

УДК 519.688:519.766.44

АЛГОРИТМ ВЫЧИСЛЕНИЯ СОВМЕСТНОГО ЧАСТОТНОГО СПЕКТРА
ДВУХ ТЕКСТОВ

В.Д.Гусев, Т.Н.Титкова

§1. Определение совместного частотного спектра

Сопоставление разных объектов (или групп объектов) друг с другом лежит в основе всех задач классификации, в том числе классификации символьных последовательностей (текстов). Для определения формальных мер сходства двух текстов T_1 и T_2 удобно ввести понятие совместной частотной характеристики 1-го порядка:

$$\Phi_1(T_1, T_2) = \{\varphi_{11}, \varphi_{12}, \dots, \varphi_{1M_1}(T_1, T_2)\}, \quad (1)$$

где $M_1(T_1, T_2)$ - количество различных 1-грамм (связных подпоследовательностей из 1 символа), общих для обоих текстов, а φ_{1r} ($1 \leq r \leq M_1(T_1, T_2)$) есть тройка: $\langle r$ -я 1-грамма, частота ее встречаемости в первом тексте - $F_1^{(1)}(r)$, частота встречаемости во втором тексте - $F_1^{(2)}(r) \rangle$. Вошедшие в $\Phi_1(T_1, T_2)$ 1-граммы могут быть упорядочены, например, по убыванию суммарной частоты встречаемости, лексикографически или в соответствии с каким-либо другим принципом, вытекающим из содержательной постановки задачи. Из определения (1) следует, что при $M_1(T_1, T_2) \neq 0$ $F_1^{(i)}(r) \geq 1$ для $i = 1, 2$ и всех значений r .

Совокупность совместных частотных характеристик

$$\Phi(T_1, T_2) = \{\Phi_1(T_1, T_2), \Phi_2(T_1, T_2), \dots, \Phi_L(T_1, T_2)\}, \quad (2)$$

где L - максимальное значение l , при котором $M_l(T_1, T_2)$ отлично от нуля, назовем совместным частотным спектром текстов T_1 и T_2 . Иными словами, параметр L характеризует длину максимального общего для обоих текстов участка. Величина L может изменяться от 0 (тексты с непересекающимися алфавитами) до $\min\{N_1, N_2\}$, где N_i ($i = 1, 2$) - длина i -го текста. Параметр L , вообще говоря, не

связан с параметрами $l_{\max}(T_1)$ и $l_{\max}(T_2)$, характеризующими длины максимальных повторов внутри текстов T_1 и T_2 соответственно.

В работе предлагается алгоритм вычисления $\Phi(T_1, T_2)$, ориентированный на тексты, длина которых не превышает размера оперативной памяти. Алгоритм итеративен по l . Каждая итерация имеет линейную (по $N_1 + N_2$) трудоемкость и требует $O(N \log_2 N)$ бит памяти, где $N = \max\{N_1, N_2\}$. Число итераций определяется степенью близости двух текстов, количественной оценкой которой в первом приближении может служить значение параметра L . Заметим, что для близких по длине случайных последовательностей L имеет порядок $\left\lceil \frac{2 \ln(N_1 + N_2)}{\ln \sum_{r=1}^n p_r^2} \right\rceil$, где p_r ($1 < r \leq n$) - вероятности встречаемости в текстах элементов алфавита $[I]$. Ниже будет упомянута одна из возможных адаптивных стратегий, позволяющая в ряде случаев уменьшить число итераций.

Близкими в идейном плане к данной проблеме являются работы [2-5]. В работах [2,3] предложены функционально эквивалентные конструкции (дерево идентификаторов и суффиксное дерево), позволяющие, в частности, с линейными затратами (по трудоемкости) отыскивать наиболее длинные общие для двух текстов цепочки. В принципе эти конструкции могут быть приспособлены для решения рассматриваемой нами задачи, но они достаточно сложны, что приводит к большой мультипликативной константе в оценках памяти и трудоемкости. Главный же недостаток их в том, что они значительно более критичны к росту мощности алфавита, чем описываемый ниже алгоритм.

