

Расстояния между строками и меры близости

Расстояние Хэмминга: число несовпадений символов.

Пример: $A = abbaeac$; $d_{\text{ХЭМ}} = 2$.

$B = abdaecc$;

Редакционное расстояние. В простейшем виде расстояние $d(A,B)$ между строками A и B определяется как **минимальное число операций** типа "вставка", "устранение" или "замена" символа, переводящих одну строку в другую.

Пример: $A = abbaeac$; $B = bdedac$

$A \rightarrow B$:
a b b a e a c $d(A,B) = 4$ (1 дел., 3 зам).
— b d e d a c

$A \rightarrow B$:
a b b a e — a c $d(A,B) = 4$ (2 дел., 1 вст., 1 зам).
— b d — e d a c

Хроника введения расстояний и мер сходства, основанных на идеях динамического программирования.

Год	Наименование (меры, расстояния, процедуры сравнения)	Авторы	Предметная область
1965	Метрика Левенштейна	Левенштейн В. И.	теория связи (двоичные коды)
1968	процедуры нелинейной нормализации речи по темпу	Слуцкер Г.С.	распознавание речи
1968	ДП - метод	Винцюк Т.К.	распознавание речи
1969	дискретный вариант меры Слуцкера	Загоруйко Н.Г., Величко В.М.	распознавание речи
1970	мера сходства AM-последовательностей	Needleman S.B., Wunch S.B.	молекулярная биология
1974	редакционное расстояние	Wagner R.A., Fisher M.J.	лингвистика
1974	эволюционное расстояние	Sellers P.	молекулярная биология

Схема динамического программирования для вычисления редакционного расстояния:

$$d(0,0) = 0;$$

$$d(i,0) = i, 1 \leq i \leq m; m = |A|;$$

$$d(0,j) = j, 1 \leq j \leq n; n = |B|;$$

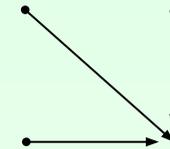
$$d(i,j) = \min \{ d(i-1, j-1) + \gamma(a_i, b_j), \\ d(i-1, j) + 1, \\ d(i, j-1) + 1 \}.$$

$$\gamma(a_i, b_j) = \begin{cases} 0, & a_i = b_j \\ 1, & a_i \neq b_j \end{cases}$$

B

A

	ϵ	a	b	b	a	e	a	c
ϵ	0	1	2	3	4	5	6	7
b	1	1	1	2	3	4	5	6
d	2	2	2	2	3	4	5	6
e	3	3	3	3	3	3	4	5
d	4	4	4	4	4	4	4	5
a	5	4	5	5	4	5	4	5
c	6	5	5	6	5	5	5	4



- Этап 1: заполнение матрицы динамического программирования.
- Этап 2: обратный ход. Построение выравнивания. $\forall d(i,j)$ на этапе 1 запоминаем указатели на элемент, от которого реализовался переход ($d[i-1, j-1]$, $d[i-1, j]$ и /или $d[i, j-1]$); либо на этапе 2 делаем проверку $d[i, j] = d[i-1, j-1] + \gamma(a_i \rightarrow b_j)?$, или $d[i, j] = d[i-1, j] + \gamma(a_i \rightarrow \varepsilon)?$, или $d[i, j] = d[i, j-1] + \gamma(\varepsilon \rightarrow b_j)?$.

Трудоемкость и затраты памяти $O(n \times m)$.

Если требуется только вычислить редакционное расстояние (без выравнивания) можно хранить две строки: текущую и предыдущую.

Для построения выравнивания требуется знать всю матрицу.

