

ЧИСЛЕННЫЙ КРИТЕРИЙ БЛИЗОСТИ ТЕКСТОВ

В.А. Леус

В [1] изложен общий подход к заданию критериев близости (расстояний) на множестве произвольных конечных текстов. Он охватывает как частные случаи расстояние Хемминга и обобщенное расстояние Хемминга. Геометрическую структуру на множестве текстов можно вводить различными способами. Пусть расстояние ρ определяет одну геометрическую структуру, а расстояние ρ' - другую. Разрешающая способность расстояния ρ' выше, чем у расстояния ρ , если для любых трех текстов L, M, N неравенство $\rho(L, M) \neq \rho(L, N)$ влечет неравенство $\rho'(L, M) \neq \rho'(L, N)$, но не наоборот. Подход дает возможность достаточно полно учитывать различие и взаиморасположение символов в сравниваемых текстах и тем самым получать высокую разрешающую способность расстояния. Однако комбинаторный характер критериев близости высокого разрешения часто сопряжен с экспоненциальным объемом вычислений.

В настоящей статье предлагается компромиссный критерий близости для текстов. Метод вычисления значения критерия для двух данных текстов использует свойство выпуклости функции, задающей расстояние между отдельными вхождениями букв. Это открывает возможность для направленного поиска, трудоемкость которого оказывается полиномиально зависящей от объемов сравниваемых текстов. Приводятся примеры работы метода в некоторых прикладных задачах.

§1. Расстояние между текстами

Введение геометрической структуры в множестве текстов основывается на понятии контакта двух текстов M и N . Пусть a и b - произвольные буквы алфавита A . Вхождение буквы a на месте μ в текст M обозначим через a_μ , а буквы b на месте ν в текст

$N - b_v$. Контакты $K(M, N)$ рассматриваемых текстов M и N есть двудольный граф, одну долю вершин которого образуют вхождения букв в текст M , другую - вхождения букв в текст N , причем каждая вершина обязательно инцидентна одному и только одному ребру. Ребра a_μ и b_v называются связями, а висячие ребра (a_μ, Λ) (или b_v, Λ), где Λ - пустой символ, называются псевдосвязями. Каждой связи и псевдосвязи ставится в соответствие неотрицательное число, называемое ее длиной. При фиксированном способе назначения длин связей и псевдосвязей за длину контакта принимается сумма длин всех его ребер. Расстояние $\rho(M, N)$ между текстами есть минимум длин контактов этих текстов.

Чем больше тонкостей в различии текстов учитывает расстояние ρ , тем обширнее множество допустимых контактов. С другой стороны, чем менее представительно множество контактов, тем легче искать контакт с минимальной длиной (минимальный контакт) для данных текстов. В качестве контактов для текстов M и N будем рассматривать только такие двудольные графы, в которых: а) каждое ребро (связь)

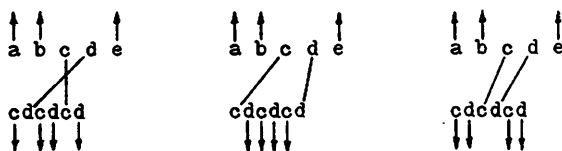


Рис. 1

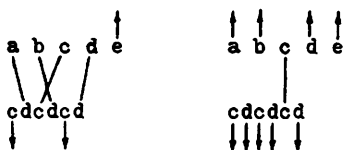


Рис. 2

инцидентно паре (a_μ, a_v) вхождений одной и той же буквы; б) все непарные вхождения (и только они) являются концами псевдосвязей. На рис. 1 показаны три разных контакта между текстами. Псевдосвязи обозначены стрелками, обращенными наружу. На рис. 2 даны примеры двудольных графов, которые контактами не являются (в левом нарушено условие "а", в правом - условие "б").

Допустим, что текст M состоит из m вхождений одной буквы, а текст N - из n вхождений той же самой буквы (для определенности $m \leq n$). Тогда число n всевозможных контактов $K(M, N)$ равно $C_n^m \cdot m! = n! / (n-m)!$. Если некая k -я буква алфавита имеет m_k

вхождений в текст M и n_k вхождений в текст N , то можно рассматривать контакты подтекстов (фрагментов текстов) M_k и N_k , состоящих из всех вхождений k -й буквы в тексты M и N . Число таких контактов $K_k(M, N)$ текстов M и N по k -й букве есть $n_k = n_k! / |n_k - m_k|!$, откуда количество всех контактов $K(M, N)$ не меньше, чем произведение Π_k по всем буквам, имеющим вхождения хотя бы в одно из сравниваемых слов. Мощност множества контактов есть очень быстро растущая функция от числа вхождений букв в сравниваемых текстах, поэтому актуальной оказывается задача сокращения перебора.

Настоящая работа посвящена методу направленного поиска контакта минимальной длины для двух конечных текстов. Основой для построения различных текстов служит целочисленная сетка, в каждом узле которой стоит один экземпляр пустого символа Λ . Текст получается в результате замены пустых символов в конечном числе позиций содержательными буквами алфавита A (все формально строгие определения приведены в [1]). При таком определении два текста, получаемые один из другого сдвигом по сетке, являются разными. Число вхождений содержательных букв в тексте назовем размером (или объемом) текста.

Количества связей и псевдосвязей по любой букве во всех контактах данной пары текстов постоянны. При любом способе назначения длины псевдосвязи сумма длин всех псевдосвязей одинакова для всех контактов данной пары текстов. В длине контакта выделяются две компоненты - суммарная длина связей и суммарная длина псевдосвязей. В этом смысле для обозначения расстояния удобно использовать символику комплексного числа $\rho = \alpha + \beta i$, где часть $\text{Re } \rho$ отвечает вкладу связей, а часть $\text{Im } \rho$ - вкладу псевдосвязей. Расстояние между текстами, имеющими одинаковые количества вхождений каждой буквы, имеет действительное значение, а расстояние любого содержательного текста от пустого - чисто мнимое.