В работах [4,5] рассмотрены соответственно задачи отыскания наиболее длинных повторов и вычисления полного частотного спектра применительно к одному тексту, длина которого существенно превышает размер доступной оперативной памяти (ОП). Отличия данной работы от [4,5], таким образом, заключаются в том, что речь идет о другом объекте (совместном частотном спектре) и введено ограничение на длины текстов (алгоритм для "оперативной памяти"). Эти отличия являются принципиальными, и хотя предлагаемый алгоритм в целом сохраняет идеологию, проводившуюся в [4,5], реализация его потребовала специфических решений.

§2. Алгоритм вычисления $\Phi(T_1, T_2)$

Предлагаемый алгоритм представляет собой итеративную процедуру (аналогично [5]), на каждом шаге которой вычисляется совмест-

ная частотная характеристика l -го порядка ($l = 1, 2, \dots, L^* \leq L$). Параметр L^* характеризует переход на другую схему обработки, что бывает целесообразно, когда число l -грамм, общих для T_1 и T_2 , уже мало и частота их встречаемости невелика.

Для подсчета частот при $l < L^*$ используется аппарат ассоциативного кодирования (хеширования [6]), оперирующий с конструкцией, называемой расстановочным полем (РП). При $l \geq L^*$ используется процедура расширения общих l -грамм до значений l , при которых они перестают быть общими.

Для указания l -грамм, не подлежащих обработке на последующих этапах, т.е. не вошедших в $\Phi_1(T_1, T_2)$, используется конструкция, названная "информационной лентой" (ИЛ). Это битовая строка длиной $N_1 + N_2$ двоичных разрядов (N_1 для T_1 и N_2 для T_2), каждый из которых соответствует l -грамме текста с тем же номером и может находиться в одном из двух состояний: 0 - l -грамма подлежит обработке, 1 - не подлежит ("лишняя" l -грамма). Перед началом обработки всем элементам ИЛ приписываются нулевые значения. Таким образом, по ходу l -й итерации текст просматривается слева направо скользящей рамкой, захватывающей l символов, которая скачками движется вдоль текста, пропуская только те позиции, которые помечены кодом 1 в ИЛ. Подлежащие анализу l -граммы хешируются, и информация об их частоте накапливается в РП.

И т е р а ц и я. Алгоритмы, описанные в [4,5], являлись "многопрогонными", т.е. каждая итерация состояла из нескольких просмотров текста. Это объяснялось их ориентацией на последовательности большой длины, не размещающиеся целиком в ОП (некоторым аналогом могут служить процедуры внешней сортировки). В рассматриваемом алгоритме каждая итерация по l осуществляется за один просмотр обоих текстов, что достигается наложением ограничений на их длину: необходимо, чтобы каждый из текстов в отдельности мог быть при обработке целиком размещен в ОП.

Введенное ограничение не следует считать слишком обременительным, поскольку многие реальные тексты, для которых возникает задача вычисления $\Phi(T_1, T_2)$, ему удовлетворяют, и класс таких текстов непрерывно расширяется в связи с быстрым ростом ОП у современных ЭВМ. Заметим, что логика однопрогонных алгоритмов значительно упрощается по сравнению с многопрогонными, поскольку отпадает необходимость в разбиении (иногда нетривиальном) текста на

части, каждая из которых обрабатывается за один прогон, и объединении результатов обработки отдельных частей.

2. Структура расстановочного поля. Введение ограничений на длины текстов позволяет (помимо отмеченных выше преимуществ) упростить и ускорить саму процедуру хеширования. Достигается это изменением структуры элементов РП по сравнению с [5].

В [5] каждый элемент РП состоял из двух зон, в одну из которых записывалась 1-грамма, а другая служила счетчиком числа вхождений ее в текст. В данном алгоритме вместо 1-граммы указывается номер ее позиции в тексте по первому символу (для восстановления 1-граммы необходимо, чтобы текст целиком размещался в ОП), вместо одного счетчика вводится два для фиксации числа вхождений 1-граммы в каждый из текстов и резервируется место под адрес отсылки для организации списка.

Такое структурирование элементов РП обусловлено тем, что при наличии одного из текстов в ОП нет необходимости дублировать каждую из M_1 его различных 1-грамм в РП. Нет необходимости иметь в ОП и оба текста одновременно, поскольку $M_1(T_1, T_2) \leq \min \{M_1(T_1), M_1(T_2)\}$, т.е. все 1-граммы, которые потенциально могут войти в $\Phi_1(T_1, T_2)$, уже содержатся в том тексте, который размещен в ОП.