–	b	d	–	e	d	a	c	или	–	b	d	e	d	a	c
a	b	b	a	e	–	a	c		a	b	b	a	e	a	c

- В общем случае число редакционных операций можно расширить, например, добавить перестановку соседних символов или блочные вставки или удаления.
- Каждая операция может иметь свой "вес". В этом случае редакционное расстояние определяется как минимальная "стоимость" перевода A в B.
- Если к редакционным операциям добавлена операция перестановки соседних символов, схема динамического программирования выглядит следующим образом:

$$d[i, j] = \min \{ d[i-1, j-1] + \gamma(a_i \rightarrow b_j), \\ d[i-1, j] + \gamma(a_i \rightarrow \varepsilon), \\ d[i, j-1] + \gamma(\varepsilon \rightarrow b_j), \\ d[i-2, j-2] + \gamma(a_{i-1} a_i \rightarrow b_{j-1} b_j). \text{ если } a_i a_{i-1} = b_{j-1} b_j \}$$

при тех же начальных условиях

$$d(i,0) = i, 1 \leq i \leq m; m = |A|;$$

$$d(0,j) = j, 1 \leq j \leq n; n = |B|;$$

Возможные обобщения

- Поиск участков текста, **близких к заданному образцу**.
- Может формулироваться в виде: найти все фрагменты текста, для которых редакционное расстояние с образцом не превышает R (некоторый заданный параметр).
- Для решения этой задачи достаточно составить матрицу D , изменив начальные условия, а именно, положить $d[0, j] = 0$, для всех $1 \leq j \leq n$.
- Пример. $P = \text{baaa}$; $T = \text{bbabbaabab}$;

	ϵ	b	b	a	b	b	a	a	b	a	b
ϵ	0	0	0	0	0	0	0	0	0	0	0
b	1	0	0	1	0	0	1	1	0	1	0
a	2	1	1	0	1	1	0	1	1	0	1
a	3	2	2	1	1	2	1	0	1	1	1
a	4	3	3	2	2	2	2	1	1	1	2

Возможные обобщения:

- Поиск «похожих» фрагментов текста.
- Схема динамического программирования, но с начальными условиями:
 $d[0, j] = 0$, для $0 \leq j \leq n$ и $d[i, 0] = 0$ для $0 \leq i \leq m$.
- Однако удобнее использовать "меру близости", т.е. величину, противоположную ред. расстоянию: чем тексты ближе, тем бóльшие значения имеет мера близости.

Мера близости определяется следующими соотношениями

$$r(i, j) = \max \left\{ \begin{array}{l} 0, \\ r(i-1, j-1) + \delta(a_i, b_j), \\ r(i-1, j) + \delta(a_i, \varepsilon), \\ r(i, j-1) + \delta(\varepsilon, b_j) \end{array} \right\}$$

$$\delta(a_i, b_j) = 1, \text{ если } a_i = b_j;$$

$$\delta(a_i, b_j) = -2, \text{ если } a_i \neq b_j.$$

$$r[0, j] = 0, r[i, 0] = 0, \quad 0 \leq j \leq n, 0 \leq i \leq m.$$

$$\delta(a_i, \varepsilon) = \delta(\varepsilon, b_j) = -2.$$

	ϵ	c	b	b	c	c	a	a	a	b	a	a	a
ϵ	0	0	0	0	0	0	0	0	0	0	0	0	0
a	0	0	0	0	0	0	1	1	1	0	1	1	1
a	0	0	0	0	0	0	1	2	2	0	1	2	2
a	0	0	0	0	0	0	1	2	3	1	1	2	3
a	0	0	0	0	0	0	1	2	3	1	2	3	3
a	0	0	0	0	0	0	1	2	3	1	2	3	4
c	0	1	0	0	1	1	0	0	1	1	0	1	2
c	0	1	0	0	1	2	0	0	0	0	0	0	0
b	0	0	2	1	0	0	0	0	0	1	0	0	0
c	0	1	0	0	2	1	0	0	0	0	0	0	0
c	0	1	0	0	1	3	1	0	0	0	0	0	0

$$r(i, j) = \max \{0, \\ r(i-1, j-1) + \delta(a_i, b_j), \\ r(i-1, j) + \delta(a_i, \epsilon), \\ r(i, j-1) + \delta(\epsilon, b_j)\}$$

