Повторы в текстах

Слова и словосочетания.

Повтор предложения (абзаца) – признак ошибки.

Мелодические обороты.

В ДНК-последовательностях:

- 1. Участки с аномальным нуклеотидным составом: (А,Т)-богатые...
- 2.Участки микросателлитной ДНК: периодичности с малой длиной периода (2-3) и достаточно высокой кратностью повторений.
- 3. Тандемные повторы с произвольной длиной периода.
- 4. Разнесенные повторы значительной длины.
- 5. Мобильные элементы

Повтор — пара совпадающих фрагментов текста.

Классификация повторов

•По способу их прочтения и использованию переименований:

Повтор (в широком смысле) — пара фрагментов текста, совпадающих с точностью до переименования элементов алфавита и (или) изменения направления считывания.

Типы повторов:

- прямые : ... AGTTC ... AGTTC...– симметричные: ... AGTTC ... CTTGA...
- с точностью до подстановки на элементах алфавита: секвентные переносы в музыке; замены $0 \leftrightarrow 1$; $A \leftrightarrow T$, $C \leftrightarrow G$.
 - прямые комплементарные: ... AGTTC... TCAAG...
 - инвертированные: ... AGTTC... GAACT...

Классификация повторов

· По наличию искажений:

Повторы могут быть

совершенные:

... AGTTC ... AGTTC...

```
и несовершенные (с заменами, вставками, делециями): ... AGTTC ... AATTC ... (замена), ... AGTTC ... AGTTTC... (вставка), в том числе с точностью до агрегирования: ... AGTTC ... GATCT ... (совпадают при заменах \{A,G\} \rightarrow Pu, \{C,T\} \rightarrow Py).
```

Классификация повторов

• По характеру расположения в тексте:

будем различать повторы

разнесенные

... AGTTC ... AGTTC...

тандемные ... AGTTC AGTTC...

с наложением: ... AGTTCAGTTC AGTTC ...

Представление текста в терминах повторов

Полный частотный спектр текста.

Пусть

∑ – конечный алфавит;

S – текст, составленный из элементов Σ ;

N = |S| — длина текста;

 $S[i] = S_i - i$ -й элемент текста $S(1 \le i \le N)$;

S[i:j] — фрагмент текста, включающий элементы с i-го по j-й $(1 \le i < j \le N)$.

l-грамма — связная цепочка текста длины *l* (S[i:i+l-1]).

$$S = s_1 s_2 s_3 s_4 s_5 \dots s_N$$

Полное число *l*-грамм: N-l+1.

Число различных *l*-грамм: $M_l \le N - l + 1$.

Частомная характеристика *l***-го порядка** текста S — совокупность элементов $\Phi_l(S) = \{ \varphi_{l1}, \varphi_{l2}, ..., \varphi_{lMl} \}$ где φ_{li} ($1 \le i \le M_l$) есть пара : « i-я l-грамма — x_i , ее частота в тексте — $F_l(x_i)$ ».

 $l_{\max}(S)$ — наибольшее l, при котором в S есть *повторяющиеся* l-граммы.

Совокупность частотных характеристик

$$\Phi(S) = \{\Phi_1(S), \Phi_2(S), ..., \Phi_{l \max + 2}(S)\}$$

называется полным частотным спектром текста S.

Усеченный спектр

$$\Phi^*(S) = \{\Phi^*_1(S), \Phi^*_2(S), ..., \Phi^*_{l max}(S)\}$$

 $\Phi^*_{l}(S)$ отличается от $\Phi_{l}(S)$ отсутствием l-грамм с единичной частотой встречаемости.

Пример. Пусть S = caabcabbca.

Тогда
$$\Phi *(S) = {\Phi *_{1}(S), \Phi *_{2}(S), \Phi *_{3}(S)},$$
 где

$$\Phi_{1}(S) = \{ \langle a, F_{1}(a) = 4 \rangle; \quad \langle b, F_{1}(b) = 3 \rangle; \quad \langle c, F_{1}(c) = 3 \rangle \};$$

$$\Phi_2^*(S) = \{ \langle ca, F_2(ca) = 3 \rangle; \langle ab, F_2(ab) = 2 \rangle; \langle bc, F_2(bc) = 2 \rangle \};$$

$$\Phi_{3}^{*}(S) = { };$$

Для повторов значительной длины спектр можно дополнить позиционной информацией.

