

# Семинар 14

Динамическое программирование (продолжение)

# Задача о нахождении расстояния редактирования

Текстовый редактор, подозревая опечатку, предлагает заменить написанное слово на *близкое* (?).

Каково расстояние между словами SNOWY и SUNNY?

S	—	N	O	W	Y	или	—	S	N	O	W	—	Y
S	U	N	N	—	Y		S	U	N	—	—	N	Y

Стоимость выравнивания: количество столбцов, в которых символы различаются. Расстояние редактирования между словами – стоимость их наилучшего выравнивания.

$$\underbrace{\begin{array}{l} \text{—} X + \text{зазоры} \text{—} \\ \text{—} Y + \text{зазоры} \text{—} \end{array}}_{\text{общая длина } \ell}$$

- пусть  $A$  - конечный алфавит.  $X = x_1x_2 \dots x_m$ ,  $Y = y_1y_2 \dots y_n$ ,

$x_i, y_j \in A$ . Сколько существует возможностей для содержимого конечного столбца оптимального выравнивания?

- а) 2
- б) 3
- в) 4
- г)  $mn$

3 варианта

1)

$x_m$

$y_n$

2)

$x_m$

-

3)

-

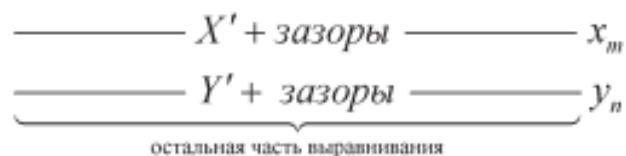
$y_n$

# Вывод рекуррентного соотношения

1)

$x_m$

$y_n$



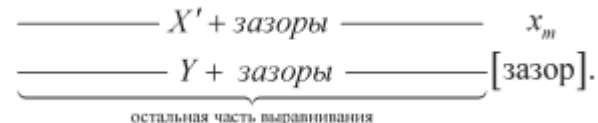
$P = P' + \alpha$ ,  $P'$  - оптимальное выравнивание  
для  $X'$  и  $Y'$

$$\alpha = \begin{cases} 0, & x_m = y_n \\ 1, & x_m \neq y_n \end{cases}$$

2)

$x_m$

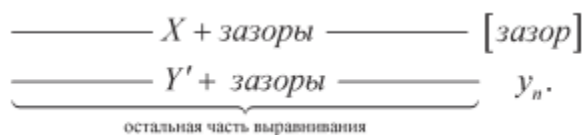
—



3)

—

$y_n$



. Оптимальное выравнивание двух непустых символьных цепочек  
 $X = x_1, x_2, \dots, x_m$  и  $Y = y_1, y_2, \dots, y_n$  равно либо:

- (i) оптимальному выравниванию  $X'$  и  $Y'$ , дополненному сочетанием  $x_m$  и  $y_n$  в последнем столбце;
- (ii) оптимальному выравниванию  $X'$  и  $Y$ , дополненному сочетанием  $x_m$  и зазора в конечном столбце;
- (iii) оптимальному выравниванию  $X$  и  $Y'$ , дополненному сочетанием зазора и  $y_n$  в конечном столбце,

где  $X'$  и  $Y'$  обозначают соответственно  $X$  и  $Y$  с удаленными конечными символами  $x_m$  и  $y_n$ .

Если одна из цепочек пуста (например,  $Y$ ),  
то расстояние редактирования равно  $X$

## Рекуррентное соотношение

- $P_{i,j} = \min\{P_{i-1,j-1} + \alpha, P_{i-1,j} + 1, P_{i,j-1} + 1\}, i = 1, 2, \dots, m, j = 1, 2, \dots, n$
- Подзадачи: вычислить минимальное расстояние редактирования первых  $i$  символов  $X$  и первых  $j$  символов  $Y$ .