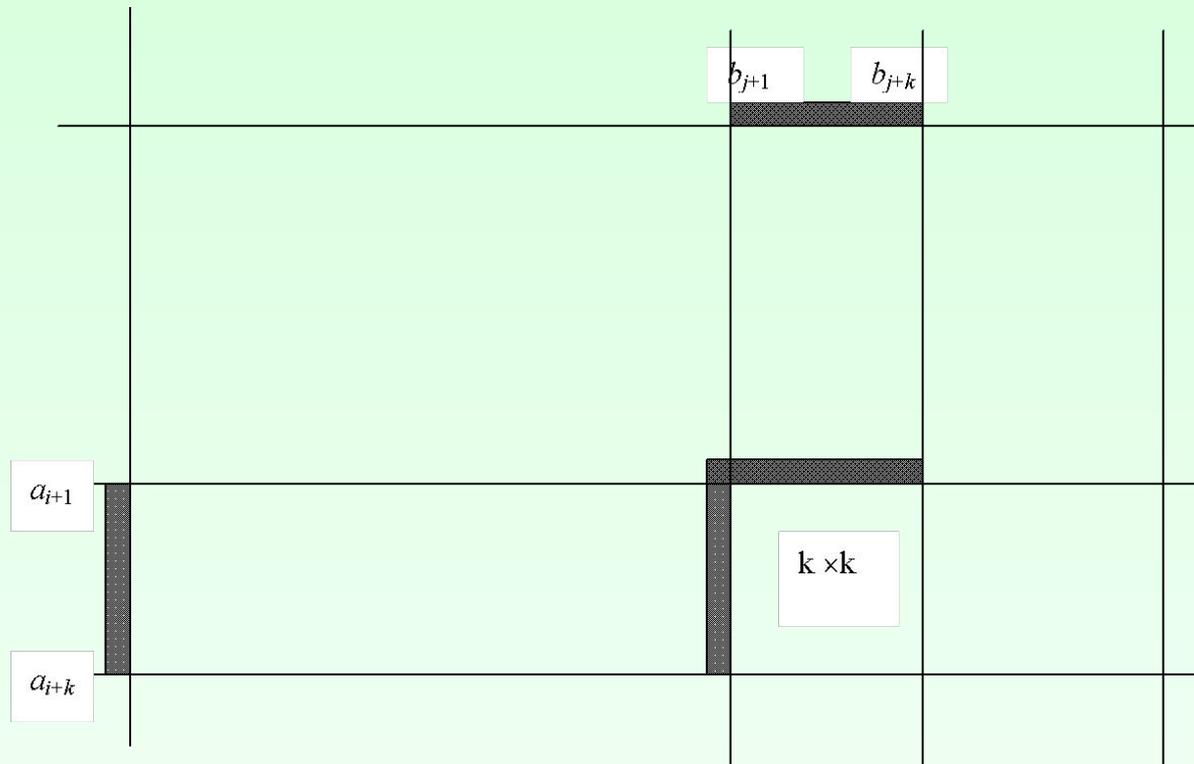
$$\delta(a_i, b_j) = 1, \text{ если } a_i = b_j \\ \delta(a_i, b_j) = -2, \text{ если } a_i \neq b_j \\ \delta(a_i, \epsilon) = \delta(\epsilon, b_j) = -2.$$

	c	b	b	c	c	a	a	a	b	a	a	a	ϵ
a	0	0	0	0	2	4	3	2	1	3	2	1	0
a	0	0	0	0	1	3	3	2	1	3	2	1	0
a	0	0	0	0	1	3	2	2	1	3	2	1	0
a	0	0	0	0	1	3	2	1	1	3	2	1	0
a	0	0	0	0	0	2	2	1	0	2	2	1	0
a	0	0	0	0	0	1	1	1	0	1	1	1	0
c	2	0	0	2	1	0	0	0	0	0	0	0	0
c	2	1	1	1	1	0	0	0	0	0	0	0	0
b	0	1	3	0	0	0	0	0	1	0	0	0	0
c	1	0	0	2	1	0	0	0	0	0	0	0	0
c	1	0	0	1	1	0	0	0	0	0	0	0	0
ϵ	0	0	0	0	0	0	0	0	0	0	0	0	0

$$r'(i, n+1) = 0; \\ r'(m+1, j) = 0;$$

$$r'(i, j) = \max \{0, \\ r'(i+1, j+1) + \delta(a_i, b_j), \\ r'(i+1, j) + \delta(a_i, \epsilon), \\ r'(i, j+1) + \delta(\epsilon, b_j)\}$$