Наиболее важными являются следующие параметры частотного спектра:

 l_{max} — длина максимального повтора в тексте.

Для случайного текста длины N с вероятностями появления элементов алфавита p_r $(1 \le r \le n = |\Sigma|)$ можно пользоваться оценкой: .

Если реальная длина l_{max} в тексте существеннен ревышает ожидаемое значение, это свидетельствует о наличии дупликативных механизмов порождения текста.

 M_l — размер словаря l-грамм.

Он фигурирует в определении комбинаторной сложности текста.

— максимальное значение частот встречаемости

$$F_l^{\max} = \max_{1 \le i \le M_l} F_l(x_i)$$
 I-грамм в тексте $(1 \le l \le l_{\max});$

— — минимальное значение частот встречаемости

 $F_l^{\min} = \min_{1 \le i \le M_l} F_l(x_i)$ — грамм в тексте $(1 \le l \le l_{\max})$; представляет интерес при малых значениях l; при больших, обычно $F_l^{\min} = 1$.

 E_l^k — Число различных l-грамм, каждая из которых встречается в тексте ровно k раз (k=1,2,...,N-l+1); зависимость E_l^k от k при фиксированном l является аналогом известной в лингвистике кривой Юла;

 $M_l - E_l^{-1}$ — число различных повторяющихся l-грамм; E_l^{0} — число l-грамм, ни разу не встретившихся в тексте. Наличие таковых в ситуации, когда $N/n^l>>1$, можно трактовать как аномальный эффект.

Имеют место простые соотношения, связывающие основные параметры:

$$M_{l} = \sum_{k=1}^{F_{l}^{\max}} E_{l}^{k}, \qquad N-l+1 = \sum_{k=1}^{F_{l}^{\max}} k \cdot E_{l}^{k}, \qquad E_{l}^{0} = n^{l} - M_{l}.$$

Алгоритмы отыскания совершенных повторов

Метод лексикографической сортировки

```
Лексикографический порядок слов u = u_1 \dots u_p и v = v_1 \dots v_q: u \le v, если u = v или (\exists j, такое что u_j < v_j и u_i = v_i \ \forall i < j) или (p < q и для всех i \le p u_i = v_i).
Пример. S = abcdabcbcbabcd; l = 3;
                  f(abc) = 3 аналог для произвольного l -
abc
                               суффиксный массив
abc
abc
bab
bcb
                  f(bcb) = 2
bcb
                   f(bcd) = 2
bcd
bcd
cba
cbc
cda
dab
```

Задача:

Дан *массив* **A** из n *элементов*: $a_1, a_2, ... a_n$

С каждым элементом a_i связан *ключ* $k_i \in K$.

На множестве ключей К задано отношение порядка — **линейного** (или **совершенного**) **упорядочивания**, для которого выполнены следующие условия:

закон трихотомии: для любых $x,y \in K$ либо x < y, либо x > y, либо x = y; *транзитивность*: для любых x,y, $z \in K$ если x < y и y < z, то x < z.

Задачей сортировки по неубыванию является нахождение перестановки элементов $p(1), p(2), \dots p(n)$ при которой ключи располагаются в порядке неубывания: $\mathbf{k}_{p(1)} \leq \mathbf{k}_{p(2)} \leq \dots \leq \mathbf{k}_{p(n)}$.

В результате работы алгоритма и применения перестановки получается отсортированный массив: $a_{p(1)}, a_{p(2)}, \ldots, a_{p(n)}$

Аналогично можно определить сортировку по невозрастанию.