Согласно условию "а" всякий контакт любых двух текстов можно разбить на непересекающиеся "контакты по букве" соответственно числу букв, имеющих вхождения в сравниваемых текстах. В минимальном контакте каждый "контакт по букве" минимален, так что процедура поиска минимального контакта сводится к независимому отысканию минимальных контактов по каждой из букв. Длина минимального контакта по некоторой букве "b" есть расстояние по букве $\rho_b(M, N)$ между текстами, а полное расстояние получается суммированием величин ρ_b по всем $b \in A$. Таким образом, для вычис-

ления расстояния кардинальной оказывается задача поиска минимального контакта текстов в алфавите A , содержащем одну букву, которую можно обозначить, например, символом i .

Связь, инцидентную вхождениям i_μ и i_ν , будем обозначать (μ, ν) , а ее длину — $s_{\mu\nu}$; псевдосвязь и ее длину естественно обозначить (μ, Λ) и s_μ (соответственно (ν, Λ) и s_ν). Величину $s_{\mu\nu}$ будем считать функцией от модуля разности $|\mu - \nu|$ индексов вхождений, а длину псевдосвязей — тождественно равной константе s . Очевидно, что расстояние между текстами обладает свойством симметрии, т.е. $\rho(M, N) = \rho(N, M)$. Если к тому же потребовать, чтобы функция $s_{\mu\nu}(|\mu - \nu|)$ принимала нулевые значения тогда и только тогда, когда $\mu = \nu$, то расстояние будет обладать свойством тождества: $\rho(M, N) = 0 \Leftrightarrow M = N$. Наконец, исключим для длин связей и псевдосвязей отрицательные значения, т.е. действительная и мнимая части расстояния всегда неотрицательны.

Пусть (μ_1, ν_1) и (μ_2, ν_2) — две связи в контакте $K(M, N)$. При $(\mu_2 - \mu_1) \cdot (\nu_2 - \nu_1) > 0$ связи назовем параллельными, в противном случае — скрещивающимися. Для любых двух вхождений μ_1 и μ_2 в тексте M и любых двух вхождений ν_1 и ν_2 в тексте N существуют две пары связей $((\mu_1, \nu_1), (\mu_2, \nu_2))$ и $((\mu_1, \nu_2), (\mu_2, \nu_1))$. Эти пары назовем сопряженными. В одной из сопряженных пар связи параллельные, в другой — скрещивающиеся.

ТЕОРЕМА I о парных связях. Если $s(x)$ — монотонная неубывающая выпуклая (вниз) функция, то сумма длин пары параллельных связей не превосходит суммы длин сопряженной пары скрещивающихся связей.

Для доказательства этой теоремы нам потребуется одно свойство выпуклых функций, которое сформулируем в виде леммы. Пусть $f(\vec{x})$ — выпуклая функция, определенная на R^n , \vec{x}_1 и \vec{x}_2 — радиус-векторы двух произвольных несовпадающих точек в R^n , t — действительное число ($0 \leq t \leq 1$), вектор $\vec{\delta} = t(\vec{x}_2 - \vec{x}_1)$, Δ_1 — приращение функции $f(\vec{x})$ при смещении из точки \vec{x}_1 в точку $(\vec{x}_1 + \vec{\delta})$, Δ_2 — приращение функции при смещении из точки $(\vec{x}_2 - \vec{\delta})$ в точку \vec{x}_2 .

ЛЕММА. Для выпуклой вниз функции $f(\vec{x})$ справедливо неравенство $\Delta_1 = f(\vec{x}_1 + \vec{\delta}) - f(\vec{x}_1) \leq f(\vec{x}_2) - f(\vec{x}_2 - \vec{\delta}) = \Delta_2$, а для выпук-

л о й в в е р х - о б р а т н о е н е р а в е н с т в о
 $\Delta_2 \leq \Delta_1$.

ДОКАЗАТЕЛЬСТВО леммы. Значение функции $f(\vec{x})$ в точке $\vec{x} = (\vec{x}_1 + \vec{\delta})$ есть

$$f(\vec{x}_1 + \vec{\delta}) = f(\vec{x}_1 + t\vec{x}_2 - t\vec{x}_1) = f((1-t)\vec{x}_1 + t\vec{x}_2).$$

Согласно определению выпуклости вниз имеем неравенство

$$f((1-t)\vec{x}_1 + t\vec{x}_2) \leq (1-t)f(\vec{x}_1) + tf(\vec{x}_2).$$

Значение функции $f(\vec{x})$ в точке $\vec{x} = (\vec{x}_2 - \vec{\delta})$ по определению выпуклости вниз оценивается следующим образом:

$$f(\vec{x}_2 - \vec{\delta}) = f(\vec{x}_2 - t\vec{x}_2 + t\vec{x}_1) = f(t\vec{x}_1 + (1-t)\vec{x}_2) \leq tf(\vec{x}_1) + (1-t)f(\vec{x}_2).$$

Отсюда выводим оценку для суммы значений

$$f(\vec{x}_1 + \vec{\delta}) + f(\vec{x}_2 - \vec{\delta}) \leq (1-t)f(\vec{x}_1) + tf(\vec{x}_1) + tf(\vec{x}_2) + (1-t)f(\vec{x}_2) = f(\vec{x}_1) + f(\vec{x}_2).$$

Или, перенося слагаемые, получим

$$\Delta_1 = f(\vec{x}_1 + \vec{\delta}) - f(\vec{x}_1) \leq f(\vec{x}_2) - f(\vec{x}_2 - \vec{\delta}) = \Delta_2.$$

Обратное неравенство $\Delta_2 \leq \Delta_1$ для функции $f(\vec{x})$, выпуклой вверх, получается аналогичным образом. Лемма доказана, переходим к доказательству теоремы о парных связях.