Алгоритм Мазека-Патерсона. Для вычисления подматрицы D размера $k \times k$ должны быть известны начальные векторы и соответствующие фрагменты исходных текстов. Идея четырех русских: Арлазаров, Диниц, Кронрод, Фараджиев. Ускорение достигается за счет предварительного вычисления и сохранения информации обо всех вариантах вызова подзадачи. Тогда при решении задачи можно по введенным аргументам (подстроки текста и начальные векторы) сразу выписать ответ: конечные векторы, которые, в свою очередь, будут начальными для следующих фрагментов. Для сокращения числа аргументов можно кодировать разности между соседними элементами. $O(n^2/\log n)$



Максимально длинная общая подпоследовательность (МДП, LCS)

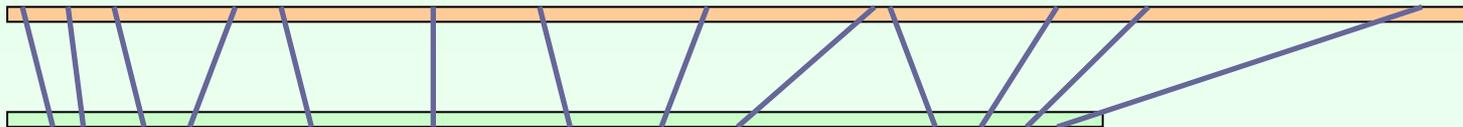
Последовательность U является подпоследовательностью A , если существует монотонно возрастающая последовательность целых чисел $r_1 \dots r_{|U|}$ такая, что $U[j] = A[r_j]$, $1 \leq j \leq |U|$, $1 \leq r_j \leq |A|$.

Если $\gamma(a_i \rightarrow b_j) \geq \gamma(a_i \rightarrow \varepsilon) + \gamma(\varepsilon \rightarrow b_j)$, то при переводе A в B используются только операции вставки и устранения, а элементы, составляющие LCS, остаются без изменения.

Если $\gamma(a_i \rightarrow b_j) = 2$ (при $a_i \neq b_j$), а $\gamma(a_i \rightarrow \varepsilon) = \gamma(\varepsilon \rightarrow b_j) = 1$,

$$D(A, B) = m + n - 2L(A, B), \quad m = |A|, n = |B|$$

$$2L(A, B) = m + n - D(A, B)$$



Вычисление длины МПД:

$$L(i, j) = \begin{cases} L(i-1, j-1) + 1, & \text{если } a_i = b_j \\ \max(L(i, j-1), L(j-1, j)) & \text{в остальных случаях} \end{cases}$$

$$L(0, j) = 0; \quad L(i, 0) = 0;$$

$$0 \leq i \leq m; \quad 0 \leq j \leq n;$$

Пример. A = aacacbb; B = ababc.

$$L(A, B) = 3.$$

Всего 12 МДП:

3 – abb,

6 – aab,

3 – aac

	ε	a	b	a	b	c
ε	0	0	0	0	0	0
a	0	1	1	1	1	1
a	0	1	1	2	2	2
c	0	1	1	2	2	3
a	0	1	1	2	2	3
c	0	1	1	2	2	3
b	0	1	2	2	3	3
b	0	1	2	2	3	3

