

*f / . .*

*Vi*

На правах рукописи

УДК 519.95

**МАЗУРЕНКО Иван Леонидович**

**АВТОМАТНЫЕ МЕТОДЫ  
РАСПОЗНАВАНИЯ РЕЧИ**

01.01.09 - дискретная математика и математическая кибернетика

**ДИССЕРТАЦИЯ**  
**на соискание ученой степени**  
**кандидата физико-математических наук**



Научный руководитель:  
доктор физико-математических наук, профессор БАБИН Д.Н.

Москва-2001

# СОДЕРЖАНИЕ

ВВЕДЕНИЕ.....	
---------------	--

3

ГЛАВА 1. ФОРМАЛЬНАЯ МОДЕЛЬ РЕЧИ.....	26
--------------------------------------	----

1.1 Речь как математический объект. Основные определения..	26
--	----

1.2 Дискретизация .....	37
-------------------------	----

1.3 Формализация задачи распознавания речи. Функция качества словаря команд .....	51
---	----

ГЛАВА 2. ДЕТЕРМИНИРОВАННЫЕ АВТОМАТНЫЕ МОДЕЛИ.....	57
---	----

2.1 Метод решения задачи локального распознавания речи.....	57
---	----

2.2 Модель распознавания последовательного кода команды с помощью детерминированных автоматов.....	62
--	----

2.3 Модель и алгоритм распознавания кода команды для случая классов звуков. Нижняя оценка качества словаря команд...	66
--	----

ГЛАВА 3. ВЕРОЯТНОСТНЫЕ АВТОМАТНЫЕ МОДЕЛИ .....	72
--	----

3.1 Понятие монотонного автономного вероятностного автомата. Модель распознавания речи с помощью вероятностных автоматов .....	72
--	----

3.2 Метрика $p_1$ на множестве стохастических словарных функций .....	85
---	----

3.3 Метрика $p_2$ на множестве стохастических словарных функций. Функция качества словаря команд .....	87
--	----

3.4 Метрика $p_3$ на множестве стохастических словарных функций .....	95
---	----

3.5 Связь между метрикой и вероятностью.....	100
--	-----

ПРИЛОЖЕНИЕ 1. ПРИМЕР АЛФАВИТА БУКВ И ЗВУКОВ, ПРАВИЛ ТРАНСКРИБИРОВАНИЯ ДЛЯ РУССКОГО ЯЗЫКА .....	104
--	-----



ПРИЛОЖЕНИЕ 2. ПРИМЕР ФОНЕМНЫХ КЛАССОВ ДЛЯ РУССКОГО ЯЗЫКА .....	113
ПРИЛОЖЕНИЕ 3. МЕТРИКА НА МНОЖЕСТВЕ ЗВУКОВ РУССКОГО ЯЗЫКА .....	114
ЛИТЕРАТУРА.....	115
ПУБЛИКАЦИИ АВТОРА ПО ТЕМЕ ДИССЕРТАЦИИ.....	119

## ВВЕДЕНИЕ

Образное восприятие - одно из свойств мозга, позволяющее правильно воспринимать информацию о внешнем мире. При этом происходит классификация воспринимаемых ощущений, т.е. разделение их на группы похожих образов. Образы обладают тем свойством, что ознакомление с некоторым конечным числом их реализаций дает возможность узнавать все остальные представители образа. Более того, разные люди, обучающиеся на различном материале наблюдений, одинаково распознают одни и те же объекты.

Процессу распознавания образов всегда предшествует процесс обучения, во время которого мы знакомимся с некоторым числом примеров, зная наперед об их принадлежности к каким-то образам. У человека процесс обучения заканчивается успешно. Следовательно, образы - это объективные характеристики внешнего мира, а не произвольные наборы изображений. Эту объективную характеристику образов и стремятся выявлять с помощью математических методов. Создаются аппаратные и программные распознающие системы, на вход которых подается информация о предъявляемых объектах, а на выходе отображается информация о классах, к которым отнесены распознаваемые объекты.