Пусть для определенности $\mu_1 < \mu_2$, $\nu_1 < \nu_2$, так что (μ_1, ν_1) и (μ_2, ν_2) — параллельные, а (μ_1, ν_2) и (μ_2, ν_1) — скрещивающиеся связи. Введем обозначения: $c = |\mu_1 - \nu_1|$, $d = |\mu_2 - \nu_2|$, $e = |\mu_1 - \nu_2|$, $f = |\mu_2 - \nu_1|$, $P = e + f - c - d$. Найдем значения величины P при различных взаиморасположениях индексов $\mu_1, \mu_2, \nu_1, \nu_2$ на числовой оси:

$$1. \quad \nu_1 < \nu_2 \leq \mu_1 < \mu_2:$$

$$P = \mu_1 - \nu_2 + \mu_2 - \nu_1 - \mu_1 + \nu_1 - \mu_2 + \nu_2 = 0;$$

$$2. \quad \nu_1 \leq \mu_1 \leq \nu_2 \leq \mu_2:$$

$$P = \nu_2 - \mu_1 + \mu_2 - \nu_1 - \mu_1 + \nu_1 - \mu_2 + \nu_2 = 2(\nu_2 - \mu_1) \geq 0;$$

$$3. \quad \mu_1 \leq \nu_1 < \nu_2 \leq \mu_2:$$

$$P = \nu_2 - \mu_1 + \mu_2 - \nu_1 - \nu_1 + \mu_1 - \mu_2 + \nu_2 = 2(\nu_2 - \nu_1) > 0;$$

$$4. \quad \nu_1 \leq \mu_1 < \mu_2 \leq \nu_2:$$

$$P = \nu_2 - \mu_1 + \mu_2 - \nu_1 - \mu_1 + \nu_1 - \nu_2 + \mu_2 = 2(\mu_2 - \mu_1) > 0;$$

$$5. \mu_1 \leq \nu_1 \leq \mu_2 \leq \nu_2 :$$

$$P = \nu_2 - \mu_1 + \mu_2 - \nu_1 - \nu_1 + \mu_1 - \nu_2 + \mu_2 = 2(\mu_2 - \nu_1) \geq 0;$$

$$6. \mu_1 < \mu_2 \leq \nu_1 < \nu_2 :$$

$$P = \nu_2 - \mu_1 + \nu_1 - \mu_2 - \nu_1 + \mu_1 - \nu_2 + \mu_2 = 0.$$

Все возможности исчерпаны, и всякий раз $P \geq 0$, т.е. всегда

$$c + d \leq e + f. \quad (1)$$

Дальнейшие рассуждения будем проводить в предположении, что $c \leq d$. Обратное неравенство рассматривать нет необходимости в силу симметрии длин параллельных и скрещивающихся связей относительно изменения направления возрастания индексов на противоположное.

При $\nu_2 \leq \mu_2$ $d = \mu_2 - \nu_2$ и заведомо $\nu_1 < \mu_2$, так что $f = \mu_2 - \nu_1$. Имеем $f - d = \mu_2 - \nu_1 - \mu_2 + \nu_2 = \nu_2 - \nu_1 > 0$, т.е.

$$d < f. \quad (2)$$

При $\mu_2 \leq \nu_2$ $d = \nu_2 - \mu_2$ и заведомо $\mu_1 < \nu_2$ и $e = \nu_2 - \mu_1$. Поэтому $e - d = \nu_2 - \mu_1 - \nu_2 + \mu_2 = \mu_2 - \mu_1 > 0$, откуда

$$d < e. \quad (3)$$

Пусть $f \leq e$, а $\nu_2 \leq \mu_2$. Тогда $f = \mu_2 - \nu_1$, $d = \mu_2 - \nu_2$. Если допустить, что $\mu_1 \leq \nu_1$, то $e = \nu_2 - \mu_1$, $c = \nu_1 - \mu_1$ и $(e - f) = (\nu_2 - \mu_1) - (\mu_2 - \nu_1) = (\nu_1 - \mu_1) - (\mu_2 - \nu_2) = (c - d) \leq 0$. Таким образом, здесь совместимым с условием $f \leq e$ является только равенство $f = e$. Но тогда согласно (2) имеем

$$d < f = e. \quad (4)$$

Предположим, что $\nu_1 < \mu_1$. Тогда $c = \mu_1 - \nu_1$ и либо $\mu_1 \leq \nu_2$, либо $\nu_2 < \mu_1$.

При $\mu_1 \leq \nu_2$ имеем $e = \nu_2 - \mu_1$, и

$$(e - f) = (\nu_2 - \mu_1) - (\mu_2 - \nu_1) = -(\mu_1 - \nu_1) - (\mu_2 - \nu_2) = -(c + d) < 0.$$

При $\nu_2 < \mu_1$ имеем $e = \mu_1 - \nu_2$ и получаем неравенство

$$(e - f) = (\mu_1 - \nu_2) - (\mu_2 - \nu_1) = (\mu_1 - \mu_2) + (\nu_1 - \nu_2) < 0.$$

В обоих случаях оказывается $e < f$, значит, неравенство $\nu_1 < \mu_1$ несовместимо с условием $f \leq e$. Таким образом, согласно (3) и (4), получается, что условия $c \leq d$ и $f \leq e$ влекут за собой неравенство $d < e$.

Пусть теперь $e \leq f$ и $\mu_2 < v_2$. Тогда $e = v_2 - \mu_1$ и $d = v_2 - \mu_2$. Если предположить, что $v_1 \leq \mu_1$, то $f = \mu_2 - v_2$, $c = \mu_1 - v_1$, откуда

$$(f - e) = (\mu_2 - v_1) - (v_2 - \mu_1) = (\mu_1 - v_1) - (v_2 - \mu_2) = (c - d) \leq 0.$$

Здесь с условием $e \leq f$ совместимо только равенство $e = f$, и согласно (3) имеем неравенство

$$d < e = f. \quad (5)$$

Предположим, что $\mu_1 < v_1$, тогда $c = v_1 - \mu_1$ и либо $v_1 \leq \mu_2$, либо $\mu_2 < v_1$. При $v_1 \leq \mu_2$ имеем $f = \mu_2 - v_1$, откуда

$$(f - e) = \mu_2 - v_1 - v_2 + \mu_1 = (\mu_1 - v_1) + (\mu_2 - v_2) = -c - d < 0.$$

При $\mu_2 < v_1$ имеем $f = v_1 - \mu_2$, откуда

$$(f - e) = v_1 - \mu_2 - v_2 + \mu_1 = (\mu_1 - \mu_2) + (v_1 - v_2) < 0.$$

Наше последнее предположение приводит, таким образом, к противоречию с условием $e \leq f$. Остается либо принять неравенство $v_1 \leq \mu_1$, из которого следует (5), либо отказаться от неравенства $\mu_2 < v_2$, что приводит к (3). Из (3) и (5) вытекает, что условия $c \leq d$ и $e \leq f$ влекут за собой неравенство $d < f$.

