



Math-Net.Ru

Общероссийский математический портал

А. И. Шкловский, Некоторые вопросы графической регулярности, *Зап. научн. сем. ЛОМИ*, 1971, том 20, 272–281

Использование Общероссийского математического портала Math-Net.Ru подразумевает, что вы прочитали и согласны с пользовательским соглашением
<http://www.mathnet.ru/rus/agreement>

Параметры загрузки:

IP: 18.97.9.173

25 марта 2025 г., 05:23:36



НЕКОТОРЫЕ ВОПРОСЫ ГРАФИЧЕСКОЙ РЕГУЛЯРНОСТИ ^{ж)}

I. Задача экстраполяции последовательностей (см., например, [1]) и связанные с ней задачи распознавания образов (см., например, [2], стр. 94 и [3]) естественным образом приводят к необходимости выявления и изучения внутренних связей (регулярностей) в имеющейся информации.

В связи с этим мы рассматриваем следующий вопрос. Пусть имеется некоторое слово (например, экстраполируемая последовательность). В какой мере можно сократить это слово путем введения новых "атомарных понятий" для его регулярных частей, т.е. для многократно встречающихся в нем достаточно длинных подслов. Точнее, мы рассматриваем следующий метод выявления регулярности: вместо непересекающихся вхождений слова B в слово X подставляется новый знак \mathcal{L} и перед получившимся словом приписывается слово $\mathcal{L} \rightarrow B^*$, где $*$ - разделительный знак, отличный от \mathcal{L} и букв слова X . Этот процесс выявления графической регулярности повторяется. Начальное слово X называется сократимым, если суммарная длина получившегося слова (включающего схемы подстановок) меньше длины X . В работе в основном исследуются характеристики наибольших по длине несократимых слов в зависимости от числа различных букв в этих словах. Использование нами аппарата порождающих бесконтекстных грамматик вызвано, во-первых, удобством его для решения поставленных задач и, во-вторых, имеющимся опытом использования этого аппарата для решения задач экстраполяции после-

ж) Основные результаты настоящей заметки были доложены на Ленинградском семинаре по конструктивной математике 2 октября 1969 года.

довательностей (см. [1]). Методы выявления регулярностей, аналогичные нашим, предлагались М.М.Бонгардом (см. [2], стр.209).

Основным результатом работы является нахождение для бесконтекстной грамматики точной верхней оценки длины несократимых слов в зависимости от числа различных букв в них. Кроме того, найдена точная нижняя оценка возможной величины сокращения слова.

Если рассматривать наши методы сокращения слов как методы кодирования, то они не дают уменьшения количества информации (в обычном для теории связи понимании) для всех слов, т.к. кодирование вызывает расширение алфавита, однако на достаточно регулярных словах выигрыш получается столь значительным, что превышает потери на расширение алфавита. Поэтому некоторые результаты работы могут быть интерпретированы как результаты теории связи. Но, как уже говорилось, основной нашей целью является рассмотрение методов выявления регулярностей для решения задач экстраполяции последовательностей, а эти задачи существенно отличаются от задач теории связи.

2. Порождающая грамматика - это четверка вида (V_T, V_N, P, S) (будем обозначать ее через G), где V_T - основной алфавит, не пересекающийся со вспомогательным алфавитом V_N ; P - список схем подстановок, т.е. пар вида (φ_i, ψ_i) , где φ_i и ψ_i - некоторые слова в алфавите $V_N \cup V_T$; S - аксиома, т.е. непустое слово в $V_N \cup V_T$. Обычно требуется, чтобы слово S было однобуквенным, но мы не накладываем ограничений на длину S .

Большими латинскими буквами из начала алфавита, возможно с индексами, будем обозначать буквы из V_N . Если U - алфавит, то через U^* будем обозначать множество всех слов в алфавите U .