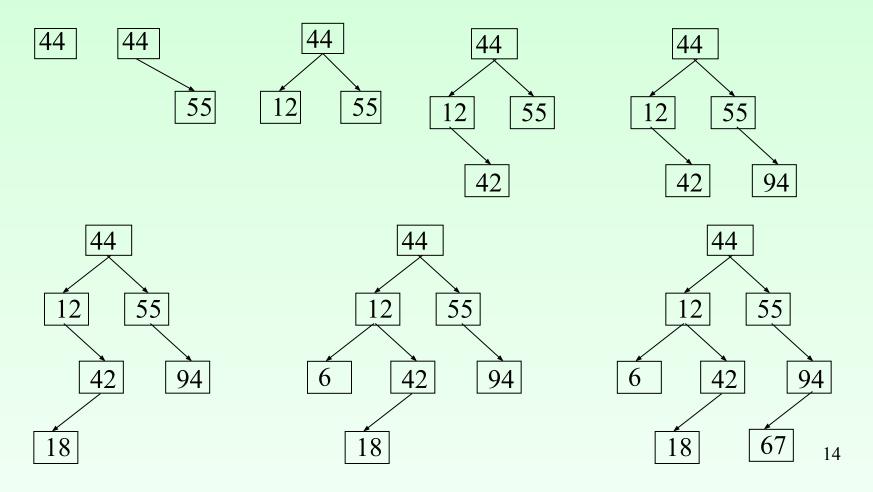
Сортировка выбором (Selection sort):

- •находим номер минимального значения в текущем списке
- •производим обмен этого значения со значением первой неотсортированной позиции (обмен не нужен, если минимальный элемент уже находится на данной позиции)
- •сортируем хвост списка, исключив из рассмотрения уже отсортированные элементы
- **5** 9 4 7 2 4 **1** 6
- 1 9 4 7 2 4 5 6
- 12479456
- 12479456
- 12449756
- 12445796
- 1 2 4 4 5 6 9 7
- 12445679

- •Сортировка пузырьком (сортировка всплыванием)
- •Элементы последовательно сравниваются попарно и, если порядок в паре неверный, выполняется обмен элементов. При каждом проходе алгоритма по внутреннему циклу, очередной наибольший элемент массива ставится на своё место в конце массива рядом с предыдущим «наибольшим элементом», а наименьший элемент перемещается на одну позицию к началу массива («всплывает»).

•Первый цикл:	второй цикл	третий цикл	четвертый цикл
•5 <u>9 4</u> 7 2 4 1 6	54 724169	4 <u>5 2</u> 4 1 6 7 9	<u>42</u> 415679
•5 4 <u>9 7</u> 2 4 1 6	4 <u>5 7</u> 2 4 1 6 9	4 2 <u>5 4</u> 1 6 7 9	2 <u>4 4 1</u> 5 6 7 9
•5 4 7 <u>9 2</u> 4 1 6	4 5 <u>7 2</u> 4 1 6 9	4 2 4 <u>5 1</u> 6 7 9	2 4 1 <u>4 5</u> 6 7 9
•5 4 7 2 <u>9 4</u> 1 6	4 5 2 <u>7 4</u> 1 6 9	4 2 4 1 <u>5 6</u> 7 9	пятый цикл
•5 4 7 2 4 <u>9 1</u> 6	4 5 2 4 <u>7 1</u> 6 9	2144567	9
•5 4 7 2 4 1 <u>9 6</u>	4 5 2 4 1 <u>7 6</u> 9	шестой цикл	
•5 4 7 2 4 1 6 9	4 5 2 4 1 6 <u>7 9</u>	1 2 4 4 5	6679

Сортировка деревом. При добавлении в дерево нового элемента его последовательно сравнивают с нижестоящими узлами. Если элемент >= корня - он идет в правое поддерево, сравниваем его уже с правым сыном, иначе - он идет в левое поддерево, сравниваем с левым, и так далее, пока есть сыновья, с которыми можно сравнить. **44 55 12 42 94 18 06 67**



Сортировка вычерпыванием:

Пусть известно, что максимальный элемент сортируемого массива не превосходит некоторое натуральное m:

- •Организовать m пустых очередей по одной для каждого натурального числа от 1 до m. Каждую такую очередь называют *черпаком*.
- •Просмотреть последовательность A слева направо, помещая элемент a_i в очередь с номером a_i
- •Сцепить эти очереди, т.е. содержимое (i+1)-й очереди приписать к концу ій очереди (i=1,2,...m-1), получив тем самым упорядоченную последовательность.

```
Пример: 1 3 2 5 2 5 1 2 1 5 4 3 4 5:11122233445555
```