Общий итог проведенных рассуждений таков. Как бы ни были расположены на числовой оси индексы четырех вхождений в сопряженных парах скрещивающихся и параллельных связей, модуль разности концевых вхождений максимален у скрещивающейся связи. Принимая во внимание это обстоятельство, нетрудно получить утверждение теоремы.

Рассмотрим случай, когда $e \leq f$, и, по доказанному, $d < f$. Здесь имеются три возможности упорядочения величин c, d, e и f .

Пусть $e \leq c \leq d \leq f$, тогда из (1) имеем

$$0 \leq c - e \leq f - d = p. \quad (6)$$

Поскольку функция $s(x)$ по условию теоремы монотонная неубывающая, то, учитывая (6), получим

$$s(c) - s(e) = s(c - e + e) - s(e) \leq s(p + e) - s(e).$$

В силу выпуклости вниз функции $s(x)$ имеем по доказанной лемме

$$s(e + p) - s(e) \leq s(f) - s(f - p) = s(f) - s(f - f + d) = s(f) - s(d).$$

Поэтому $s(c) - s(e) \leq s(f) - s(d)$ или, что то же, $s(c) + s(d) \leq s(e) + s(f)$.

Пусть теперь $c \leq e \leq d \leq f$. Дважды используя только монотонное неубывание функции $s(x)$, получаем

$$s(c) + s(d) \leq s(e) + s(d) \leq s(e) + s(f).$$

Пусть, наконец, $c \leq d \leq e \leq f$. Поскольку c и d по отдельности и вместе не превосходят величин e и f , то в силу монотонного неубывания $s(x)$ очевидно, что

$$s(c) + s(d) \leq s(e) + s(f).$$

Случай, когда $f \leq e$ и $d < e$, также допускает три возможности, которые мы кратко рассмотрим по аналогии.

а) $f \leq c \leq d < e$. Из (I) имеем

$$0 \leq c - f \leq e - d = q. \quad (7)$$

Используя неубывание $s(x)$ и учитывая (7), имеем

$$s(c) - s(f) = s(c - f + f) - s(f) \leq s(q + f) - s(f),$$

а в силу выпуклости вниз получим

$$s(f + q) - s(f) \leq s(e) - s(e - q) = s(e) - s(d).$$

Отсюда $s(c) - s(f) \leq s(e) - s(d)$ или, что то же, $s(c) + s(d) \leq s(e) + s(f)$.

б) $c \leq f \leq d < e$. В силу неубывания $s(x)$ имеем

$$s(c) + s(d) \leq s(f) + s(d) \leq s(f) + s(e).$$

в) $c \leq d \leq f \leq e$. Также в силу неубывания

$$s(c) + s(d) \leq s(f) + s(e).$$

Во всех возможных случаях мы убедились в том, что по сумме длин параллельные связи не превосходят скрещивающихся связей сопряженной пары. Теорема доказана.

§2. Алгоритмы вычисления расстояния

Рассмотрим тексты M и N , состоящие из вхождений одной и той же буквы. Для определенности примем, что число вхождений в текст M не превосходит числа вхождений в текст N , т.е. $m \leq n$. Пусть M' есть некоторый подтекст текста M ($m' < m$). Имеет смысл говорить о расстоянии $\rho(M', N)$ между M' , рассматриваемым как отдельный текст, и текстом N . Пусть контакт K' реализует это расстояние, а π' — множество тех концов его связей, которые входят в N . Назовем π' множеством, ближайшим к подтексту M' (для краткости — просто ближайшим).

Всякий текст является объединением (в теоретико-множественном смысле) своих подтекстов. Выберем непересекающиеся подтексты M_1, M_2, \dots, M_I так, что $M' = \bigcup_{i=1}^I M_i$. Обозначим через $\pi_1, \pi_2, \dots, \pi_I$ множества, ближайшие к M_1, M_2, \dots, M_I соответственно. Эти ближайшие могут иметь непустые пересечения между собой. В [1] доказано существование такого π' , ближайшего к подтексту M' , которое включает в себя объединение ближайших, т.е. $\pi' \supseteq \bigcup_{i=1}^I \pi_i$.

Пусть N' — подтекст текста N , причем γ — наименьший, δ — наибольший индексы вхождения этого подтекста. Подтекст N' называется сегментом текста N , если всякое вхождение a_v текста N при $\gamma \leq v \leq \delta$ является также и вхождением подтекста N' . В [1] показано, что если объединение ближайших $\bigcup_{i=1}^I \pi_i$ есть сегмент, то среди содержащих его ближайших к объединению $\bigcup_{i=1}^I M_i$ найдется сегмент. Эти два утверждения из [1] и сформулированная в предыдущем параграфе теорема о парных связях положены в основу метода направленного поиска минимального контакта двух текстов.

Алгоритмы поиска минимального контакта сводится к повторяющемуся процессу отыскания ближайшего к объединению подтекстов M_i . На первом этапе работы алгоритма за исходное принимается разбиение текста M на отдельные вхождения. Для каждого подтекста M_i (состоящего здесь из одного вхождения a_μ) ищется ближайшее из вхождений a_v текста N . В силу условия теоремы о парных связях длина связи $s_{\mu v}$ как функция $s_\mu(v)$ переменного индекса v при фиксированном индексе μ является выпуклой вниз функцией, поэтому вычисленные для двух соседних аргументов значения $s_\mu(v_1)$ и $s_\mu(v_2)$ этой функции указывают направление ее убывания. Минимальная связь выявляется сменой убывания функции $s_\mu(v)$ на возрастание. Когда найдены все ближайшие π_i (на первом этапе состоящие из одного вхождения каждое), может оказаться, что они попарно не пересекаются, тогда искомым минимальный контакт состоит из найденных минимальных связей. В противном случае объединение всех π_i разбивается на минимально возможное число J сегментов seg_j ($j=1, 2, \dots, J < I$). Вхождения сегмента seg_j связаны с вхождениями текста M , образующими некоторый подтекст M_j . Подтексты M_j попарно не пересекаются. исчерпывают все M , и их совокупность берется в качестве исходного разбиения ($\bigcup_j M_j = M$) для второго этапа работы алгоритма.