Две фундаментальные гипотезы, лежащие в основе работы большинства систем автоматического распознавания речи (АРР), заключаются в следующем ([1, 6]):

- 1) информация в речевом сигнале переносится в виде изменений во времени его амплитудного спектра

- 2) речь - это сложный иерархически организованный сигнал, так что более простые образы одного уровня однозначно определенными правилами объединены в более сложные образы следующего уровня (аллофоны, фонемы, дифоны, трифоны, слоги, слова, предложения и т.п.)

Все это позволяет эффективно использовать автоматные методы[34] для описания речевых сигналов.

При этом решаются задачи синтеза порождающих образ автоматов (детерминированных и вероятностных), распознавания объектов, оценки сложности этих алгоритмов. Заметим, что участвующие в задании объектов массивы чисел хотя и конечные, но необозримые, и не представляется возможным хранение достаточного запаса примеров этих объектов. К примеру, одна секунда речевой записи занимает около 10000 байт, а сами примеры этих записей, как бы мы их не накапливали, не только не повторяются и не обладают близостью в 10000-мерном пространстве, но даже имеют принципиально разную длину.

Практически все известные методы распознавания речи обладают рядом основных общих свойств:

- для распознавания используется метод сравнения с эталонами;
- сигнал может быть представлен либо в виде непрерывной функции времени, либо в виде слова в некотором конечном алфавите;
- для сокращения объема вычислений используются методы динамического программирования.

Методы распознавания речи можно условно разделить на две большие группы — непараметрические — с использованием непараметрических мер близости к эталонам (к ним можно отнести методы на основе формальных грамматик и методы на основе метрик на множестве речевых сигналов) и параметрические (вероятностные - на основе скрытых марковских моделей, нейросетевые).

Первые разработки в области создания устройств автоматического распознавания речи можно отнести к концу 40-х годов прошлого столетия. Первые устройства ([2, 3, 4]) были аналоговыми и использовали пороговую логику. Исследователи сразу же столкнулись с тем, что распознавание речи — достаточно сложная задача в силу неравномерности

растяжения речевых сигналов во времени и сложной связи сигнала с фонемной структурой произносимой речи. В силу этого первые системы распознавания речи не обладали высокой надежностью и были узкоспециализированными.

Ситуация резко изменилась к началу 70-х годов в результате послевоенного развития электроники и вычислительной техники и появления лингвистической теории речи, представляющей устную речь как производную фонетической транскрипции произносимого текста. К числу первых проектов такого рода можно отнести проект Агентства перспективных исследовательских работ ARPA ([5]). В этих системах сигнал разбивался на сегменты, которые обозначались фонемоподобными символами, а затем получившийся код сигнала распознавался с помощью формальных грамматик. Узким местом такого подхода оказалось то, что сегментация и локальное распознавание является задачей, трудно поддающейся точному автоматическому решению, несмотря на то, что человек с этой задачей вполне справляется ([7]).

Следующим этапом в развитии теории распознавания речи стало развитие непараметрических систем, основанных на мерах близости на множестве речевых сигналов как функций времени. Революционный вклад в развитие этого направления оказал подход Винцюка (1960-е г.г., [8]), предложившего использовать новый для того времени метод динамического программирования (Беллман, 1950-е г.г. [9]) для быстрого вычисления меры близости между двумя функциями, задающими изменение во времени параметров речевых сигналов. Подход Винцюка, развитый Итакурой ([10]) и другими исследователями, позволил сократить время вычисления значений функции близости к эталонам с экспоненциального (от длины сигнала) до квадратичного. В силу того, что основной спецификой метода являлось нелинейное искажение временной оси одной из сравниваемых функций, метод получил название «динамической

деформации времени» (ДДВ).

Метод ДДВ обладает рядом достоинств и недостатков. К очевидным достоинствам метода относятся простота его реализации и обучения - для работы метода достаточно в качестве эталона команды использовать одно из ее произнесений. Основным недостатком метода является сложность вычисления меры близости (пропорционально квадрату длины сигнала) и большой объем памяти, необходимый для хранения эталонов команд (пропорционально длине сигнала и количеству команд в словаре).