Порождающая грамматика G называется бесконтекстной, если каждая схема из P имеет вид (A, ψ) , где ψ - непустое слово из $(V_N \cup V_T)^*$. Запись $\varphi_1 \rightarrow \varphi_2$ означает, что φ_1 представимо в виде $W_1 A W_2$, а φ_2 представимо в виде $W_1 \psi W_2$, где $(A, \psi) \in P$, а W_1 и W_2 принадлежат к $(V_N \cup V_T)^*$. Запись $\varphi = \psi$ означает, что су-

существует последовательность $\varphi_0, \dots, \varphi_n$ (называемая выводом Ψ из φ такая, что $\varphi_0 = \varphi$, $\varphi_n = \psi$ и $\varphi_{i-1} \vdash \varphi_i$, $i = 1, 2, \dots, n$. Множество $L(G) = \{X : X \in V_T^* \text{ и } S \equiv X\}$ будем называть языком, порожденным грамматикой G . В дальнейшем, говоря о какой-либо грамматике, будем иметь в виду бесконтекстную порождающую грамматику, язык которой состоит из одного слова.

3. Следующим образом определим запись $P : \{L, P, S\} = \varphi_e \psi_e^* \dots \varphi_1 \psi_1^*$, где ℓ - число схем подстановок в P . Через $|X|$ будем обозначать длину слова X .

$\tau(X)$ - мерой графической регулярности слова X - будем называть величину $\min_{L(G)=X} |z_{L, P, S}|$

Ту грамматику G , на которой достигается этот минимум, назовем оптимальной для X . Очевидно, что $\tau(X) \leq |X|$ для всякого X .

Пусть P состоит из ℓ пар $(A_1, \psi_1), \dots, (A_\ell, \psi_\ell)$. Будем говорить, что P (или G) имеет приведенный вид, если

- (а) $A_i \neq A_j$ при $i \neq j$
- (б) если $j > i$, то A_i не входит в ψ_j
- (в) $|\psi_i| \geq 2$ ($1 \leq i \leq \ell$)

Теорема I. Для всякого слова X существует оптимальная для него грамматика, имеющая приведенный вид.

Замечание. Не всякая грамматика, имеющая приведенный вид, является оптимальной.

4. Назовем слово X сократимым грамматикой G (грамматику G - сокращающей слово X), если $L(G) = X$ и $|z_{L, P, S}| < |X|$. Будем говорить, что слово X сократимо, если существует сокращающая его грамматика.

Через U^m будем обозначать слово $UU \dots U$, составленное из m слов U .

Определим оператор P_{τ} (оператор проекции) на всех словах вида X/Y , где Y непусто, следующим образом. Положим $P_{\tau}(X/Y) = Y^m$, где m - наибольшее число, для которого существует представление X в виде $w_1 Y w_2 \dots w_m Y w_{m+1}$.

Теорема 2. Слово X тогда и только тогда сократимо грамматикой, множество схем подстановок которой состоит только из одной пары когда

$$\frac{P_{\tau}(X/Y)}{|Y|} > \frac{|Y| + 2}{|Y| - 1} \quad (*)$$

Доказательство. Пусть выполнено условие (*). Положим $P = \{(A, Y)\}$. Тогда $|z_{\perp} P_{\perp}| = |Y| + 2$.

Пусть $P_{\tau}(X/Y) = Y^m$. Тогда X можно представить в виде $w_1 Y w_2 \dots w_m Y w_{m+1}$. Через S обозначим результат замены в этом слове всех выделенных вхождений Y на A . Нетрудно видеть, что

$$|S| = |X| - m(|Y| - 1).$$

Отсюда и из (*) получаем:

$$|z_{\perp} P_{\perp}| + |S| = |Y| + 2 + |X| - m(|Y| - 1) < |X|$$

Обратное утверждение доказывается с помощью аналогичных выкладок.

Теорема 3. Если слово X сократимо, то оно сократимо и некоторой грамматикой, множество схем подстановок которой состоит только из одной пары.

