

УДК 62-5:007:621.391

АВТОМАТИЧЕСКОЕ РАСПОЗНАВАНИЕ ОГРАНИЧЕННОГО
НАБОРА УСТНЫХ КОМАНД

В.И.Величко, Н.Г.Загоруйко

Задача автоматического распознавания ограниченного набора устных команд привлекает внимание исследователей речевого сигнала, главным образом, с точки зрения получения практически приемлемых решений таких частных задач, как управление голосом автоматическими устройствами, устный ввод исходных данных и программ в ЭВМ и т.д. (см., например, [1]).

В данной статье описываются эксперименты по распознаванию словаря из 168 слов (таблица I). В качестве исходного материала взяты все термины языка "Алгол-60", дополненные некоторыми терминами входного языка α -транслятора и названиями некоторых элементарных функций. Большинство названий букв русского и латинского алфавитов заменено собственными именами (как это принято в системах связи); кроме того, вместо "ми" используется "ми" и вместо "пси" - "пси-штрих".

Для описания речевого сигнала в данном эксперименте была принята следующая система признаков. Акустический сигнал проpusкался через систему из 5 фильтров второго порядка (резонансные контуры LCR). Центральные частоты фильтров были выбраны равными 112,5 гц; 450 гц; 900 гц; 1800 гц; 7200 гц [2]. Полоса пропускания каждого фильтра по уровню 0,5 равнялась центральной частоте, что соответствует добротности фильтров $\sim 2,45$.

Слово разбивалось на сегменты фиксированной длительности $T = 4$ мсек. Длительность сегмента выбрана превышающей максимально возможный период основного тона. В каждом сегменте подсчитывалась энергия сигнала в общей полосе E_0 и на выходе каждого из 5 фильтров E_i ($i = 1, \dots, 5$). В качестве параметров, ха-

Таблица I

| | |
|---------------------|-----------------------------|
| I. Один | 40. Влево |
| 2. Два | 41. Вправо |
| 3. Три | 42. Вопросительный |
| 4. Четыре | 43. Вниз |
| 5. Пять | 44. Плюс-минус |
| 6. Шесть | 45. Набла |
| 7. Семь | 46. Альфа |
| 8. Восемь | 47. Бета |
| 9. Девять | 48. Гамма |
| I0. Ноль | 49. Дельта |
| I1. Плюс | 50. Эпсилон |
| I2. Минус | 51. Дзета |
| I3. Разделить | 52. Эта |
| I4. Запятая | 53. Тата |
| I5. Точка | 54. Каппа |
| I6. Пробел | 55. Лямбда |
| I7. Десять | 56. Ми |
| I8. Степень | 57. Ну |
| I9. Скобка | 58. Кси |
| 20. Открыть | 59. Пи |
| 21. Закрыть | 60. Ро |
| 22. Круглая | 61. Сигма |
| 23. Умножить | 62. Тау |
| 24. Равно | 63. Фи |
| 25. Точка с запятой | 64. Хи |
| 26. Квадратная | 65. Пси-итрих |
| 27. Звездочка | 66. Омега |
| 28. Кавычки | 67. Аш |
| 29. Не | 68. Исправление и выделение |
| 30. Больше | 69. Примечание |
| 31. Меньше | 70. Целый |
| 32. Двоеточие | 71. Вещественный |
| 33. Надчеркнуть | 72. Комплексный |
| 34. Подчеркнуть | 73. Логический |
| 35. Восклицательный | 74. Собственный |
| 36. Твердый | 75. Начало |
| 37. Градус | 76. Конец |
| 38. Минута | 77. Масса |
| 39. Секунда | 78. Стоп |

| | |
|-----------------------|-----------------|
| 79. Для | II7. Арктангенс |
| 80. Шаг | II8. Логарифм |
| 81. До | II9. Экспонента |
| 82. Цикл | I20. Корень |
| 83. Рав | I21. Русская |
| 84. Пока | I22. Латинская |
| 85. Если | I23. Греческая |
| 86. То | I24. Большая |
| 87. Иначе | I25. Малая |
| 88. На | I26. Аля |
| 89. Процедура | I27. Борис |
| 90. Результат | I28. Владимир |
| 91. Значение | I29. Геннадий |
| 92. Метка | I30. Дмитрий |
| 93. Барабан | I31. Елена |
| 94. Лента | I32. Женя |
| 95. Скаляр | I33. Зинаида |
| 96. Вектор | I34. Иван |
| 97. Матрица | I35. И краткая |
| 98. Страна | I36. Клим |
| 99. Функция | I37. Лидия |
| I00. Истина | I38. Михаил |
| I01. Ложь | I39. Николай |
| I02. После | I40. Ольга |
| I03. Печать | I41. Павел |
| I04. Заменить | I42. Роман |
| I05. Число повторений | I43. Сергей |
| I06. Идентификатор | I44. Татьяна |
| I07. Конец отладки | I45. Ульяна |
| I08. Переключатель | I46. Фёдор |
| I09. Присвоить | I47. Харитон |
| I10. Дифференциал | I48. Цецилия |
| III. Интеграл | I49. Чарли |
| I12. Троеточие | I50. Шура |
| I13. Фигурная | I51. Шукарь |
| I14. Синус | I52. И |
| I15. Косинус | I53. Мягкий |
| I16. Арксинус | I54. Знак |