Ряд исследователей (Бейкер — CMU — система «Драгон» [11] и, независимо, Йелинек — IBM [12], 1970-е годы) для распознавания речи применили теорию скрытых марковских моделей (СММ, скрытых марковских процессов, вероятностных функций цепей Маркова), созданную Баумом и коллегами в конце 60-х - начале 70-х г.г ([13]). Скрытые марковские процессы представляют из себя дважды стохастические процессы — марковские цепи ([14]) по переходам между состояниями и множества стационарных процессов в каждом состоянии цепи. Основы теории СММ были опубликованы в ряде научных журналов ([13,19]), однако получили развитие среди разработчиков систем распознавания речи лишь после выхода серии обзоров, посвященных популярному изложению теории СММ ([15-18]). Для обучения моделей и вычисления функции близости к эталону (здесь — вероятности наблюдения слова на выходе СММ) был также применен метод динамического программирования (алгоритмы прямого-обратного хода [19], Баума-Уэлча, или ЕМ-алгоритм [20], Виттерби [21-22]).

Достоинствами метода СММ являются достаточно быстрый способ вычисления значений функции расстояния (вероятности) и существенно меньший объем памяти, необходимый для хранения эталонов (пропорционально количеству фонем, Трифонов и т.п. в языке), а основны-



ми недостатками — достаточно большая сложность его реализации, а также необходимость использования больших фонетически сбалансированных речевых корпусов (баз данных) для обучения параметров СММ. Последнее обстоятельство привело к тому, что вплоть до начала 90-х годов методы ДДВ были более эффективными и точными при решении практических задач, чем методы СММ.

В настоящее время распознавание речи на основе СММ переживает период бурного роста. Десятки крупных коммерческих компаний (IBM, Dragon, L&N, Philips, Microsoft, Intel...) создали и активно развивают коммерческие системы распознавания речи. Эксперты в области компьютерных технологий называют распознавание речи одной из важнейших задач XXI века.

Основными характеристиками современных систем АРР являются следующие:

- словари размером в десятки тысяч слов;
- распознавание слитной речи;
- возможность работы как с предварительной настройкой на голос диктора, так и без настройки;
- надежность работы 95-98%.

СММ, возникшие как обобщение цепей Маркова, тесно связаны с понятием вероятностного автомата. Вероятностные автоматы, впервые введенные в общей форме Дж Карлайлом (1963, [23]) и, независимо от него, Р.Бухараевым (1964, [24]) и П.Штарке (1965, [25]), представляют из себя в практическом плане устройства с конечной памятью, перерабатывающие информацию с входных каналов в выходные, переходы и выходы которых происходят на основе вероятностных законов ([26]). Скрытые марковские модели являются частным случаем вероятностных автоматов, а именно, вероятностными автоматами без входа. СММ, используемые в системах распознавания речи, обладают дополнительно тем свой-

ством, что на каждом такте работы автомата переход осуществляется в состояние с тем же или большим номером. Такие модели, предложенные впервые Бакисом ([27, 12]), называются лево-правыми (left-right), или моделями Бакиса. Автор предложил при изложении на автоматном языке называть соответствующие этим моделям вероятностные автоматы монотонными.

Распознавание речи с помощью автономных монотонных вероятностных автоматов сводится к следующему. Каждой речевой команде сопоставляется эталон в виде вероятностного автомата. Вычисляется вероятность наблюдения распознаваемого кодового слова на выходе каждого из таких автоматов, и выбирается автомат, максимизирующий эту вероятность. Соответствующая найденному автомату команда считается результатом распознавания.

Метод СММ на уровне распознавания отдельных команд (т.е. без учета контекстной зависимости слов внутри предложения на естественном языке) дает надежность распознавания, близкую к надежности распознавания речевых команд человеком. Это позволяет сделать вывод о том, что общая задача распознавания отдельно произносимых команд решается методом СММ почти с максимально возможной точностью. Тем не менее, в практических приложениях возникает необходимость построения систем распознавания речи с ограниченным словарем, которые должны обладать сравнительно более высокой надежностью работы.