Достоинства подобного структурирования заключаются:

а) в экономии памяти (при большой мощности алфавита): вместо $1 \cdot \lceil \log_2 n \rceil$ двоичных разрядов (n — мощность алфавита), требующихся для записи самой 1-граммы, теперь используются $\lceil \log_2 N_1 \rceil$ разрядов ($i = 1, 2$) для записи ее позиции в тексте. К примеру, при $n = 256$ и $N = 32000$ получаем экономию памяти, уже начиная со значения $1 = 3$, причем тем большую, чем больше 1 ;

б) в переходе от динамической схемы определения числа элементов РП и размера одного элемента к статической, не меняющейся от итерации к итерации (длина записи номера позиции 1-граммы в тексте не зависит от 1);

в) в упрощении механизма работы с ИЛ (см. ниже).

3. Механизм работы с информационной лентой (устранение "лишних" 1-грамм). Роль единичных [5] или "лишних" 1-грамм в данном алгоритме играют 1-граммы, присутствующие только в одном из текстов (их частота может быть произвольной). Обозначим множество 1-грамм, присутствующих только в T_1 или в T_2 соответственно через $X_1(T_1)$ и $X_1(T_2)$, а

множество общих для T_1 и T_2 1-грамм - через $X_1(T_1, T_2)$. Нетрудно видеть, что множества $X_1(T_1)$ и $X_1(T_2)$ образуются из левосторонних и правосторонних расширений (1-1)-грамм из $X_{1-1}(T_1)$ и $X_{1-1}(T_2)$ и таких расширений (1-1)-грамм из $X_{1-1}(T_1, T_2)$, которые в результате добавления 1-го символа перестают быть общими. Лишние 1-граммы первого типа естественно назвать потомственными, а второго типа - вновь образованными. В предлагаемой версии алгоритма все потомственные 1-граммы исключаются с помощью ИЛ автоматически, т.е. без хеширования. Вновь образующиеся лишние 1-граммы могут быть выявлены, лишь пройдя сито отбора (по частотам) в РП.

В многопрогонных алгоритмах информацию о единичных 1-граммах, выявленных после очередного прогона, невозможно отобразить в ИЛ без дополнительного прогона, поскольку в РП хранится сама 1-грамма, а не адрес ее вхождения в текст. Дополнительный же прогон, хотя и совмещен с обработкой следующей порции 1-грамм, но требует перехеширования уже выявленных единичных 1-грамм. Иное структурирование элементов РП в рассматриваемом алгоритме позволяет отобразить в ИЛ информацию о выявленных лишних 1-граммах без перехеширования.

Из определения $\Phi(T_1, T_2)$ вытекает, что на каждой итерации в РП следует фиксировать (путем указания адреса) 1-граммы лишь одного текста за исключением потомственных). Пусть это будет текст T_1 . Его 1-граммы обрабатываются в режиме: поиск (в РП) - размещение (если поиск оканчился неудачно). Все подлежащие обработке 1-граммы второго текста обрабатываются только в режиме "поиск". Неудачный поиск свидетельствует о том, что данная 1-грамма присутствует только в T_2 , т.е. является лишней. В соответствующий разряд ИЛ заносится код I, и анализируется следующая 1-грамма. При удачном поиске в счетчик, соответствующий T_2 , добавляется единица.

Таким образом, за один просмотр обоих текстов выявляются и устраняются из дальнейшего рассмотрения все 1-граммы, принадлежащие $X_1(T_2)$, а также часть 1-грамм, принадлежащих $X_1(T_1)$. Это 1-граммы с нулевым значением второго счетчика, причем лишь те из них, адреса которых указаны в РП и соответствуют первым вхождениям указанных 1-грамм в текст. Все последующие вхождения этих 1-грамм в T_1 остаются неидентифицированными на данном прогоне, но они (точнее, (1+1)-граммные расширения их) выявляются и устраняются из дальнейшего рассмотрения уже на следующей (1+1)-й итерации.

С этой целью на $(i+1)$ -й итерации тексты T_1 и T_2 меняются местами, т.е. текст T_2 располагается в ОП (и, следовательно, его l -граммы обрабатываются в режиме "поиск - размещение"), а текст T_1 располагается (если это необходимо) на внешнем носителе и небольшими порциями прокачивается через ОП. По окончании $(i+1)$ -й итерации устраняются из дальнейшего анализа все $(i+1)$ -граммы, принадлежащие $X_{i+1}(T_1)$ и часть $(i+1)$ -грамм из $X_{i+1}(T_2)$. Тексты T_1 и T_2 вновь меняются местами, и процесс повторяется.