Продолжение таблицы I

- | | |
|--------------|--------------------|
| I55. Эдуард | I62. Дизъюнция |
| I56. Юрий | I63. Конъюнция |
| I57. Яков | I64. Импликация |
| I58. Жи | I65. Целое деление |
| I59. Ку | I66. От |
| I60. Дубльва | I67. Тождественно |
| I61. Зэт | I68. Отрицание |

рактеризующих сегмент, брались величины $\ln \frac{E_0}{E_i}$, т.е. сегмент описывался точкой в пятимерном пространстве. Слово описывалось как последовательность сегментов. Частичное обоснование выбора логарифмических параметров состоит в том, что слух подчиняется известному в физиологии закону Вебера-Фехнера [3], согласно которому минимально заметный прирост ΔJ внешнего воздействия на органы чувств пропорционален величине воздействия J

$$\Delta J \sim \Delta \beta J,$$

откуда $\beta \sim \ln J$ — сила восприятия, пропорциональная логарифму внешнего воздействия.

Для сравнения различных сегментов выбрана мера сходства

$$\alpha = \frac{\alpha^2}{\alpha^2 + \rho^2}. \quad (I)$$

Здесь

$$\rho^2 = \sum_{i=1}^5 \left[\ln \frac{E_0^{(i)}}{E_i^{(i)}} - \ln \frac{E_0^{(k)}}{E_k^{(i)}} \right]^2$$

—квадрат евклидова расстояния между сегментами в пятимерном пространстве, α^2 принято равным 2 на основе предыдущих экспериментов [4]. Мера сходства α может принимать значения от 1 (для одинаковых сегментов) до 0 (для бесконечно удаленных).

Чтобы перейти от меры сходства α для сегментов к мере сходства λ для слов, использована идея, изложенная в работе [5]. Пусть сравниваются две реализации одного и того же слова "восемь" (рис. I). По координатным осям отложим время от начала слова и условно отметим границы фонем. Слова из-за различия в произношении имеют, как правило, разную временную структуру — одни фонемы удлиняются, другие укорачиваются, причем в разной степени. Интуитивно очевидно, что для получения максималь-

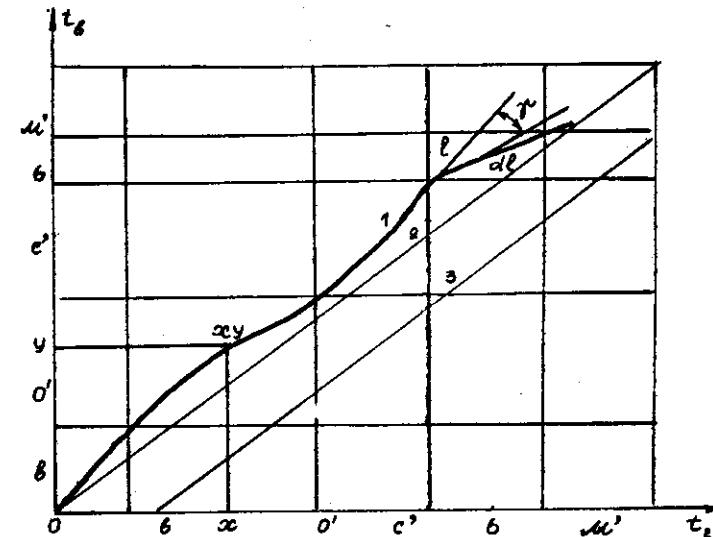


Рис. I.