Второй этап состоит в повторении циклов отыскания ближайшего к объединению при известном объединении ближайших. Любой подтекст M_j , в свою очередь, состоит из подтекстов M_{ij} , для каждого из которых в предыдущем цикле алгоритма было найдено ближайшее π_i (используем опять индексацию по i , хотя это уже другие подтексты). Задача состоит в отыскании ближайшего π_j к объединению $\bigcup_i M_{ij} = M_j$. Поскольку ближайшее к объединению содержит объединение ближайших, то минимальный контакт подтекста M_j с текстом N нет необходимости искать среди всех возможных контактов. Область поиска можно еще сузить, если учесть, что каждое $\bigcup_i \pi_i$ есть сегмент seg_i , и среди содержащих его ближайших к объединению найдется сегмент, т.е. достаточно просматривать только такие контакты, вершины которых в тексте N образуют сегменты. Наконец, в силу теоремы о парных связях минимальный контакт следует искать среди параллельных контактов, содержащих только параллельные связи. Если подтекст $M_j = \bigcup_i M_{ij}$ содержит m_j вхождений, а объединение ближайших $\bigcup_i \pi_i$ содержит n_j вхождений ($n_j \leq m_j$), искомый минимальный контакт находится среди $(m_j - n_j) + 1$ параллельных контактов по числу сегментов с m_j вхождениями, содержащими сегмент из n_j вхождений. Множество этих контактов обозначим K_j . На рис. 3 для случая $m_j = 7$, $n_j = 3$ показаны два "зацепленных" сегмента из пяти возможных.



Рис. 3

В множестве K_j возможен направленный поиск контакта минимальной длины. Все контакты в K_j имеют одни и те же вершины в тексте M , и каждый из них индивидуализирован младшим индексом \tilde{v} своей вершины из числа вхождений текста N . Все множество K_j упорядочивается по возрастанию \tilde{v} . Длина контакта $K(\tilde{v})$ складывается из длин его параллельных связей. Длина каждой связи по условию теоремы I является выпуклой вниз функцией от \tilde{v} , и так как эти связи входят в параллельные контакты, вершины которых в тексте M фиксированы, длина каждой связи есть выпуклая вниз функция и от \tilde{v} . Это обуславливает возможность направленного перебора множества K_j ; вычисление длин контактов $K(v_1)$ и $K(v_2)$ для двух соседних значений младших индексов указывает направление убывания длины. Смена убывания на возрастание говорит о том, что минимальный контакт найден. Возможно применение и более эффективных приемов поиска экстремума выпуклой функции.

После того как в каждом K_j найден минимальный контакт, совокупность всех их вхождений в текст N образует объединение ближайших к подтекстам M_j . Если число этих вхождений равно m , т.е. найденные ближайшие попарно не пересекаются, то объединение минимальных контактов по всем j и дает искомый минимальный контакт для текстов M и N .

На рис. 4 приведен пример работы алгоритма при поиске минимального контакта текстов $M(m=39)$ и $N(n=40)$. Полностью связи показаны лишь для первого этапа, а в циклах второго этапа, чтобы не загромождать рисунок, связи показаны частично. Минимальный контакт длиной 255 найден на шестом цикле второго этапа. Имеется в виду действительная часть расстояния, а минимая часть здесь равняется единице ($40 - 39 = 1$).

Трудоемкость поиска объединения ближайших на первом этапе работы алгоритма пропорциональна произведению $m \times n$, так как для каждого вхождения в тексте M достаточно перебрать все вхождения в текст N , чтобы отыскать ближайшее. Для одноциклового второго этапа в худшем случае объединение ближайших для всех вхождений в текст M состоит из одного вхождения в текст N . Тогда множество $K_j (j=1)$ имеет мощность $(m-1)+1=m$ и минимальный контакт отыскивается перебором длин m контактов по m связей в каждом.

В многоцикловом варианте второго этапа подтексты M_1, M_2, \dots, M_J имеют соответственно количества вхождений m_1, m_2, \dots, m_J . В худшем случае на каждом цикле второго этапа происходит присоединение только одного из J исходных ближайших. Сначала оказываются самыми ближайшими для $j=1$ и $j=2$, и минимальный контакт в первом цикле ищется самое большое либо среди $(m_1 + m_2) - m_1 + 1$ при $m_2 < m_1$, либо среди $(m_1 + m_2) - m_2 + 1$ при $m_1 < m_2$ контактов с $(m_1 + m_2)$ связями каждый. Трудоемкость здесь не превосходит величины $(m_1 + m_2) \times [m_2 + 1]$. Во втором цикле к объединению $\pi_1 \cup \pi_2$ присоединяется π_3 и минимальный контакт ищется либо среди $(m_1 + m_2 + m_3) - (m_1 + m_2) + 1$ контактов, если $m_3 < (m_1 + m_2)$, либо среди $(m_1 + m_2 + m_3) - m_3 + 1$ контактов, если $(m_1 + m_2) < m_3$, с числом связей $(m_1 + m_2 + m_3)$. Во всех случаях трудоемкость не превосходит величину, пропорциональную $(m_1 + m_2 + m_3) \cdot [m_3 + 1]$. По аналогии трудоемкость второго этапа можно оценить величиной

$$(m_1 + m_2) \cdot [m_2 + 1] + (m_1 + m_2 + m_3) \cdot [m_3 + 1] + \dots + \sum_1^J m_j [m_j + 1].$$

Максимируя каждое слагаемое в круглых скобках суммой всех m_j и вы-

нося ее, получим выражение $[m_2 + m_3 + \dots + m_j + m - 1] \cdot \Sigma m_j$, которое очевидным образом оценивается сверху величиной $2 \Sigma m_j \cdot \Sigma m_j = 2m^2$. Таким образом, доказана

ТЕОРЕМА 2 о трудоемкости вычисления расстояния между текстами. Алгоритм вычисления расстояния для двух текстов с m и n вхождениями имеет по числу операций не более чем билинейную по m и n трудоемкость.