Dist Red

Вход: цепочки  $X = x_1x_2 \dots x_m, Y = y_1y_2 \dots y_n$  над алфавитом  $A$

Выход:  $P$  – расстояние редактирования

//решения подзадач(индексируемых с 0)

$A := (m + 1)(n + 1)$  двумерный массив

// базовый случай 1 ( $j = 0$ )

For  $i := 0$  to  $m$  do  $A[i][0] = i$

// базовый случай 2 ( $i = 0$ )

For  $j := 0$  to  $n$  do

$A[0][j] = j$

//систематическое решение всех подзадач

For  $i := 0$  to  $m$  do

For  $j := 0$  to  $n$  do

$A[i][j] := \min\{A[i-1][j-1] + \alpha, A[i-1][j] + 1, A[i][j-1] + 1\}$

Return  $A[m][n]$  //решение самой крупной подзадачи

**$O(m \cdot n)$**

# Пример вычисления расстояния редактирования

Вычислить расстояние редактирования между словами «МУСОР» и «ССОРА»

а	5	5	5	5	4	3
р	4	4	4	4	3	2
о	3	3	3	3	2	3
с	2	2	2	2	3	4
с	1	1	2	2	3	4
	0	1	2	3	4	5
		м	у	с	о	р

—	С	С	О	Р	А
М	У	С	О	Р	—

Алгоритм реконструкции (обратный проход)

↓ Зазор в слове «МУСОР»

← Зазор в слове «ССОРА»

↖ Сверху и снизу нет зазора

# Выравнивание последовательностей

Сравниваем 2 участка одного или нескольких геномов.

Алфавит  $A=\{A, C, G, T\}$  A: [Аденин](#); G — Г: [Гуанин](#); C — Ц: [Цитозин](#); T — Т: [Тимин](#)

Допустим, что за несовпадение символов штраф=2, за отсутствие (зазор) штраф=1

Сравним ACG и AGCT

G	3	2	1	2	3
C	2	1	2	1	2
A	1	0	1	2	3
	0	1	2	3	4
		A	G	C	T

A G C T  
A - C G

или

A G C T  
A - C G -

или

A G C T -  
A - C - G

Минимальный штраф за выравнивание (=3) называется отметкой Нидлмана-Вунша (NW-отметка символьных цепочек)

# Кратчайшие пути в ориентированном графе с отрицательными длинами ребер (повторно)

- Задача: кратчайшие пути с единственным истоком

**Вход:** ориентированный граф  $G = (V, E)$ , истоковая вершина  $s \in V$  и вещественная длина  $\ell_e$  для каждого ребра  $e \in E$ .

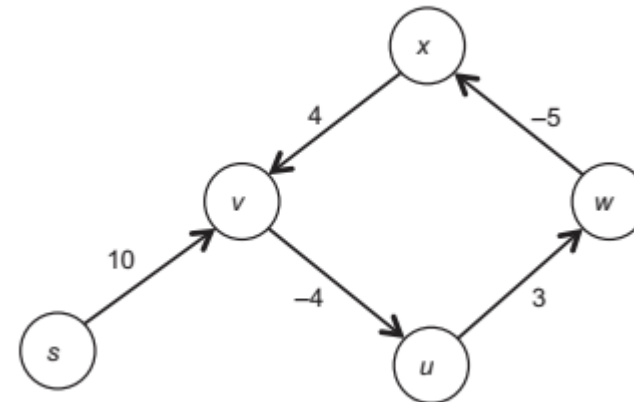
**Выход:** один из следующих:

- (i) расстояние кратчайшего пути  $dist(s, v)$  для каждой вершины  $v \in V$  либо
- (ii) заявление о том, что  $G$  содержит отрицательный цикл.