ной меры сходства между одинаковыми словами (не определяя пока понятия меры сходства) желательно сравнивать похожие участки слов: "в" — "в", "о" — "о" и т.д. Пусть момент времени X первой ("горизонтальной") реализации слова соответствует моменту времени Y второй ("вертикальной") реализации того же слова. Тогда сходство между параметрами "горизонтального" и "вертикального" слов в этих точках будет велико (например, в смысле меры сходства (I)). Построим в плоскости t_1, t_2 геометрическое место точек, подобных точке X, Y кривую I. Каждая точка этой кривой проецируется в "похожие" точки "горизонтального" и "вертикального" слов, а вся кривая является как бы "кривой максимального сходства". Пусть $B(\ell)$ — мера сходства между параметрами слов в момент, соответствующий точке ℓ кривой I. Тогда мера сходства B между словами можно определить как сумму мер сходства для всех точек кривой I (в непрерывном приближении сумма заменяется интегралом):

$$B = \int B(\ell) d\ell. \quad (2)$$

Выражение (2) может быть скорректировано с учетом общей длины кривой I и наклона кривой I в различных точках. Ясно, что чем длиннее слова, тем больше получается B . Такое влияние длины слова на величину сходства нежелательно. Далее, чем больше на-

клон кривой в какой-то точке отличается от $\pi/4$, тем больше разница в темпе произнесения слов, что должно быть учтено при вычислении меры сходства. Поэтому может быть принято уточненное определение меры сходства между словами:

$$B = \frac{1}{L} \int f(r) B(\ell) d\ell , \quad (3)$$

где L - длина кривой I;

$f(r)$ - убывающая функция угла $|f|$ между биссектрисой координатного угла и кривой I, например, $\cos 2r$.

Мера сходства B существенно зависит от кривой, вдоль которой она вычисляется. Нам заранее не известна "кривая максимального сходства", поэтому в соответствии с [5] определим B как максимум функционала (3) вдоль всех возможных кривых ℓ , соединяющих начала и концы сравниваемых слов. На кривые наложено единственное ограничение - они должны монотонно возрастать, т.е. $-\frac{\pi}{4} \leq r \leq \frac{\pi}{4}$ (4)

Мера сходства (4) введена для непрерывного случая. Её вычисление сводится к решению вариационной задачи нахождения максимума функционала B . Нами принято дискретное описание слова в виде последовательности точек шестимерного пространства параметров. Переформулируем задачу для дискретного случая. Построим матрицу $\{\alpha_{ik}\}$, элементами которой являются меры сходства (в смысле формулы (1)) между сегментами сравниваемых слов:

$$\alpha_{ik} = \frac{\alpha^2}{\alpha^2 + \rho_{ik}^2} \quad (5)$$

Здесь ρ_{ik} - квадрат евклидова расстояния между k - м сегментом первого ("горизонтального") и i - м сегментом второго ("вертикального") слова. Нумерация строк матрицы совпадает с нумерацией сегментов "вертикального" слова, нумерация столбцов - с нумерацией сегментов "горизонтального" слова (см. рис.2). Матрица $\{\alpha_{ik}\}$ является дискретным аналогом непрерывного поля мер сходства $B(\ell)$ (рис.1). Точка A на рис.2 соответствует началу координат на рис.1. Введем правило вычисления функционала сходства F между словами по матрице $\{\alpha_{ik}\}$: функционал F вдоль пути $ABC...PQRST$ есть сумма тех элементов матрицы $\{\alpha_{ik}\}$, через которые проходит выбранный путь, деленная на длину более длинного слова. В приведенном примере

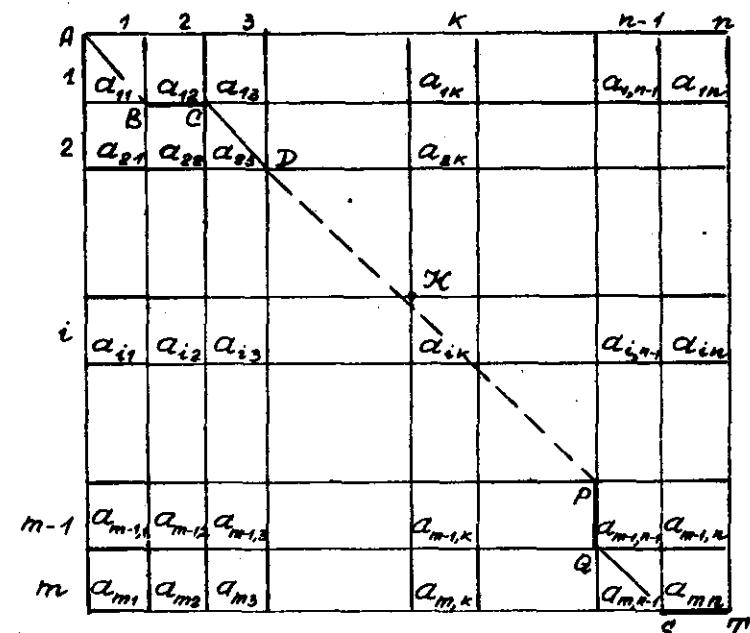


Рис.2.