§3. Практическое применение

Рассмотрим несколько примеров сравнения текстов, описывающих структуры различной природы. Длину псевдосвязи считаем единичной, а за длину связи принимаем величину модуля разности индексов связанных вхождений.

Если индексы всех содержательных вхождений текста M_0 изменить на одну и ту же величину Δ (смещение), то получится новый текст M_Δ , который мы будем рассматривать как результат сдвига текста M_0 на величину смещения (влево - отрицательное, вправо - положительное). Для двух сравниваемых текстов M_0 и N_0 расстояние $\rho(M_\Delta, N_0)$ равно расстоянию $\rho(M_0, N_{-\Delta})$. Условимся "двигать" текст с меньшим числом вхождений, т.е. считать $N(m \leq n)$ неподвижным.

ПРИМЕР I. Для текстов $M(m=39)$ и $N(n=40)$, показанных на рис.4, были вычислены расстояния при различных смещениях слова M в пределах от -3 до 17 с единичным шагом. На рис.5 слева приведен график изменения действительной части в зависимости от смещения, откладываемого по оси абсцисс. При нулевом смещении действительная часть принимает значение 255, для получения которого и был реализован процесс поиска минимального контакта, изображенный на рис.4. Приведенный график содержит характерную особенность зависимостей $Re_p(\Delta)$, построенных для любой пары текстов, а именно неограниченный рост при стремлении к бесконечности абсолютной величины смещения. Достаточно сдвинуть текст M вправо так, чтобы его наименьший индекс превысил наибольший индекс текста N , как все связи минимального контакта при дальнейших смещениях будут только удлиняться. Аналогичная ситуация достигается и при сдвигах влево. Где-то в промежуточном положении имеется расстояние с минимальной действительной частью, называемое расстоянием между свободными

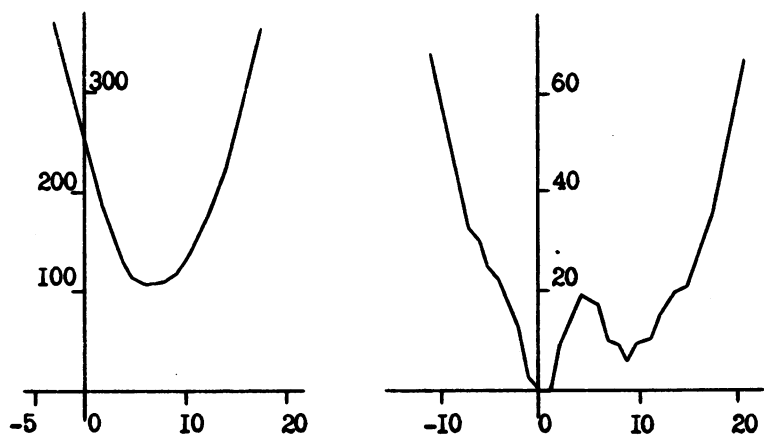


Рис. 5

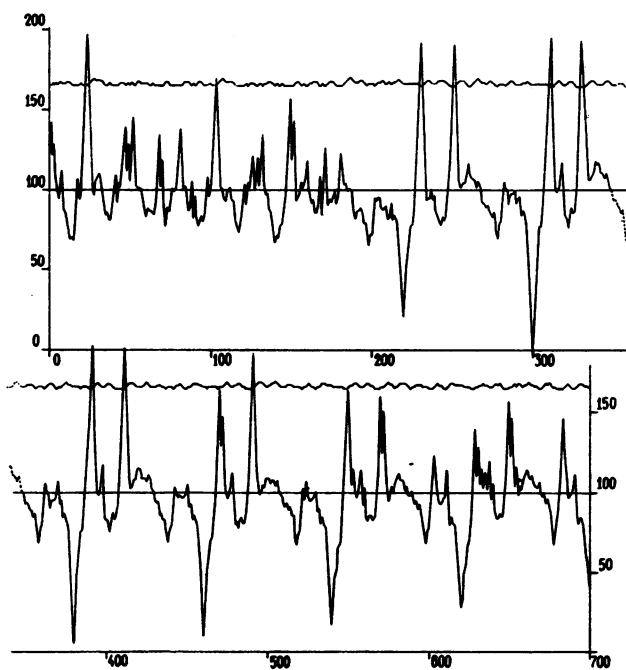


Рис. 6

текстами M и N . В нашем примере это расстояние со значением $Re_p \approx 105$ достигается при смещении $\Delta = 6$.

ПРИМЕР 2. Зависимость $Re_p(\Delta)$ необязательно имеет вид выпуклой (вниз) кривой с единственным минусом. Так, для текстов $M(IIII00000IIII00000III)$ и $N(IIII0000IIII0000IIII00000000IIII)$ на рис.5 справа приведена зависимость $Re_p(\Delta)$, которая имеет два минимума, разделенных максимумом, равным 19, в точке $\Delta = 4$. В точке $\Delta = 9$ находится локальный минимум со значением 6. Глобальный минимум с нулевым значением достигается в двух соседних точках $\Delta = 0$ и $\Delta = 1$, а минимая часть $Im_p = 9$. Расстояние между текстами в этом примере имеет, таким образом, чисто мнимое значение. Это является формальным выражением того факта, что текст $M(\Delta)$ при определенном смещении оказывается подтекстом текста M .

ПРИМЕР 3. В табл. I приведены короткий текст M и длинный текст N , вдоль которого перемещается короткий (вверху для удобства показана нумерация позиций). При каждом смещении из текста $M(\Delta)$ "вырезается" связный кусок - окно, включающее все позиции текста $M(\Delta)$. Вхождения текста N , попадающие в окно, образуют слово \tilde{N} . Размеры окна могут быть расширены за пределы слова $M(\Delta)$ так, чтобы выполнялось неравенство $m < n$. Для различных смещений вычислялось расстояние между текстами N и $M(\Delta)$. Результаты графически представлены на рис.6. Кривая большой амплитуды - график действительной части, а вверх вынесен график мнимой части, которая (при выбранном окне) изменяется в пределах от 0 до 5. Начиная со смещения $\Delta = 200$, на кривой $Re_p(\Delta)$ прослеживается повторяемость с периодом, равным 80.