- 1: 1 1 1
- 2: 2 2 2
- 3: 3 3
- 4: 4 4
- 5: 5 5 5 5

Лексикографическая сортировка на основе сортировки вычерпыванием

- l-граммы сортируем по позиции l
- Полученный список сортируем по позиции l-1
- ...
- Полученный список сортируем по позиции 1

Пример: abcdabcbcbabcd; *l* = 3; abc, bcd, cda, dab, abc, bcb, cbc, bcb, cba, bab, abc, bcd;

Сортировка по 3-й позиции:

a: cda, cba

b: dab, bcb, bcb, bab

c: abc, abc, cbc, abc

d: bcd, bcd

Список *l*-грамм после первого этапа: cda, cba, dab, bcb, bcb,bab, abc, abc, cbc, abc, bcd, bcd

Сортировка по 2-й позиции:

a: dab, bab

b: cba, abc, abc, cbc, abc

c: bcb, bcb, bcd, bcd

d: cda

Список l-грамм после второго этапа:

dab, bab, cba, abc, abc, cbc, abc, bcb,

bcb, bcd, bcd, cda

Сортировка по 1-й позиции:

a: abc, abc, abc

b: bab, bcb, bcd, bcd

c: cba, cbc, cda

d: dab

Список *l*-грамм после 3 этапа:

abc, abc, abc, bab, bcb, bcd, bcd, cba, cbc, cda, dab

Алгоритмы отыскания совершенных повторов

Метод, основанный на хешировании.

Хеширование (ассоциативная адресация) — отображение, которое ставит в соответствие l-грамме текста $x_i = T[i:i+l-1]$ $(1 \le i \le N-l+1)$ число $H(x_i)$ (адрес, по которому хранится информация об x_i).

Нумерующая функция

$$ns: \Sigma \to [0..|\Sigma| - 1]$$
 — порядок символов в Σ (s = $|\Sigma|$). $H(x_i) = ns(t_i) \times s^{l-1} + ns(t_{i+1}) \times s^{l-2} \dots ns(t_{i+l-2}) \times s + ns(t_{i+l-1})$. Рекуррентное хеширование:

$$H(x_{i+1}) = (H(x_i) - ns(t_i) \times s^{l-1}) \times s + ns(t_{i+1}).$$
Пример. $\Sigma = \{acgt\}, T = ggacataccaggac;$
 $H(T[1:4]) = 2 \times 4^3 + 2 \times 4^2 + 0 \times 4^1 + 1 = 161;$

 $H(T[2:5]) = 2 \times 4^3 + 0 \times 4^2 + 1 \times 4^1 + 0 = 132 = (161 - 2 \times 4^3) \times 4 + 0;$

 $H(T[3:6]) = 0 \times 4^3 + 1 \times 4^2 + 0 \times 4^1 + 3 = 19 = (132-2 \times 4^3) \times 4+3;$

...

$$H(T[11:14]) = 2 \times 4^3 + 2 \times 4^2 + 0 \times 4^1 + 1 = 161;$$

Недостаток этого отображения — большой (порядка $|\Sigma|^l$) диапазон изменения чисел $H(x_i)$ (сильно разреженный массив адресов).

Достоинство — отображение H взаимно-однозначное и достаточно просто вычислимо.

Пример функции расстановки с наложениями:

$$h_2(x_i) = H(x_i) \bmod M$$

 $h_2(x_i) = H(x_i) \bmod M$ — простое число (размер поля).

Пример списковой схемы устранения наложений

Информа-		Расстановочное	e	Дополнительное				
ционный		поле		поле (ДП)				
массив	No	I(x)	A(x)		No	I(x)	A(x)	
x_1	1	$I(x_4)$	*		1	$I(x_3)$	•	\Box
x_2	2				2	$I(x_5)$	*	
x_3	3	$I(x_1)$			3	$I(x_7)$	*	┢┚
x_4	4	$I(x_6)$	*		4			
x_5	5							
x_6	6	$I(x_2)$						
x_7	7	$I(x_8)$	*					
x_8	8							

Алгоритмы отыскания совершенных повторов

Хеширование. Пример.