при $n > m$ это будет сумма $\frac{1}{n} (\alpha_{11} + \alpha_{12} + \dots + \alpha_{m,n-1})$, соответствующая участкам пути AB, CD, \dots, QS . Участки пути BC, \dots, PQ, ST проходят между элементами матрицы и на сумму не влияют. На возможные пути накладывается ограничение: разрешенными являются только участки пути по горизонтали слева направо, по вертикали сверху вниз и по диагонали слева-справа направо-вниз.

Мерой сходства A между словами будем считать максимум функционала F , вычисленного вдоль всех возможных путей на матрице $\{\alpha_{ik}\}$: $A = \max_{\text{путь}} F$. Нетрудно видеть, что матрица $\{\alpha_{ik}\}$ представляет собой граф или сеть (см., например, [6]), а задача максимизации функционала сходства есть задача нахождения максимального пути на графе. В данном эксперименте длина максимального пути вычислялась методом динамического программирования (см., например, [7]). Ниже приведено краткое описание примененного алгоритма.

Пронумеруем узловые точки матрицы $\{\alpha_{ik}\}$ (т.е. точки A, B, \dots, T) по следующему правилу. Точка, лежащая левее и выше элемента α_{ik} , имеет номер (i, k) . На рис. 2 точка X имеет номер (L, K) , точка A - $(1, 1)$, B - $(2, 2)$ и т.д. Точки, лежащие на правой границе матрицы, имеют номера $(i, n+1)$, на нижней границе -

($m+1, n$), где $i = 1, \dots, m+1$; $k = 1, \dots, n+1$. Длина пути A_{ik} от точки (i, k) до точки $T(m+1, n+1)$ есть максимальное из следующих трех чисел: $A_{i,k+1}, A_{i+1,k}$ и $A_{i,k} + A_{i+1,k+1}$. Требуется найти длину пути от точки (1,1) до ($m+1, n+1$), т.е. $A_{m,n}$. Границные условия задаются в следующем виде: $A_{i,n+1} = A_{m+1,k} = 0$. Вычисления ведутся последовательно, начиная с точки (m, n):

$$1) A_{m,n} = \max[A_{m,n+1}, A_{m+1,n}, A_{mn} + A_{m+1,n+1}] = \alpha_{mn},$$

$$2) A_{m-1,n} = \max[A_{m-1,n+1}, A_{m,n}, \alpha_{m-1,n} + A_{m,n+1}] = \max(\alpha_{mn}, \alpha_{m-1,n})$$

$$3) A_{m-i,n} = \max[A_{m-i+1,n+1}, A_{m-i+2,n}, A_{m-i,n} + A_{m-i+1,n+1}]$$

$$m) A_{1,n} = \max[A_{2,n}, A_{1,n+1}, \alpha_m + A_{2,n+1}]$$

$$m+1) A_{m,n-1} = \max[A_{m,n}, A_{m,n-1}, A_{mn} + A_{m+1,n}] = \max(\alpha_{mn}, \alpha_{m,n-1})$$

$$m,n) A_{11} = \max[A_{2,1}, A_{1,2}, \alpha_{11} + A_{2,2}]$$

На последнем шаге вычислений мы получаем требуемую длину пути A_{11} . После нормировки по длине слова, т.е. после деления A_{11} на длину более длинного слова, получаем λ — меру сходства между словами.

Вышеописанный алгоритм вычисления точной меры сходства между словами довольно громоздок и требует большого расхода машинного времени — количество операций пропорционально квадрату длины слов n^2 . Для значительного сокращения времени счета применён комбинированный метод принятия решения [8], в котором сначала вычисляются приближенные меры сходства между словами. Приближенная мера сходства определялась следующим образом: совмещаются начала и концы сравниваемых слов и суммируются (с последующей нормировкой) меры сходства между сегментами, одинаково удаленными от начала слов. Это соответствует нахождению функционала вдоль прямой 2 на рис. 1. Для корректировки возможных ошибок при автоматическом определении границ слова эта процедура повторяется со словами, сдвинутыми друг относительно друга на некоторое число сегментов (см. прямую 3 на рис. 1). В нашем эксперименте сдвиг начала одного слова относительно другого составлял -8, -4, 0, +4, +8 сегментов. Из полученных мер сходства истинной считалась максимальная.