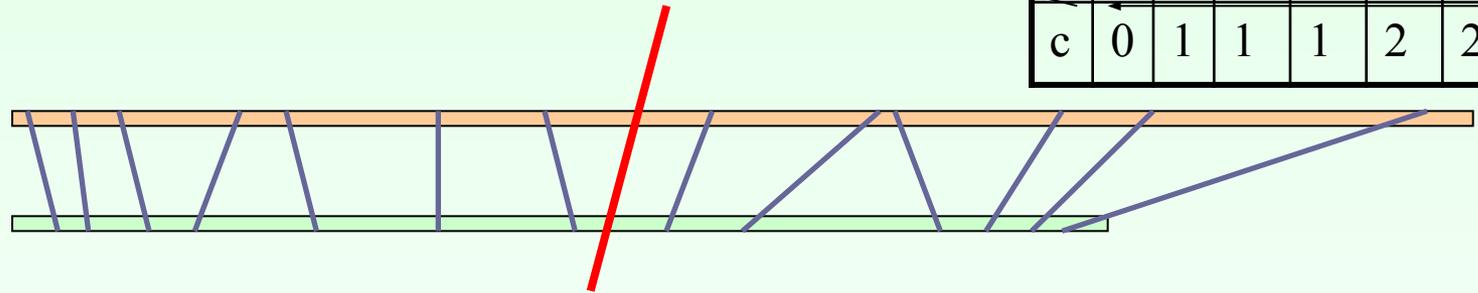
Алгоритм Хишберга (Hirschberg D.S.)

Пусть $L^*(i, j)$ – длина МДП текстов $A[i + 1 : m]$ и $B[j + 1 : n]$.

Тогда для любого i : $L(m, n) = \max_j \{L(i, j) + L^*(i, j)\}$

<p>$A = \text{aaca'cbb}$</p> <p>$B = \text{ababc}$</p> <p>$L(4,0) + L^*(4,0) = 0 + 2$</p> <p>$L(4,1) + L^*(4,1) = 1 + 2$</p> <p>$L(4,2) + L^*(4,2) = 1 + 1$</p> <p>$L(4,3) + L^*(4,3) = \underline{2 + 1}$</p> <p>$L(4,4) + L^*(4,4) = 2 + 1$</p> <p>$L(4,5) + L^*(4,5) = 3 + 0$</p>	<p>$A1 = \text{aa'ca}$</p> <p>$B1 = \text{aba}$</p> <p>$L(2,0) + L^*(2,0) = 0 + 1$</p> <p>$L(2,1) + L^*(2,1) = \underline{1 + 1}$</p> <p>$L(2,2) + L^*(2,2) = 1 + 1$</p> <p>$L(2,3) + L^*(2,3) = 2 + 0$</p>	<p>$A2 = \text{cb'b}$</p> <p>$B2 = \text{bc}$</p> <p>$L(2,0) + L^*(2,0) = 0 + 1$</p> <p>$L(2,1) + L^*(2,1) = \underline{1 + 0}$</p> <p>$L(2,2) + L^*(2,2) = 1 + 0$</p>
--	---	---

	ϵ	c	b	a	c	a
ϵ	0	0	0	0	0	0
b	0	0	1	1	1	1
b	0	0	1	1	1	1
c	0	1	1	1	2	2



Адаптивные алгоритмы вычисления длины МПД:

Hunt - Shimanski

Q_{ik} – наим. j , такое что $(L(A[1 : i], B[1 : j])) = k$

$$Q_{i+1,k} = \begin{cases} \text{наим. } j, & \text{такое что } Q_{i,k-1} < j \leq Q_{ik} \text{ и } a_{i+1} = b_j \\ Q_{ik}, & \text{если такого } j \text{ не существует} \end{cases}$$

Начальные условия:

$$Q_{i,0} = 0, \quad 0 \leq i \leq n;$$

$$Q_{0k} = n + 1, \quad 1 \leq k \leq \min(m,n);$$

Пример. $A = \text{a a s a c s b b}$,

$B = \text{a b a b c}$.