---

Рассмотрим экземпляр задачи о кратчайшем пути с единственным истоком с  $n$  вершинами,  $m$  ребрами, стартовой вершиной  $s$  и без отрицательных циклов. Что из перечисленного является истинным? Выберите самое подходящее утверждение:

- а) Для каждой вершины  $v$ , достижимой из истока  $s$ , существует кратчайший путь  $s-v$  не более чем с  $n - 1$  ребрами;
- б) Для каждой вершины  $v$ , достижимой из истока  $s$ , существует кратчайший путь  $s-v$  не более чем с  $m$  ребрами;
- г) Нет конечной верхней границы (как функции от  $n$  и  $m$ ) на наименьшем числе ребер в кратчайшем пути  $s-v$ .

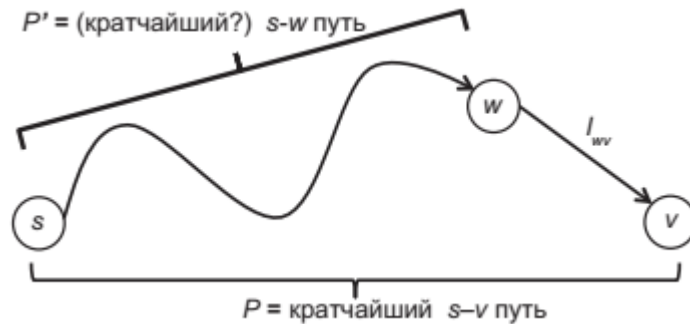




# Алгоритм Беллмана-Форда с использованием стратегии динамического программирования

- Если задача решена и  $P$  – кратчайший путь  $s \rightarrow v$ , содержащий не более, чем  $i$  ребер, то

$$L_{i,v} = \min \left\{ \begin{array}{l} L_{i-1,v} \\ \min_{(w,v) \in E} \{ L_{i-1,w} + \ell_{wv} \} \end{array} \right. \quad *$$



\_\_\_\_\_

**Вход:** ориентированный граф  $G = (V, E)$ , представленный в виде списков смежности, истоковая вершина  $s \in V$  и вещественная длина  $\ell_e$  для каждого  $e \in E$ .

---

$A := (n + 1) \times n$  двумерный массив

$$A[0][s] := 0$$
$$A[0][v] := +\infty$$

```
for  $i = 1$  to  $n$  do // размер подзадачи
```