Доказательство. Пусть G - сокращающая X грамматика, оптимальная для X и имеющая приведенный вид. Пусть S - ее аксиома, а $(A_e, \psi_e), \dots, (A_1, \psi_1)$ - список ее схем подстановок. Через Z_e обозначим S . Результат применения (A_1, ψ_1) к S , т.е. результат подстановки в S вместо всех вхождений

A_1 слова ψ_1 , обозначим через Z_1 . Результат применения (A_2, ψ_2) к Z_1 обозначим через Z_2 и т.д. Тогда $Z_e = X$. Найдем наибольшее j такое, что

$$|Z_j| + |A_e \psi_e * A_{e-1} \psi_{e-1} * \dots * A_{j+1} \psi_{j+1} *| < |X|.$$

Поскольку

$$|Z_{j+1}| + |A_e \psi_e * \dots * A_{j+2} \psi_{j+2} *| \geq |X|,$$

то отсюда и из предыдущего неравенства получаем, что

$$|A_{j+1} \psi_{j+1}| + |Z_j| < |Z_{j+1}|.$$

Следовательно, должно быть выполнено условие (*) из теоремы 2 при $X = Z_{j+1}$, $Y = \psi_{j+1}$.

Применим к слову ψ_{j+1} всюду, где возможно, схему (A_{j+2}, ψ_{j+2}) , затем схему (A_{j+3}, ψ_{j+3}) и т.д. до схемы (A_e, ψ_e) включительно. Получившееся в результате слово обозначим через $\overline{\psi_{j+1}}$. Очевидно, что $|\overline{\psi_{j+1}}| \geq |\psi_{j+1}|$ и всякому вхождению ψ_{j+1} в Z_{j+1} соответствует вхождение $\overline{\psi_{j+1}}$ в X . Отсюда видно, что условие (*) выполнено при $Y = \overline{\psi_{j+1}}$. Представим X в виде $W_1 \overline{\psi_{j+1}} W_2 \dots \overline{\psi_{j+1}} W_{m+1}$ с наибольшим возможным m . В силу теоремы 2 слово X сократимо грамматикой, единственной схемой подстановок которой является $(A, \overline{\psi_{j+1}})$, а аксиомой - слово $W_1 A W_2 \dots A W_{m+1}$. Теорема доказана.

5. Рассмотрим последовательность расширяющихся алфавитов \mathcal{b} $\Delta_1, \Delta_2, \dots, \Delta_n, \dots$, где $\Delta_1 = \{a_1\}$, $\Delta_{n+1} = \Delta_n \cup \{a_{n+1}\}$.

Теорема 4. Если $X \in \Delta_n^*$ и $|X| > 4n^2 + 3n + 1$, то X сократимо.

Доказательство. Определим функцию $\text{def}_\ell(X)$, которую будем называть дефицитом ℓ -буквенных блоков в слове X , следующим равенством:

$$\text{деф}_\ell(X) = \max\left\{0, |X| - \ell + 1 - \frac{1}{\ell} \sum_{|Z|=\ell} |P_\tau(X/Z)|\right\}.$$

Понятие дефицита дает возможность учесть, например, различие между словами $X_1 = aaaaaaaaaa$ и $X_2 = ababababa$, заключающееся в том, что из X_1 можно выделить только 4 непересекающихся вхождения под слова aa , а из X_2 - либо 4 вхождения под слова ab , либо 4 вхождения под слова ba . В дальнейшем $\text{деф}_2(X)$ будем называть дефицитом в слове X и обозначать через $\text{деф}(X)$.

(а) Докажем, что каково бы ни было слово X , $\text{деф}(X) = 0$ тогда и только тогда, когда X не содержит ни одного под слова вида a_i^3 . Действительно, если X пустое, то $\text{деф}(X) = 0$, а если X непустое и не содержит ни одного такого под слова, то все двухбуквенные под слова, у которых вхождение последней буквы одного совпадает с вхождением первой буквы другого - различны. В этом случае

$$\sum_{|Y|=2} |P_\tau(X/Y)| = 2(|X| - 1) \quad \text{и, следовательно, } \text{деф}(X) = 0.$$

Аналогично доказывается и обратное утверждение.