Трудоёмкость: $O(r \times \log n)$,

где
$$r = \sum_{i=1}^{|\Sigma|} f_i(A) \cdot f_i(B)$$

число потенциально возможных парных соответствий символов

	0	1	2	3	4
ε	0	6	6	6	6
a	0	1	6	6	6
a	0	1	3	6	6
c	0	1	3	5	6
a	0	1	3	5	6
c	0	1	3	5	6
b	0	1	2	4	6
b	0	1	2	4	6

Адаптивные алгоритмы вычисления длины МПД: Nakatsu-Kambayashi-Yajima

$R_{i,k}$ – наиб. j , такое что $(L(A[i:m], B[j:n])) = k$

$$R_{i,k} = \begin{cases} \text{наиб. } j, \text{ такое что } R_{i+1,k} \leq j < R_{i+1,k-1} \text{ и } a_i = b_j \\ R_{i+1,k}, & \text{если такого } j \text{ не существует} \end{cases}$$

Начальные условия:

$$R_{i,0} = n + 1, \quad 0 \leq i \leq m;$$

$$R_{m+2-k,k} = 0, \quad 1 \leq k \leq \min(m,n);$$

Пример. $A = \text{a a c a s c b b}$,

$B = \text{a b a b c}$.

Трудоёмкость: $O(n \times (m - L))$,

R	0	1	2	3	4
a	6	5	3	1	0
a	6	5	3	1	0
c	6	5	3	1	0
a	6	5	3	1	0
c	6	5	2	0	0
b	6	4	2	0	
b	6	4	0		
ε	6	0			

- **Близкие задачи:**
- задача о наикратчайшей **надпоследовательности**
- задача о **медиане** (string merging): построение текста T3, сумма переходов к которому от T1 и T2 минимальная. В зависимости от весов ред. операция в качестве T3 можно получить МДП, наименьшую надпоследовательность, один из текстов (T1 или T2), наиболее вероятного предка и т.п.
- поиск максимально длинной общей подпоследовательности для группы текстов.
- задача о максимально длинной возрастающей (убывающей) подпоследовательности для числовой перестановки

Меры сходства

а) Пусть T_1 и T_2 два текста.

Назовем **совместной частотной характеристикой** l -го порядка текстов T_1 и T_2 совокупность элементов

$$\Phi_l(T_1, T_2) = \{\varphi_{l1}(T_1, T_2), \varphi_{l2}(T_1, T_2), \dots, \varphi_{lM_l}(T_1, T_2)\}$$

где $M_l = M_l(T_1, T_2)$ — число l -грамм, общих для обоих текстов,

а элемент φ_{li} ($1 \leq i \leq M_l$) есть тройка:

\ll i -я общая l -грамма — x_i ,

частота ее встречаемости в T_1 — $F(T_1, x_i)$ и в T_2 — $F(T_2, x_i)$ \gg .

Простейший **набор мер сходства**, упорядоченный по возрастанию l имеет вид:

$$q_l(T_1, T_2) = \frac{M_l(T_1, T_2)}{M_l(T_1) + M_l(T_2) - M_l(T_1, T_2)}$$

$$l = 1, 2, \dots, l_{\max}(T_1, T_2)$$

б) более сложный вариант, учитывающий частоты встречаемости l -грамм:

$$\lambda(T_1, T_2) = \frac{\sum_{\alpha} \min\{F(T_1, \alpha), F(T_2, \alpha)\} \cdot |\alpha|}{\sum_{\alpha} \max\{F(T_1, \alpha), F(T_2, \alpha)\} \cdot |\alpha|}$$

где α – произвольная цепочка текстов T_1 и (или) T_2 ,
 $|\alpha|$ – ее длина.

(Findler N.V., Van Leeuten, 1979)

г) ранговая мера близости:

Пусть l -граммы в $\Phi_l(T_1)$ и $\Phi_l(T_2)$ упорядочены по убыванию частот; порядковое место l -граммы x_i в упорядочении определяет ее ранг – $r(T_1, x_i)$ (соответственно, $r(T_2, x_i)$).