Информация о части вновь образовавшихся лишних l -грамм в тексте, просматриваемом в первую очередь, может быть перенесена в ИЛ, например, при дополнительном просмотре РП (его не избежать ввиду необходимости очистки РП перед следующей итерацией). При этом выявляются l -граммы с нулевым значением второго счетчика и заносится код I (по адресам, указанным в РП) в элементы ИЛ, соответствующие тексту T_1 .

Выявление и фиксацию лишних потомственных l -грамм удобно осуществлять по ходу каждой итерации. При этом используется ИЛ, полученная в итоге предыдущей $(i-1)$ -й итерации. Те разряды этой ленты, которые содержат код I , сохраняют его: лишняя $(i-1)$ -грамма порождает при расширении вправо на один символ лишнюю же l -грамму. Те разряды ИЛ, которые содержат код 0 , вслед за которым непосредственно следует единица, меняют код 0 на код I : лишняя $(i-1)$ -грамма порождает лишнюю l -грамму и при расширении влево на один символ. Собственно, обработке, т.е. хешированию с последующим поиском в РП, подвергаются лишь l -граммы, помеченные кодом 0 , непосредственно за которым также следует код 0 .

Итак, формирование ИЛ включает в себя три процесса: процесс автоматического выявления и фиксации потомственных лишних l -грамм по ходу всей итерации, процесс выявления и фиксации вновь образующихся лишних l -грамм из текста, обрабатываемого во вторую очередь (по ходу просмотра этого текста), и процесс фиксации части вновь образовавшихся лишних l -грамм из текста, обработанного в первую очередь (по окончании прогона).

4. Схема устранения наложений. В алгоритмах [4,5] наложения устранялись за счет "многoproгонности": l -граммы, не идентифицированные на данном прогоне из-за наложений, идентифицировались на одном из следующих прогонов. Такая схема была обусловлена спецификой задачи (объем ОП $S \ll N$) и приводила

к существенному упрощению структуры РП (отсутствие дополнительного поля, адресов отсылок, повторного хеширования с целью поиска свободной позиции).

В данном алгоритме, ввиду того, что каждой итерации соответствует точно один прогон, наложения приходится устранять классическими средствами. Используется списковая схема устранения наложений с разбиением РП на основное и дополнительное. Специфика решаемой задачи используется на этапе выбора оптимальных размеров основного и дополнительного полей, а также при определении стратегии поведения в ситуациях, когда дополнительное РП оказывается заполненным.

Оценки требуемого размера РП Q и рекомендации по разбиению его на основное (Q_0) и дополнительное (Q_d) поля могут быть получены из следующих соображений.

Если выбрать $Q = \max\{N_1, N_2\} = N$, то 1-граммы любого из чередующих текстов могут быть обработаны за один просмотр, поскольку $M_1(T_1) \leq N-1+1$ ($i = 1, 2$).

С другой стороны, при фиксированном коэффициенте загрузки α списковый метод требует для своей реализации в среднем $Q_0(\alpha + e^{-\alpha})$ единиц памяти. Объединение двух этих условий приводит к уравнению

$$Q_0 + Q_d = N = Q_0(\alpha + e^{-\alpha}), \quad (3)$$

из которого (при заданном α) Q_0 может быть выражено через N и, следовательно, определено желаемое соотношение между размерами основного и дополнительного РП.

При задании α целесообразно ориентироваться не на потенциально возможный наихудший случай, а на средние экспериментально наблюдаемые значения случайной величины

$$q = \frac{1}{N-1+1} \max(M_1 - |X_{1-1}(T)|)$$

для данного класса текстов. В экспериментах с различными генетическими текстами значения q , например, колебались в диапазоне от 0,6 до 0,9, причем максимум приходился на значение $1 = 7$. Ориентируясь на среднее значение $q = 0,75$, получим коэффициент загрузки

$$\alpha = \frac{0,75 \cdot N}{Q_0}, \text{ с учетом чего (3) может быть переписано в виде } \alpha = 0,75(\alpha + e^{-\alpha}) \text{ или } \frac{1}{3}\alpha = e^{-\alpha}, \text{ откуда } \alpha \approx 1,05, Q_0 = 0,71 N \text{ и } Q_d =$$

$= N - Q_0 = 0,29N$. На классах текстов с меньшим коэффициентом загрузки Q_0 соответственно возрастет, а $Q_{\text{д}}$ уменьшится.