(б) Непосредственным вычислением нетрудно убедиться, что

$$\text{деф}(a_i^k) = \left\lfloor \frac{k-1}{2} \right\rfloor.$$

(в) Пусть $X, Y \in \Delta_n^*$ и последняя буква слова X отлична от первой буквы слова Y . Всем, кроме одного, вхождениям двухбуквенных слов в слово XY соответствуют вхождения этих же слов либо в X , либо в Y . Единственное новое вхождение двухбуквенного под слова в XY появилось на границе X и Y . Но по условию это двухбуквенное слово состоит из разных букв и, значит, учтется при суммировании. Следовательно,

$$\sum_{|Z|=2} |P_\tau(XY/Z)| = \sum_{|Z|=2} |P_\tau(X/Z)| + \sum_{|Z|=2} |P_\tau(Y/Z)| + 2$$

и поэтому $\text{деф}(XY) = \text{деф}(X) + \text{деф}(Y)$.

(г) Докажем, что $\text{деф}(X) \leq 3n$ для любого несократимого X из Δ_n^* . Пусть X - несократимое слово из Δ_n^* и

$$X = Z_1 a_{i_1}^{l_1} Z_2 a_{i_2}^{l_2} \dots Z_{N+1} \quad (**),$$

где ни одно из слов Z_j ($1 \leq j \leq N+1$) не содержит подслов вида a_i^3 , все $l_k \geq 3$ ($1 \leq k \leq N$) и взяты максимально возможными, т.е. a_i^k не является ни последней буквой слова $Z_1 a_{i_1}^{l_1} \dots Z_k$ ни первой буквой слова $Z_{k+1} a_{i_{k+1}}^{l_{k+1}} \dots Z_{N+1}$.

В силу п(а) $\text{деф}(Z_j) = 0$ ($1 \leq j \leq N+1$). Отсюда и из п(в) вытекает, что

$$\text{деф}(X) = \sum_{k=1}^N \left[\frac{l_k - 1}{2} \right] \quad (***)$$

Будем говорить, что представление (**) выделяет под- слова $a_{i_k}^{l_k}$. Так как X несократимое, то в представлении (**) не может быть 3 или более выделенных вхождений подслов, содержащих букву a_i , т.к. иначе будет выполнено условие (*) из теоремы 2 с a_i^3 в качестве U . Следовательно, для всякой буквы a_i^3 из Δ_n имеет место один из трех случаев.

(г 1) Среди выделенных подслов нет ни одного, содержащего a_i^3 . В этом случае буква a_i не дает никакого "вклада" в величину $\text{деф}(X)$, вычисляемую по (**).

(г 2) Есть ровно одно выделенное подслово вида a_i^k ($k \geq 3$). Тогда, по условию (*) из теоремы 2 $\text{деф}(a_i^k) \leq 3$. Значит в (***) есть только одно слагаемое, обусловленное вычислением дефицита слова вида a_i^k , и оно не превосходит 3.

(г 3) Есть два выделенных подслова: a_i^k и a_i^p . Возможны только следующие значения k и p :

$$k=3, p=3;$$

$$k=4, p=3;$$

$$k=3, p=4;$$

$$k=4, p=4;$$

$$k=5, p=4;$$

$$k=4, p=5;$$

$$k=3, p=5;$$

$$k=5, p=3;$$

поскольку для наборов k и p , хотя бы по одному из параметров больших, чем приведенные в таблице, слово X будет сократимым, т.к. выполнится условие (*) из теоремы 2 с a_i^3 в качестве Y . Вычислив значения $\text{деф}(a_i^k) + \text{деф}(a_i^p)$ для k и p , перечисленных в таблице, нетрудно убедиться, что "вклад" буквы a_i в величину $\text{деф}(X)$ при вычислении ее по (***) , также не превосходит 3 .

Таким образом "вклад" всякой буквы a_i ($1 \leq i \leq n$) в величину $\text{деф}(X)$ не превосходит 3 и, значит, $\text{деф}(X) \leq 3n$. По условию (*) из теоремы 2 для всякого подслова Z при $|Z|=2$ выполняется неравенство $|P_z(X/Z)| \leq 8$. Значит,

$$\sum_{|Z|=2} |P_z(X/Z)| \leq 8n^2,$$

и поэтому $|X| \leq 4n^2 + 3n + 1$. Теорема доказана.

6. Докажем, что оценка, полученная в теореме 4, точна. Определим отношение строгой принадлежности слова к алфавиту: $X \in \Delta_n^*$ означает, что $X \in \Delta_n^*$ и всякая буква из Δ_n входит в X .