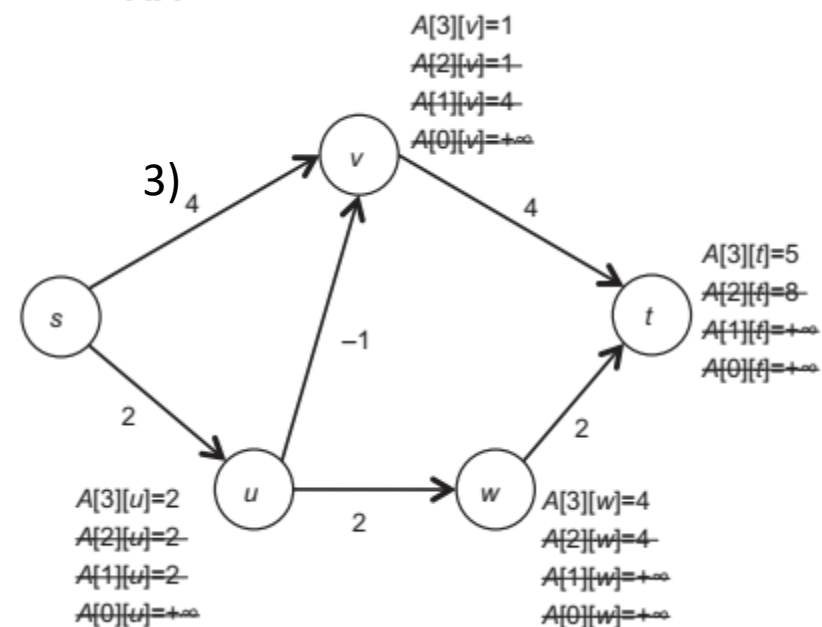
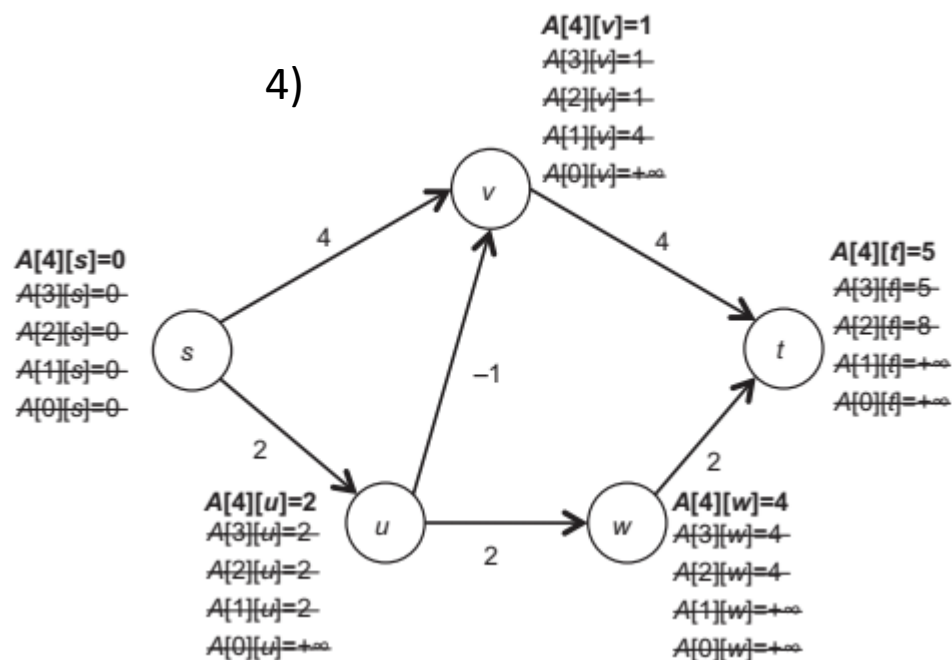
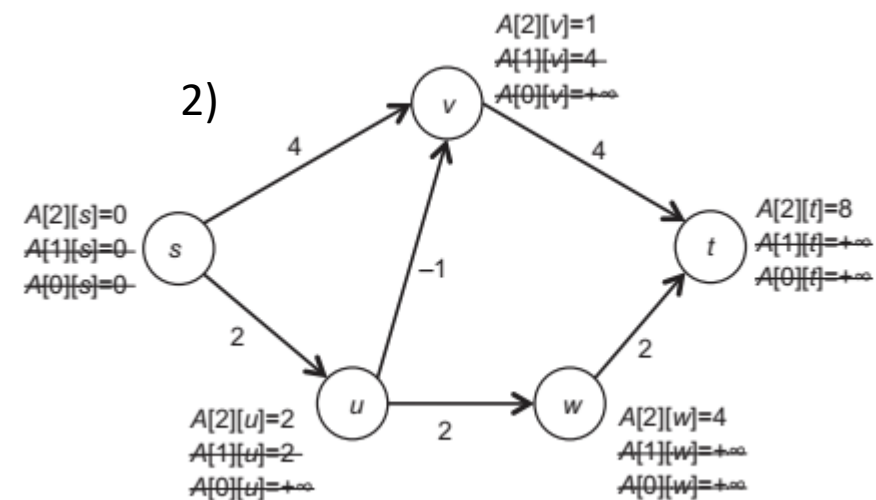
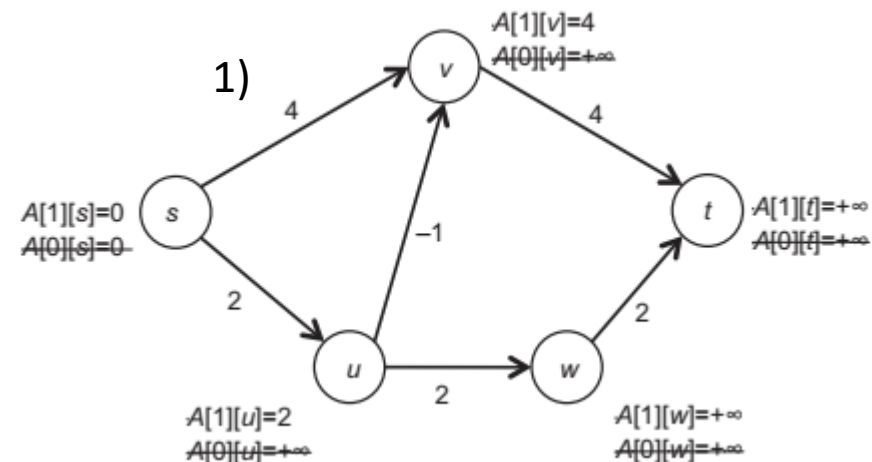
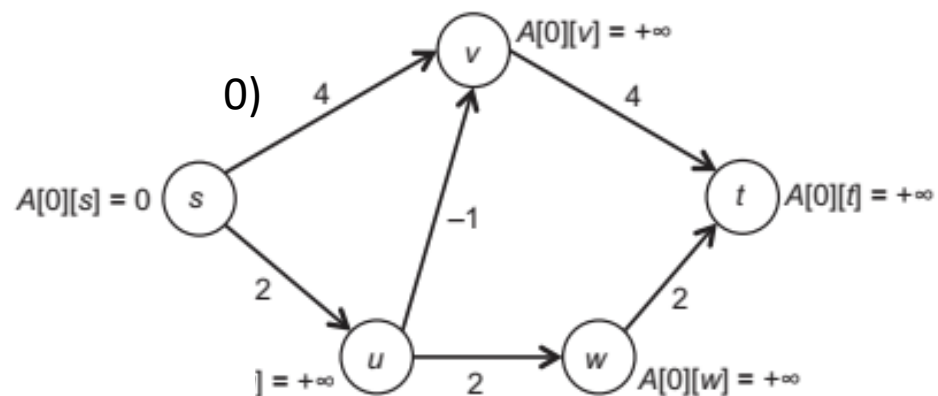
**for**  $v \in V$  **do**
$$A[i][v] :=$$