Пример. S = abcdabcbcbabcd; $l = 3; h_1(x) = H(x)$ ns(a) = 0; ns(b) = 1; ns(c) = 2; ns(d) = 3 $h_1(abc) = 0*4^2 + 1*4^1 + 2 = 6;$ $h_1(bcd) = 1*4^2 + 2*4^1 + 3 = 27;$ $h_1(\text{cda}) = 2*4^2 + 3*4^1 + 0 = 44;$ $h_1(\text{dab}) = 3*4^2 + 0*4^1 + 1 = 49;$ $h_1(abc) = 0*4^2 + 1*4^1 + 2 = 6;$ $h_1(bcb) = 1*4^2 + 2*4^1 + 1 = 25;$ $h_1(cbc) = 2*4^2 + 1*4^1 + 2 = 38;$ $h_1(bcb) = 1*4^2 + 2*4^1 + 1 = 25;$ $h_1(\text{cba}) = 2*4^2 + 1*4^1 + 0 = 36;$ $h_1(bab) = 1*4^2 + 0*4^1 + 1 = 17;$ $h_1(abc) = 0*4^2 + 1*4^1 + 2 = 6;$ $h_1(bcd) = 1*4^2 + 2*4^1 + 3 = 27;$

- 0:
- •
- 6: abc, abc, abc
- 7:
- ---
- 17: bab
- ---
- 25: bcb, bcb
- 26:
- 27: bcd, bcd
- ---
- 36: cba
- 37:
- 38: cbc
- ---
- 44: cda
- ---
- 49: dab
- ---
- 63:

Алгоритмы отыскания совершенных повторов Хеширование. Пример. Модульная функция

```
Пример. S = abcdabcbcbabcd;
l = 3; h_2(x_i) = h_1(x_i) \mod M; M = 11;
ns(a) = 0; ns(b) = 1; ns(c) = 2; ns(d) = 3
h_1(abc) = 6 \mod 11 = 6;
h_1(bcd) = 27 \mod 11 = 5;
h_1(cda) = 44 \mod 11 = 0;
h_1(\text{dab}) = 49 \mod 11 = 5;
h_1(abc) = 6 \mod 11 = 6;
h_1(bcb) = 25 \mod 11 = 3;
h_1(\text{cbc}) = 38 \mod 11 = 5;
h_1(bcb) = 25 \mod 11 = 3;
h_1(cba) = 36 \mod 11 = 3;
h_1(bab) = 17 \mod 11 = 6;
h_1(abc) = 6 \mod 11 = 6;
h_1(bcd) = 27 \mod 11 = 5;
```

Пример списковой схемы устранения наложений

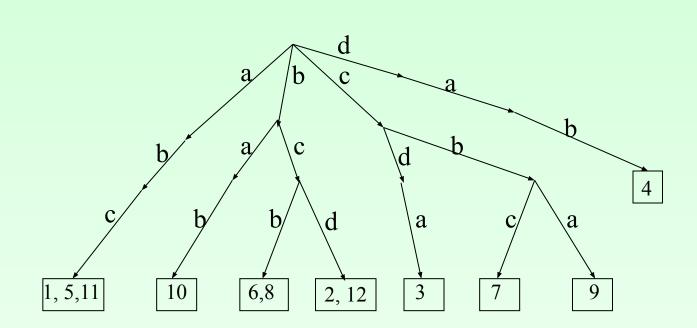
				f(x)	pos.					f(x)	pos.		
abc	1	0	cda	1	3	*	1	0	dab	1	5	•	\vdash
bcd		1						1	cbc	1	7	*	 -
cda		2						2	cba	1	3	*	
dab	1	3	bcb	1+1	6,8	•	1	3	bab	1	10	*	
abc	W	4				/							
bcb		5	bcd	1+1	2,12	•							
cbc		6	abc	1+1+1	1,5,11								-
bcb		7											
cba		8						ab	cda	bcb	cbak	occ	1
bab		9						65	056	353	3665	5	
abc		10											
bcd													

l-граммные деревья — это структура данных, представляющая все l-граммы в виде дерева.