ПРИМЕР 4. Осадконакопление в морском бассейне является существенно неравномерным процессом. Сносимый с разрушающихся гор в море обломочный материал накапливается сначала в верхней части материкового склона, а когда масса превысит некоторую критическую величину, срывающийся мутьевой поток перестраивает материал на дно бассейна. История осадконакопления фиксируется в виде последовательности чередующихся слоев различных по толщине и зернистости. Для кодирования структуры Большекарской свиты - древнего осадочного района [2] - была выбрана линейная единица, равная толщине минимального слоя, а все слои были разбиты по составу на семь типов, обозначаемых буквами. Структура разреза описывается текстом, в котором каждому слою поставлено в соответствие столько вхождений буквы данного типа, сколько линейных единиц укладывается на толщине слоя.

Т а б л и ц а 1

123456789 123456789 123456789 123456789 123456789 123456789 123456789

M 0010001000000111100101110000000011111000011110001011000001000011100001001111100

N 10000100011110000001000001101111000000000000111111010001000001010011000110101
 0100010000111100000100001110001111000000000001111011100000100001100010100111001
 010000100001111000010001110000011100000000011110011100000100001100100100111010
 001000100000111100010011100000001111000000111100011100000100001101000100111100
 001000100000011110010111000000000111100001111000101100000100001110000100111100
 010000100000011110001111000000000111100010111000101010000100001110000010111100
 10000010000001111000111010000000001111000101111000100110000100001101000010111010
 10000100000001011100011011000000000111100011011100010011000100000101110000101111001
 1000100000000110110001011100000000101110100111011000010110010000000111000001111001
 1001000000010101010010111010000001011011001111010000011101000000000111000011100111

В табл.2 приведен один из участков, где присутствуют слои типов А,В,С. К четвертому типу были отнесены слои неопределенной природы, которые кодировались нулем, соответствующим пустому символу. Этот участок рассматривался в качестве текста М, а в качестве текста N взят некоторый фрагмент из М длиной в 80 позиций, который смещался вдоль М. Для каждого положения текста N восьмидесятипозиционным окном из М "вырезался" фрагмент, который и сравнивался с текстом N. На рис.7 изображены зависимости (от смещения Δ) общего расстояния $\rho(N, M')$ (жирная кривая) и расстояний по буквам (пунктир - А, тонкая сплошная - В, пунктир с точкой - С). При смещении $\Delta = 600$ все они проходят через ноль, так как фрагмент N взят из текста М именно здесь.

На рис.8 приведены результаты локального сравнения двух текстов. В качестве одного из них взят текст М, а второй получен периодическим повторением того же, что и в предыдущем случае, фрагмента N. Окно на 80 позиций смещалось вдоль этих текстов, и для вырезаемой пары слов вычислялось расстояние. Зависимости расстояний от смещения даны на рис.8 в той же символике, что на рис. 7. Естественно, соответствующие кривые совпадают через каждые 80 позиций.

ПРИМЕР 5. Последовательностями символов описываются биомолекулы. С точки зрения вычисления критерия близости была рассмотрена транспортная рибонуклеиновая кислота (тРНК), молекулярная структура которой приведена в [3]. Эта переносящая аминокислоту валин тРНК имеет сравнительно небольшую молекулу из 77 нуклеотидных звеньев, которые кодируются буквами А,Г,Ц,У. В молекулу входят также редко встречающиеся так называемые неканонические нуклеотиды, которым сопоставлены символы нуля. Закодированная таким образом структура тРНК представлена в табл.3. Для сравнения биологических структур имеет значение не только близость по схожести, но и близость по комплементарности, когда нулевое расстояние приписывается парам разных символов. В нашем примере взаимно комплементарными являются А-У и Г-Ц.

В качестве текста N взята закодированная (табл.3) структура валиновой тРНК. Начальный участок этой структуры, включающий семь звеньев, имеет вид ГАУУУЦГ, а комплементарный ему код ЦУАААГЦ взят в качестве текста М. На рис. 9 точками показана зависимость действительной части расстояния $\rho(M, N)$ от смещения текста М. В

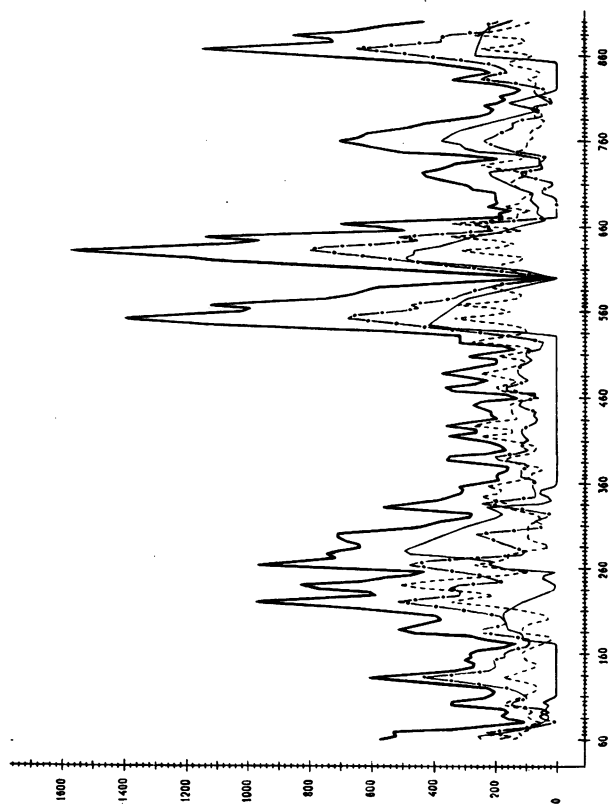


Рис. 7. Зависимость общего расстояния и расстояния по буквам от смещения Δ .