Полностью алгоритм распознавания выглядел следующим об-

разом. Предъявленное для распознавания слово сравнивалось приближенным методом с каждым из 168 слов-эталонов обучающей последовательности. Выбирались 16 слов-эталонов, наиболее похожих на предъявленное слово. Затем методом динамического программирования проводилось точное определение меры сходства предъявленного слова с каждым из 16 эталонных слов и наиболее похожее слово-эталон выбиралось как результат распознавания. Комбинированная схема распознавания позволяет существенно сократить время принятия решения при незначительном снижении надежности распознавания.

Статистические испытания проводились на следующем материале: один диктор (мужчина), 4 последовательности слов, из которых 3 включают по 168 слов, а 4-я — 166 слов (при вводе 4-й последовательности в ЭВМ были допущены ошибки в двух словах). Весь материал был записан на бытовом магнитофоне "Астра-4" с микрофона типа МД-47 в тихой комнате с обычной акустикой. Все 4 последовательности с магнитофона той же марки с помощью 9-разрядного аналого-цифрового преобразователя (частота квантования 20 кГц) были введены в ЭВМ БЭСМ-6, где производилось автоматическое определение границ слов, фильтрация с помощью аппарата Z — преобразования (см., например, [9]), выделение логарифмических параметров слов и запись этих параметров на магнитную ленту БЭСМ-6. Обработка слова продолжается в 4 раза дольше длительности слова. Затем производилось распознавание, причем в качестве обучающих (эталонных) последовательностей использовались поочередно 1-я, 2-я и 3-я последовательности, а в качестве контрольных — все остальные. Результаты эксперимента приведены в таблице 2. Общая надежность распознавания, усредненная по 1506 реализациям, составляет выше 95%.

Таблица 2

| №пп №оп | Число ошибок | | | | надежность, % |
|------------|--------------|----|---|----|------------------|
| | 1 | 2 | 3 | 4 | |
| 1 | | 3 | 9 | 12 | 95,6 |
| 2 | 3 | | 8 | 5 | 96,8 |
| 3 | 16 | 13 | | 4 | 93,4 |

Распознавание одного слова на ЭВМ БЭСМ-6 занимает в среднем около 8 сек без учета времени работы внешних устройств (печать, магнитные ленты, магнитные барабаны) и около 10 сек с учетом этих устройств.

Результаты экспериментов показывают, что с помощью данного метода можно получить практически приемлемую надежность распознавания 150-200 слов, если имеется возможность производить подстройку системы под диктора, состоящую в однократном произнесении всех слов словаря.

Дальнейшие исследования должны быть направлены, главным образом, на сокращение времени принятия решения.

Л и т е р а т у р а

1. Г.Я.ВЫСОЦКИЙ, Б.Н.РУДНЫЙ, В.Н.ТРУНИН-ДОНСКОЙ, Г.И.ЦЕМЕЛЬ . Алгоритм опознавания 40 слов на ЦВМ БЭСМ-3М. В сб. "Работы по технической кибернетике", М., вып.2, ВЦ АН СССР, 1968, стр.3-33.
2. Б.М. КУРИЛОВ, Б.П.ГАВРИЛКО. Членение речевого потока на фонемы. Данный сборник, стр.
3. Физический энциклопедический словарь, т.1, 1960, стр.242.
4. Н.Г. ЗАГОРУЙКО, В.И.ВЕЛИЧКО, Г.Я.ВОЛОШИН, В.Д.ГУСЕВ, В.Н.ЕЛКИНА, И.В.БАХМУТОВА, А.Г.ХАЙРЕТИНОВА, Л.С.ЮДИНА.Эксперименты по автоматическому распознаванию речевых сигналов.-Доклад на ІУ школе-семинаре по автоматическому распознаванию слуховых образов (АРСО-ІУ). Канев, 1968.
5. А.А.ПИРОГОВ, Г.С.СЛУЦКЕР. К фонетической теории речи. Доклад на VI Всесоюзной акустической конференции, М., 1968.
6. А.КОФМАН, Г.ДЕБАЗЕЙ. Сетевые методы планирования и их применение. Изд. "Прогресс", 1968.
7. Е.С.ВЕНТИЦЕЛЬ. Элементы динамического программирования.Изд. "Наука", 1964.
8. Н.Г.ЗАГОРУЙКО. Комбинированный метод принятия решения.Вычислительные системы. Новосибирск, 1966, вып.22, Изд."Наука", Сиб.отд., стр.
9. Приспособливающиеся автоматические системы.Под ред.Э.Мишкина и Л.Брауна.Ил., 1963.

Поступила в редакцию
6.1.1969г.