[t]=0$ and $M[v]=\infty$ for all other $v \in V$

For $i=1, \dots, n-1$

 For $w \in V$ in any order

 If $M[w]$ has been updated in the previous iteration then

 For all edges (v, w) in any order

$M[v] = \min(M[v], c_{vw} + M[w])$

 If this changes the value of $M[v]$, then $first[v]=w$

 Endfor

 Endfor

 If no value changed in this iteration, then end the algorithm

 Endfor

Return $M[s]$



Asynchronous shortest path

Asynchronous-Shortest-Path(G, s, t)

n = number of nodes in G

Array $M[V]$

Initialize $M[t] = 0$ and $M[v] = \infty$ for all other $v \in V$

Declare t to be active and all other nodes inactive

While there exists an active node

 Choose an active node w

 For all edges (v, w) in any order

$M[v] = \min(M[v], c_{vw} + M[w])$

 If this changes the value of $M[v]$, then

$first[v] = w$

v becomes active

 Endfor

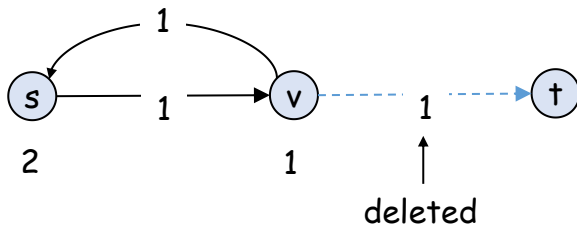
w becomes inactive

EndWhile



Distance Vector Protocol

- Distance vector protocol.
 - Each router maintains a vector of shortest path lengths to every other node (distances) and the first hop on each path (directions).
 - Algorithm: each router performs n separate computations, one for each potential destination node.
 - "Routing by rumor."
- Ex. RIP, Xerox XNS RIP, Novell's IPX RIP, Cisco's IGRP, DEC's DNA Phase IV, AppleTalk's RTMP.
- Caveat. Edge costs may **change** during algorithm (or fail completely).



"counting to infinity"



Path Vector Protocols

- Link state routing.
 - Each router also stores the entire path.
 - Based on Dijkstra's algorithm.
 - Avoids "counting-to-infinity" problem and related difficulties.
 - Requires significantly more storage.
- Ex. Border Gateway Protocol (BGP), Open Shortest Path First (OSPF).

not just the distance and first hop
↙