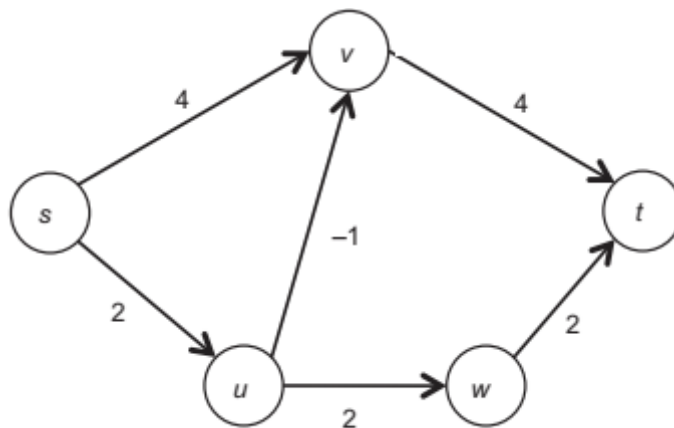
**if**  $A[i][v] \neq A[i-1][v]$  **then**

```
if stable = TRUE then // выполнено леммой 18.3
```

```
// не удалось стабилизироваться на n итерациях
```

# Пример





$v \backslash i$	0	1	2	3	4
s	0	0	0	0	0
v	$\infty$	4	1	1	1
u	$\infty$	2	2	2	2
w	$\infty$	$\infty$	4	4	4
t	$\infty$	$\infty$	8	5	5



Не меняется

$O(mn)$