Пусть Q_0 и Q_g зафиксированы. Опишем стратегию занесения в РП и поиска в нем 1-грамм, обеспечивающую возможность вычисления $\Phi_1(T_1, T_2)$ за один просмотр текстов. Отличие от классического спискового метода устранения наложений возникает лишь в том случае, когда дополнительное РП оказывается заполненным ("плохая" последовательность). Поскольку суммарный объем РП ориентирован на наилучший случай, количество 1-грамм, оставшихся необработанными после заполнения дополнительного РП, не превышает количества оставшихся незаполненными позиций основного РП. Эти позиции последовательно выявляются простейшей процедурой типа "чистки мусора" и используются для расширения дополнительного РП. Ниже описана реализация этой схемы.

Для каждой 1-граммы x вычисляется адрес позиции основного РП в соответствии с функцией расстановки $h(x) = x \bmod Q_0$. Если указанная позиция пуста, в ее первую зону заносится адрес 1-граммы x в анализируемом тексте (T_1 или T_2). Одновременно в счетчик, соответствующий T_1 (если $x \in T_1$) или T_2 (если $x \in T_2$), заносится единица. Если позиция уже занята, 1-грамма x сравнивается с 1-граммой z текста, адрес которой указан в первой зоне данной позиции. При $x = z$ добавляется единица в соответствующий счетчик, в противном случае осуществляется переход по адресу отсылки и сравнение продолжается.

Если поиск закончился неудачно, информация об 1-грамме x заносится в свободную позицию, адрес которой определяет указатель очередной свободной позиции W . В последний из просмотренных до этого элементов списка заносится адрес отсылки на новый элемент списка.

Если дополнительное поле не заполнено, переменная W дает номер первой свободной позиции в этом поле ($Q_0 + 1$, $Q_0 + 2$ и т.д. до N). Если дополнительное поле заполнено, начинается линейный просмотр основного РП с целью определения первой пустой позиции. Переменная W принимает значение, равное номеру этой позиции. При каждом последующем обращении за свободной позицией просмотр основного РП начинается с элемента с номером $W + 1$. Это гарантирует, что суммарная трудоемкость поиска свободных позиций в основном РП не превысит $O(Q_0)$.

Описанный метод устранения наложений является гибридным между списковым методом, использующим дополнительную память, и списковым методом, не использующим такую ("coalesced chaining" [6]). Средняя трудоемкость удачного поиска в первом методе $t_1 = 1 + \frac{\alpha}{2}$, а во втором $t_2 = 1 + \frac{1}{4}\alpha + \frac{1}{8\alpha}(e^{2\alpha} - 1 - 2\alpha)$ [6]. Первый метод является более быстрым ($t_1 = 1,5$ при $\alpha = 1$), но требует дополнительных затрат памяти, второй — более медленным ($t_2 \approx 1,8$ при $\alpha = 1$), но не требует дополнительной памяти.

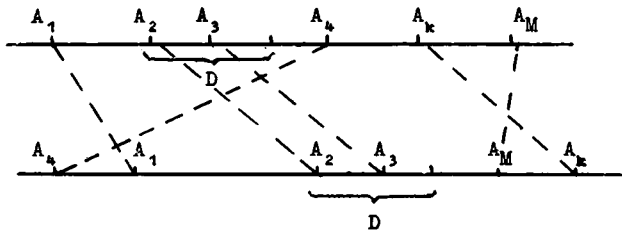
Ориентируясь при заказе памяти на наихудший случай ($q = N$), мы получаем возможность обрабатывать любые последовательности за один просмотр текстов. В то же время в подавляющем большинстве случаев такая память окажется избыточной ($M_1 - |X_{1-1}(T)|$ будет существенно меньше N), что делает целесообразным разбиение этой памяти на основную и дополнительную и обеспечивает наиболее быстрый поиск ($t = t_1$). Когда же на вход приходит "плохая" последовательность, дополнительное поле может оказаться заполненным и мы автоматически переходим на режим спискового метода без дополнительной памяти ($t \leq t_2$).