S = abcdabcbcbabcd;

$$l = 3;$$

- 1. abc
- 2. bcd
- 3. cda
- 4. dab
- 5. abc
- 6. bcb
- 7. cbc
- 8. bcb
- 9. cba
- 10. bab
- 11. abc
- 12. bcd



Префиксное и суффиксное деревья

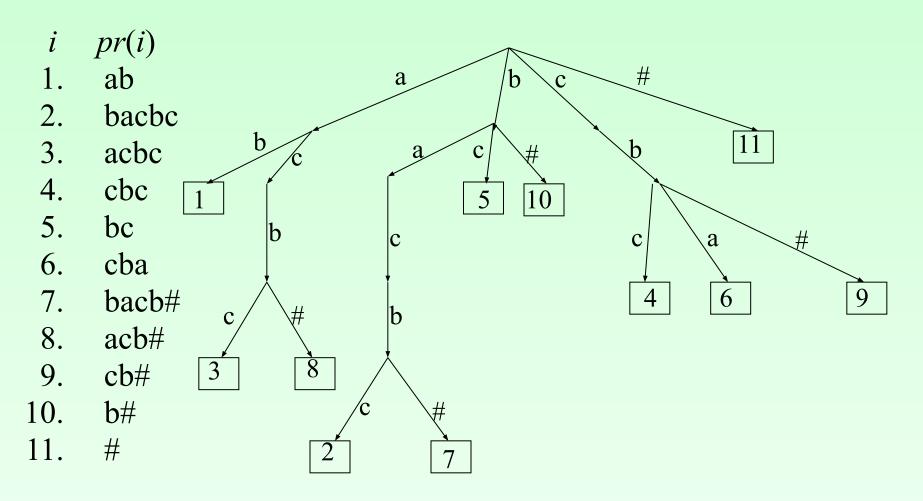
Если v = xyz, то x - префикс v, z - суффикс, <math>y - подслово. Оптимальные алгоритмы отыскания совершенных повторов основаны на построении префиксного дерева, суффиксного дерева или графа подслов текста (DAtG).

Первая конструкция разработана Вайнером (teiner P., 1973), вторая — Мак-Крейгом (McCreight, 1976), третья — (A.Blumer, J.Blumer, A.Ehrenfeucht, et al., 1984). Все конструкции функционально эквивалентны и реализуются за линейное (в зависимости от длины текста) время с линейными затратами памяти.

Префикс—идентификатор pr(i) позиции i в тексте T — кратчайшее подслово, начинающееся в позиции i и встречающееся в T # только один раз (# - конечный маркер).

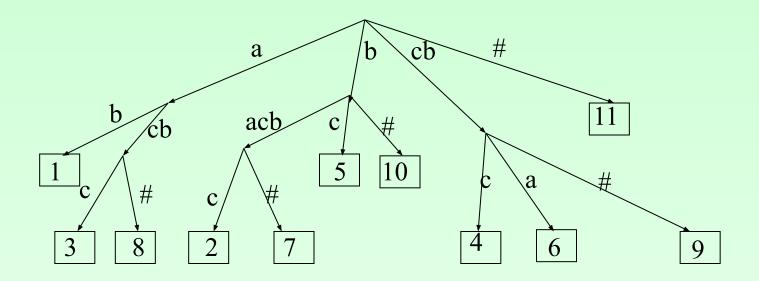
Префиксное дерево = дерево префикс-идентификаторов.