Группы равночастотных l -грамм представляются усредненным рангом. Введем l -граммный аналог расстояния:

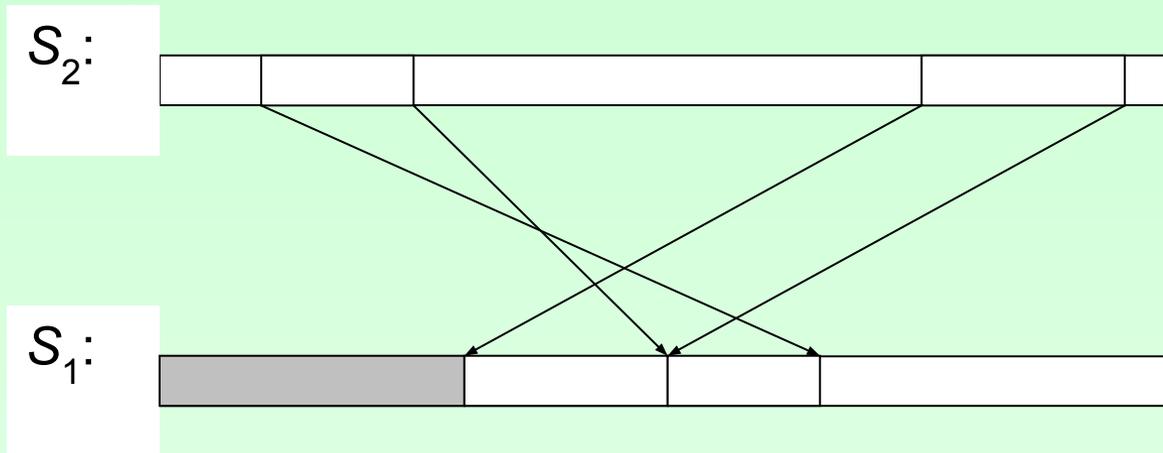
$$S_l(T_1, T_2) = \sum_{x_i \in \Sigma_l} (r(T_1, x_i) - r(T_2, x_i))^2$$

где Σ_l – совокупность всевозможных цепочек длины l ; $|\Sigma_l| = R_l = n^l$. Аналогом коэффициента Спирмэна для характеристики l -го порядка ($l = 1, 2, \dots$) является

$$\rho_l(T_1, T_2) = 1 - \frac{6S_l(T_1, T_2)}{R_l(R_l^2 - 1)}$$

При наличии равночастотных l -грамм в (***) вносится поправка на "связанность" рангов. (Кендел М. Ранговые корреляции, М., Статистика, 1975)

Обобщение подхода Лемпеля-Зива



- Представление S_1 в виде конкатенации фрагментов из S_2 назовем **сложностным разложением** S_1 по S_2 .
- На каждом шаге копируется максимальный фрагмент S_2 , совпадающий с префиксом непокрытого участка S_1
- Если такого фрагмента нет, используется операция генерации символа
- $c(S_1 / S_2)$ – сложность S_1 относительно S_2 определяется числом компонентов в разложении S_1 по S_2