Теорема 5. Для любого положительного n можно построить такое несократимое слово X , строго принадлежащее к Δ_n^* , что $|X| = 4n^2 + 3n + 1$.

Доказательство. Требуемое слово, обозначаемое ниже посредством X_n , является результатом индуктивного построения, определяемого следующими (графическими) равенствами:

$$X_1 = a_1^3,$$

$$X_{i+1} = a_{i+1} X_i a_{i+1}^3 (a_i a_{i+1})^3 (a_i a_{i+1})^4 \dots (a_{i-1} a_{i+1})^4.$$

Индукцией нетрудно показать, что X_n имеет длину, указанную в формулировке теоремы. Проследив за процессом построения X_n , можно убедиться, что всякое двухбуквенное подслово X_n имеет только 4 непересекающихся вхождения в X_n и, следовательно, при

$|Y|=2$ имеем $|Pr(X/Y)| \leq 8$. Аналогично можно убедиться в том, что для любого слова $Y \in \Delta_n^*$ такого, что $|Y| \geq 3$ выполняется неравенство

$$Pr(X_n/Y) \leq \begin{cases} 6 & \text{при } |Y| = 3 \\ 8 & \text{при } |Y| = 4 \\ |Y| & \text{при } |Y| > 4 \end{cases}$$

Из полученных оценок для длины проекции X_n на различные слова и теоремы 2 вытекает, что X_n несократимо. Теорема доказана.

7. Следующая теорема дает достижимую нижнюю оценку для $\tau(X)$

Теорема 6. Для всякого слова X

$$\tau(X) \geq 3 \log_2 |X| - 2.$$

Для всякого i существует слово X такое, что $|X| > i$ и $\tau(X) = 3 \log_2 |X| - 2$.

Доказательство. Пусть G - оптимальная для X грамматика, имеющая приведенный вид, S - ее аксиома, $\{(A_1, \Psi_1), \dots, (A_\ell, \Psi_\ell)\}$ - множество схем подстановок. Обращаясь к доказательству теоремы 1, нетрудно видеть, что вывод можно проводить так: применить к аксиоме S везде, где только возможно, пару (A_1, Ψ_1) , затем к результату применить везде, где это возможно, пару (A_2, Ψ_2) и т.д. После применения последней пары (A_ℓ, Ψ_ℓ) результатом будет X . Следовательно,

$$|S| \cdot |\Psi_1| \dots |\Psi_\ell| \geq X.$$

Для всех натуральных a выполняется неравенство $a + 2 \geq 3 \log_2 a$. Следовательно,

$$\tau(X) = |S| + \sum_{i=1}^{\ell} (|\Psi_i| + 2) \geq 3 \log_2 |S| -$$

$$- 2 + \sum_{i=1}^{\ell} 3 \log_2 |\Psi_i| = 3 \log_2 (|S| \cdot |\Psi_1| \dots |\Psi_\ell|) -$$

$$- 2 \geq 3 \log_2 |X| - 2.$$

Для доказательства второго утверждения теоремы рассмотрим слова X_i , определяемые равенством $X_i = a^{(4^i)}$. Построим грамматику, языком которой является X_i .

$$V_T = \{a\}. \quad V_N = \{a_0, a_1, \dots, a_{i-1}\}$$

$$S = a_0 a_0 a_0 a_0$$

$$P = \{(a_0, a_1^4); (a_1, a_2^4); \dots; (a_{i-1}, a^4)\}.$$

Нетрудно видеть, что

$$\tau(X_i) = |S| + |P| = 4 + 6 \cdot (i-1) = 3 \cdot 2i - 2,$$

что и требовалось доказать.

Автор благодарен Ю.В.Матиясевичу и А.О.Слисенко за обсуждения работы и полезные советы.

ЛИТЕРАТУРА

1. Solomonoff R.I., A Formal Theory of Inductive Inference, "Information and Control", 1964, 7, №1, 1-22; 224-254.
2. Бонгард М.М., Проблема узнавания. М., 1967.
3. Загоруйко Н.Г., Самохвалов К.Ф. Природа проблемы распознавания образов. В сб. "Вычислительные системы". 1966 г. вып. 36, 3-19.