С учетом того, что средняя трудоемкость поиска одной 1-граммы $t_1 \leq t \leq t_2$, т.е. не зависит от N , трудоемкость 1-й итерации в среднем есть $O(N_1 + N_2)$. Мультипликативная константа C_1 в оценке трудоемкости определяет время обработки одной 1-граммы. В наихудшем случае — это время хеширования 1-граммы плюс затраты на сравнение пар 1-грамм при прохождении списка. Напомним, что отнюдь не все 1-граммы хешируются (их тем меньше, чем больше l), не все хешируемые подлежат сравнению (не сравниваются, к примеру, 1-граммы, попадающие в свободную позицию), сравнение не всегда идет до конца (1-граммы могут отличаться уже по первому символу).

В принципе зависимость от l может быть устранена почти полностью, если использовать свойство "перекрываемости" соседних 1-грамм и организовать рекуррентную схему вычисления хеш-адресов в виде $h_1(x_{i+1}) = f(h_1(x_i))$, где $h_1(x_i)$ — значение функции расстановки для i -й 1-граммы текста, f — таблично вычисляемая функция, характеризующая связь между числовыми кодами соседних 1-грамм.

5. О к о н ч а н и е и т е р а ц и й. Если обрабатываемые последовательности очень близки, то параметр L , определяющий длину максимального повтора из обоих текстов, может оказаться большим.

В этой ситуации нежелательно продолжать итерации до значения $l=L$. Целесообразно прекращать итерации при значении $l = L^* = \max\{l_{\max}(T_1), l_{\max}(T_2)\}$. Начиная с данного значения l , все 1-граммы, принадлежащие $X_1(T_1, T_2)$, имеют частоту встречаемости, равную 2, т.е. встречаются по разу в каждом из текстов. Обозначим их число через M . Картины их расположения внутри текстов можно условно представить следующей схемой:



Места расположения 1-грамм A_1, A_2, \dots, A_M определяются по нулевым кодам в ИЛ. При $l \geq L^*$ число нулей в ИЛ(T_1) и ИЛ(T_2) одинаково. Поскольку каждая из 1-грамм встречается в каждом из текстов по одному разу, подсчитывать их дальше по прежней схеме не имеет смысла. Требуется установить взаимно-однозначное соответствие между каждой парой эквивалентных 1-грамм (см. пунктирные линии на схеме) и расширить 1-граммы до появления первого несовпавшего символа. Если зона расширения D (см. схему) полностью охватывает (кроме расширяемой) еще какие-либо равноудаленные от нее 1-граммы из множества $\{A_1, A_2, \dots, A_M\}$, они оказываются автоматически обработанными.

Чтобы избежать сопоставления каждой пары 1-грамм из множеств $\{A_1, A_2, \dots, A_M | A_k \in T_1\}$ и $\{A_1, A_2, \dots, A_M | A_k \in T_2\}$, что потребовало бы затрат $O(M^2)$, можно вновь использовать процедуру ассоциативного кодирования, зафиксировав за один прогон все пары эквивалентных 1-грамм. Процесс обнаружения каждой пары эквивалентных 1-грамм (а это происходит при прогоне "второго" текста через ОП) может быть сразу совмещен с процедурой их расширения.

§3. Пример использования совместного частотного спектра

В работе [7] введена следующая мера сходства двух последовательностей:

$$\lambda(u, v) = \frac{\sum \min \{F(u:\alpha), F(v:\alpha)\} \cdot |\alpha|}{\sum \max \{F(u:\alpha), F(v:\alpha)\} \cdot |\alpha|} \quad (4)$$

Здесь $F(u:\alpha), F(v:\alpha)$ - частоты встречаемости произвольной 1-граммы α в текстах u и v соответственно, а $|\alpha| = 1$ - длина 1-граммы. Авторы, предложившие данную меру, не указывают алгоритма ее вычисления, но нетрудно показать, что информации, содержащейся в совместном частотном спектре $\Phi(u, v)$, достаточно для получения (4). Действительно, из определения $\Phi(u, v)$ следует, что в совместном частотном спектре содержатся те и только те 1-граммы, для которых $\min\{F(u:\alpha), F(v:\alpha)\} \neq 0$ (именно эти 1-граммы дают ненулевой вклад в числитель (4)). Таким образом, числитель выражения (4), который для краткости обозначим через S_{\min} , может быть получен линейным просмотром всех частотных характеристик совместного спектра.