Пример дерева префикс-идентификаторов для T# = abacbcbacb#

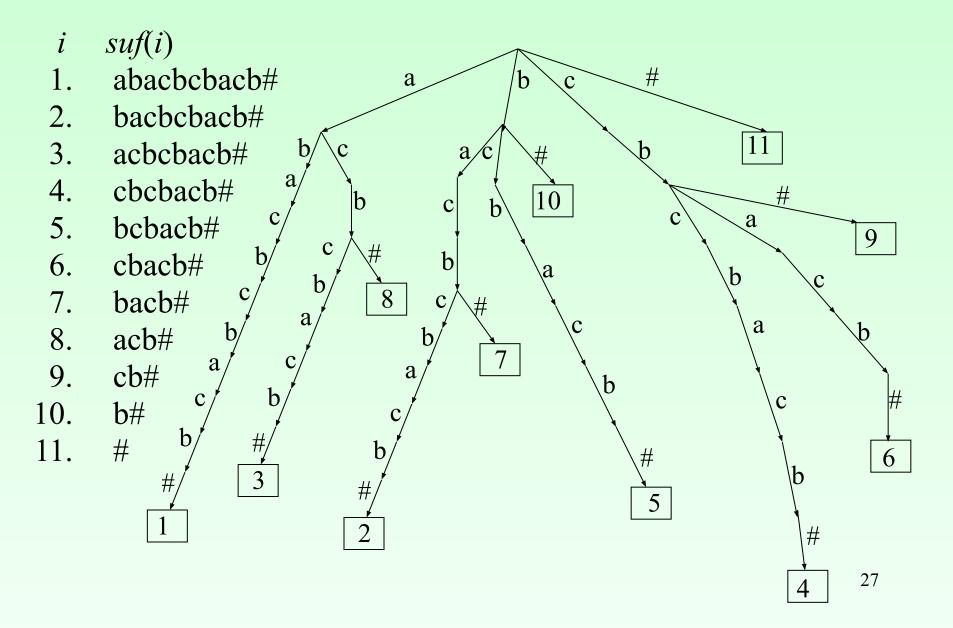


Алгоритм Мартинеца на примере T# = abacbcbacb# abacbcbacb# a 2,5,7,10 4,6,9 1,3,8 b 3,8 5 4,6,9 C b 3,8 9 4 6 8

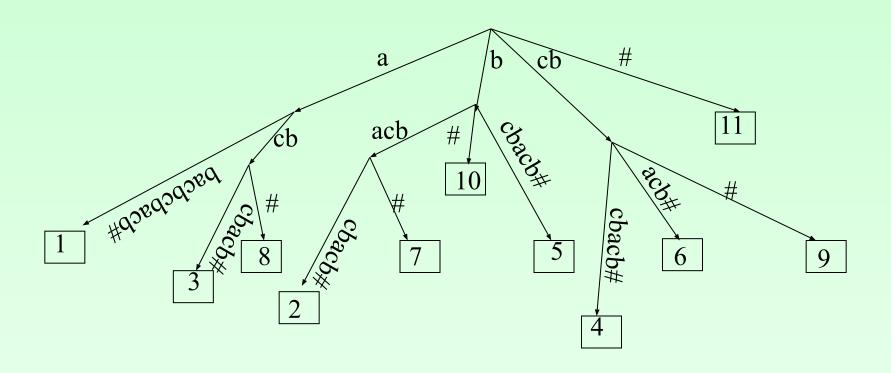
Пример компактного префиксного дерева для T# = abacbcbacb#



Пример дерева всех суффиксов для T# = abacbcbacb#



Суффиксное дерево для Т# = abacbcbacb#



Задачи, решаемые с помощью суффиксного дерева:

- Вычисление параметров полного частотного спектра;
- Поиск образца;
- Последовательный поиск множества образцов;
- Поиск образца во множестве строк;
- Наибольшая общая подстрока двух строк;
- Общие подстроки более чем двух строк;
- Задача загрязнения ДНК. Даны строки S_1 и S_2 : S_1 вновь расшифрованная ДНК, S_2 комбинация источников возможного загрязнения. Найти все подстроки S_2 , которые встречаются в S_1 и длина которых не меньше заданного l;
- Суффиксно-префиксные совпадения всех пар строк (из заданного множества строк);
- Обнаружение всех «нерасширяемых» повторов;
- Задача о наибольшем общем «продолжении». Найти длину наибольшего общего префикса i-го суффикса строки S_1 и j-го суффикса строки S_2
- Выявление всех «нерасширяемых» палиндромов.