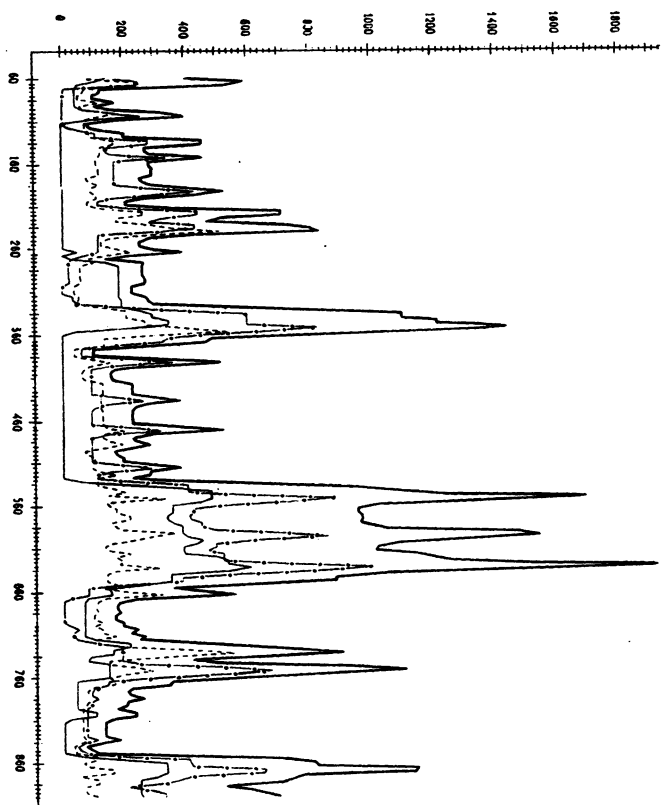


Рис. 8. Локальные сравнения двух текстов.

положениях $\Delta = 34$ и $\Delta = 66$ на этой кривой имеются минимумы, но нулевое значение в них не достигается.

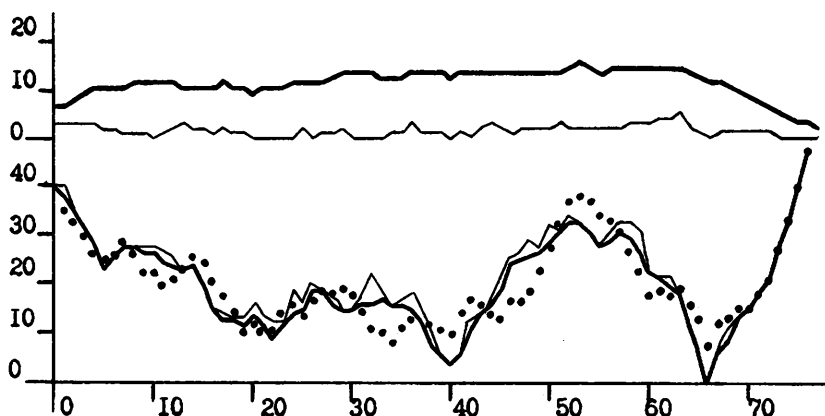


Рис. 9

Затем в качестве текста M был взят тот же комплементарный код, но в обратном порядке - ЦГАААУЦ. Сплошные кривые на рис. 9 изображают зависимости Re_p (внизу) и Im_p (вверху) от смещения для этого варианта слова M . Тонкая кривая соответствует минимальному окну, жирная - окну, расширенному на 8 позиций в обе стороны. При расширении окна количество псевдосвязей в контакте может только возрасти, а сумма длин связей, напротив, только уменьшиться. По этой причине на рис. 9 жирная кривая для Im_p лежит выше тонкой, а для Re_p , наоборот, ниже. В положении $\Delta = 66$ оба варианта имеют нулевую действительную часть (в случае малого окна вообще $\rho(M, N) = 0$). Этого и следовало ожидать, поскольку реальная молекула имеет пространственную структуру "клеверного листа", где начальный участок ГАУУУЦГ соединен водородными связями в обратной последовательности с участком ЦГАААУЦ, начинающимся с 67-й позиции [3].

ПРИМЕР 6. Дискретизация плоского изображения сеткой раstra превращает его в двумерный текст (см. [1]), где каждый элемент расположения - пиксель - становится вхождением, принимая значение из конечного алфавита (например - градаций яркости). Расстояние между двумерными текстами можно вводить различными способами, но экспоненциальный рост времени вычисления неприемлем для практическо-

[illegible]

го применения способа. Доступным в этом смысле является способ расщепления. Он сводит задачу вычисления расстояния между двумерными текстами к вычислению множества расстояний между одномерными текстами. На рис.10 приведены примеры текстов изображений, для

Т а б л и ц а 4

	80	ИХ	
BC	349 + 96 1	1612 + 2241	(x)
	46 + 202 1	1180 + 3261	(y)
	395 + 298 1	2792 + 5501	(p)
80	0	2032 + 1741	(x)
	0	1148 + 4481	(y)
		3180 + 6221	(p)

которых производилось расщепление по строкам и по столбцам. Между соответствующими строками, рассматриваемыми как одномерные тексты, вычислялись расстояния по вышерассмотренной методике. Затем то же производилось со столбцами, и все вычисленные "одномерные" расстояния суммировались. В

табл.4 сведены результаты. Как и следовало ожидать, минимальное расстояние получилось в паре "BC,80", а время вычисления на ЭВМ ЕС-1060 одного расстояния между двумерными словами размерами 25x80 пикселей - примерно 10 с.

З а к л ю ч е н и е

Близость текстов определяется как минимум функционала, заданного на множестве контактов текстов. Мощность множества контактов экспоненциально растет при возрастании размеров сравниваемых текстов. В статье предложен компромиссный подход к вычислению критерия близости, использующий тривиальную метрику на алфавите и метрику числовой оси для сравнения индексов вхождения букв. Это позволяет расщеплять тексты на подтексты в унарном алфавите и применять для последних методы направленного поиска минимума выпуклого функционала. Таким образом достигается полиномиальная зависимость трудоемкости вычислений от размеров сравниваемых текстов.

Л и т е р а т у р а

1. ЛЕУС В.А. Формальные критерии близости слов // Проблемы управления и теории информации. - 1979. - Т.8, №4. - С. 313-325.
2. ЕГАНОВ Э.А., СОВЕТОВ Д.К. Каратау - модель региона фосфоритонакопления. - Новосибирск: Наука, 1979. - 192 с.
3. ВОЛЬКЕНШТЕЙН М.В. Физика и биология. - М.: Наука, 1980. - 151 с.

Поступила в ред.-изд.отд.
23 октября 1987 года