Для вычисления знаменателя (4) представим каждое слагаемое в виде: $\max\{F(u:\alpha), F(v:\alpha)\} = F(u:\alpha) + F(v:\alpha) - \min\{F(u:\alpha), F(v:\alpha)\}$ и учтем, что

$$\sum_{\alpha} F(u:\alpha) \cdot |\alpha| = \frac{1}{6} N_u (N_u + 1) (N_u + 2), \quad \sum_{\alpha} F(v:\alpha) \cdot |\alpha| = \frac{1}{6} N_v (N_v + 1) (N_v + 2).$$

Тогда (4) запишется в виде:

$$\lambda(u, v) = \frac{S_{\min}}{\frac{1}{6} \{N_u (N_u + 1) (N_u + 2) + N_v (N_v + 1) (N_v + 2)\} - S_{\min}}$$

требуем для своего вычисления лишь знания $\Phi(u, v)$.

Меру (4), по-видимому, можно рекомендовать лишь для относительно коротких последовательностей, ибо для больших N она дает заниженные (с интуитивной точки зрения) значения.

§4. Заключение

I. Введено понятие совместного частотного спектра двух последовательностей, ориентированное на решение ряда классификационных задач (например, задач обнаружения знаков пунктуации в генетических текстах, вычисления эволюционного расстояния и т.д.). Предложен алгоритм вычисления совместного частотного спектра, ос-

нованный на применении процедуры хеширования. Использование хеширования приводит к очень простой логической структуре алгоритма, что, в свою очередь, обеспечивает малую мультипликативную константу в оценке трудоемкости. Этой же цели служат детально описанные вспомогательные структуры и приемы (использование ИЛ, чередование текстов, размещаемых в ОП и РП, стратегия поведения в случае заполнения дополнительного РП).

Понятие совместного частотного спектра допускает естественное обобщение на случай $m > 2$ последовательностей, однако "конструкция" счетчиков и стратегия работы с ними должна претерпеть существенные изменения (стать адаптивной).

2. Алгоритм может быть легко модифицирован для случая, когда лишь один из текстов помещается в ОП, а другой постоянно размещен на внешнем носителе. С такой ситуацией приходится сталкиваться в задачах информационного поиска, когда анализируемая последовательность должна сравниваться на предмет выявления локального сходства с гораздо более длинными текстами. В этом случае тексты не чередуются и размер РП определяется длиной короткого текста.

3. Возможны и другие (кроме рассматриваемых в § 2) адаптивные стратегии уменьшения числа итераций. Одна из них - выписывание в явном виде общих 1-грамм, как только число последних сократится до некоторого порога, затем лексикографическая сортировка их с последующим расширением групп эквивалентных 1-грамм.

Л и т е р а т у р а

1. ЗУБКОВ А.М., МИХАЙЛОВ В.Г. Предельные распределения случайных величин, связанных с длинными повторениями в последовательности независимых испытаний. - Теория вероятностей и ее приложения, 1974, т. XIX, № I, с. 173-181.
2. WEINER P. Linear pattern matching algorithms.-Proc. 14th Annual Symposium on Switching and Automata Theory, 1973, p.1-11.
3. MCCREIGHT E.M. A space economical suffix tree construction algorithm.-J.ACM, 1976, v.23, p.262-272.
4. ГУСЕВ В.Д., КОСАРЕВ Ю.Г., ТИТКОВА Т.Н. Отыскание статистических закономерностей текстов методом ассоциативного кодирования. - В кн.: Вычислительные системы. Вып. 62. Ассоциативное кодирование. Новосибирск, 1975, с. 72-89.
5. ГУСЕВ В.Д., КОСАРЕВ Ю.Г., ТИТКОВА Т.Н. О задаче поиска повторяющихся отрезков текста. - Там же, с. 49-71.
6. КНУТ Д. Искусство программирования для ЭВМ. Т.3. - М.: Мир, 1976. - 845 с.

7. FINDER N.V., LECUWEN J.van. A family of Similarity measures between two strings.- IEEE Trans.of Pattern Analysis and Mach.Intell., 1979, v.PAMI-1, N 1, p.116-118.

Поступила в ред.-изд. отд.
15 ноября 1983 года