Относительная сложность и редакционное расстояние

S_2 = аааа а сссс с тттттттттттт - асасасас а атататат

S_1 = аааа г сссс г тттттттттттт г асасасас г атататат

$d(S_1, S_2)$ = 4

$H(S_1/S_2)$ = аааа*г*сссс*г*тттттттттттт*г*асасасас*г*атататат

$c(S_1/S_2)$ = 9

- $c(S_1/S_2) \leq 2d(S_1, S_2) + 1$

- S_2 = -----тттттттттттттттттттаааааааа

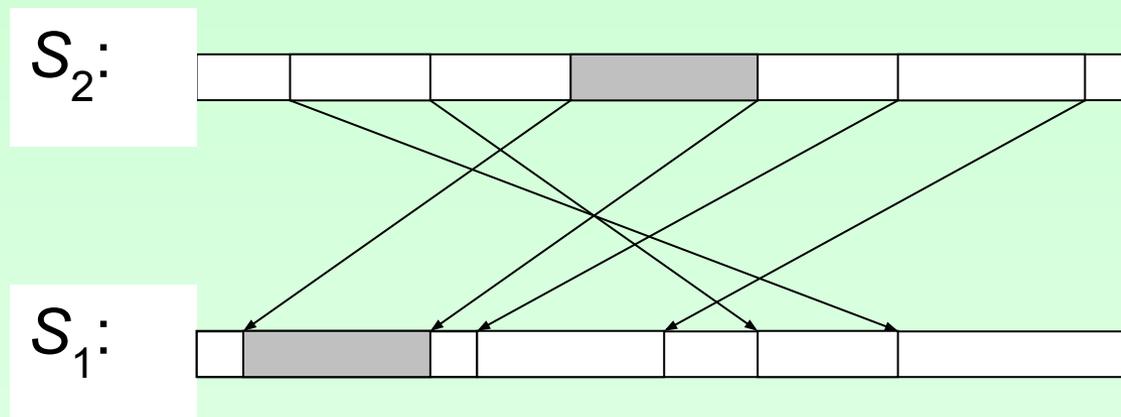
S_1 = ааааааааттттттттттттттттттт-----

$d(S_1, S_2) = 16$

$H(S_1 / S_2) = ааааа * ттттттттттттттттттт$

$c(S_1 / S_2) = 2$

Трансформационное расстояние



- Трансформационное расстояние и относительная сложность идейно близки.
- Операция «вставка сегмента» используется, если посимвольная генерация фрагмента «дешевле» его копирования.
- Порядок покрытия S_1 предполагает оптимизацию по всем парам межтекстовых повторов и промежуткам между ними. $O(N^6)$.

J.-S.Varré, J.-P.Delahaye, E. Rivals: Transformation Distances: a Family of Dissimilarity Measures Based on Movements of Segments. // Bioinformatics 15(3): 194-202 (1999)

Инверсионное расстояние

$$\left(\begin{array}{cccccccc} 1 & 2 & \boxtimes & i-1 & i & i+1 & \boxtimes & j-1 & j & j+1 & \boxtimes & N \\ & & & & \longleftarrow & \longrightarrow & & & & & & \\ 1 & 2 & \boxtimes & i-1 & j & j-1 & \boxtimes & i+1 & i & j+1 & \boxtimes & N \end{array} \right)$$

- Инверсионное расстояние $d_I(\pi, \sigma)$ между последовательностями π и σ определяется минимальным числом инверсий, переводящих одну из них в другую
Задача вычисления инверсионного расстояния для перестановок является NP-полной
- В случае "знаковых" перестановок существуют полиномиальные решения

$$+1 \ [+2 \ -4 \ -5] \ +3 \ +6 \quad \rightarrow \quad +1 \ +5 \ +4 \ -2 \ +3 \ +6$$

Hannenhalli, S. and Pevzner, P. Transforming Cabbage into Turnip (Polynomial Algorithm for Sorting Signed Permutation by Reversals). Proc. 27th Ann. ACM Symposium on the Theory of Computing, 1995, pp. 178–189

Точки разрыва

$$\pi_0 = 0 \text{ and } \pi_{N+1} = N + 1$$

π and σ – произвольные перестановки.

Разрыв между $\pi_i = a$ и $\pi_{i+1} = b$ фиксируется, если в σ нет биграмм ab и ba .

$$\pi = 0 | 6 4 | 1 8 5 | 3 | 2 9 7 | 10 \quad r(\pi, \sigma) = 5$$

$$\sigma = 0 | 5 8 1 | 2 9 7 | 6 4 | 3 | 10$$

σ – тождественная перестановка (1 2 ... N).

Число точек разрывов (**breakpoint distance**) $r(\pi, \sigma)$

определяется количеством позиций π таких что $|\pi_i - \pi_{i+1}| \neq 1$.

$$1 2 3 | 8 7 6 | 4 5 | 9$$

Инверсионное расстояние, число точек разрыва и относительная сложность

- $r(\pi, \sigma) \leq 2d_I(\pi, \sigma)$
- Точки разрыва однозначно соответствуют границам компонентов сложностного разложения

$$r(\pi, \sigma) = c(\pi / \sigma) - 1$$

$$\sigma = 1\ 2\ 3\ 4\ 5\ 6\ 7\ 8\ 9$$

$$H(\pi / \sigma) = 1\ 2\ 3\ *\ 8\ 7\ 6\ *\ 4\ 5\ *\ 9$$

- Сложностные разложения позволяют перейти от исходных перестановок к "знаковым" и сократить размерность задачи вычисления инверсионного расстояния