Традиционно эта задача решается методом подбора словаря команд. Подбор словаря может осуществляться путем проведения экспериментов с распознаванием команд из различных словарей и сравнения результатов этих экспериментов, что является «дорогим» решением. Поэтому представляет интерес задача автоматического подбора словаря команд без проведения экспериментов для обеспечения требуемой надежности распознавания.

Для решения этой задачи автор предложил использовать функцию качества словаря команд, связанную со средней вероятностью ошибки при распознавании и вычисляемую непосредственно по эталонам команд без проведения экспериментов с помощью введенной метрики на множестве монотонных вероятностных автоматов.

В статистике традиционно задача нахождения расстояния между случайными величинами решается на основе статистических мер близости типа расстояний Махаланобиса, Кульбака-Лейбнера, Бхатачария и др. ([30]). Однако, данные метрики не имеют эффективного обобщения на случай вероятностных автоматов. В настоящей работе предлагается новая адекватная метрика, которая вычисляется эффективно.

Целью настоящей работы является построение формальной математической модели для распознавания речи с целью исследования вопроса оценки качества работы различных автоматных алгоритмов распознавания речи для различных словарей команд.

В работе имеются следующие достижения:

1. Предложена формальная математическая модель для задачи распознавания речи, в которой выделены множества (алфавиты букв, звуков, сигналов и кодов) и действующие на них функции (транскрибирования, произнесения, описания, кодирования и распознавания). Данная модель охватывает все случаи распознавания речевых сигналов в произвольных внешних условиях. В эту модель укладывается, в частности, задача распознавания речи при наличии дополнительных неакустических источников речевой информации в условиях шумов.

2. Сформулирована постановка задачи распознавания речи в данной модели и предложена функция для оценки качества словаря команд для фиксированного алгоритма распознавания.

3. Предложен универсальный детерминированно-автоматный алго-

ритм распознавания в рамках введенной модели, для которого найдена нижняя оценка функции качества словаря команд.

4. Известный метод распознавания речи на основе скрытых марковских моделей переведен на язык вероятностных автоматов.

5. Введено новое понятие автономного монотонного вероятностного автомата. На множестве таких автоматов введены операции декартового и скалярного произведений автоматов и на их основе построена метрика, эффективно вычисляемая по автоматам.

6. Введено понятие качества словаря команд для вероятностного случая, получена точная формула для вычисления качества произвольного словаря команд через метрику на множестве монотонных автономных вероятностных автоматов.

**В первой главе** вводится формальное определение речевой модели как математического объекта. Основой ее является понятие **звукового сигнала**, под которым понимается функция  $s : \mathbb{R} \rightarrow \mathbb{R}$ , непрерывная на некотором отрезке  $[t_0, t_1] \subset \mathbb{R}$  и принимающая вне этого отрезка нулевые значения. Переходя к дискретному времени, **дискретным звуковым сигналом** называется любое отображение  $s' : \mathbb{Z} \rightarrow \mathbb{K}$ , такое что  $\exists z_0, z_1 \in \mathbb{Z} : s'(z) = 0 \ \forall z \notin [z_0, z_1]$ . Дискретные звуковые сигналы получаются из непрерывных с помощью операции **дискретизации**:  $s'(z) = s(z \cdot 5)$ , где 5 называется **периодом дискретизации**. Для простоты далее будем опускать штрих при обозначении звукового сигнала.

**Носителем** сигнала  $s : \mathbb{Z} \rightarrow \mathbb{R}, s \not\equiv 0$  назовем отрезок  $[\inf\{z | s(z) \neq 0\}, \sup\{z | s(z) \neq 0\}]$ , вне которого он принимает нулевые значения. Границы носителя назовем **началом** и **концом** сигнала  $s$  соответственно. Определим операцию **конкатенации**  $S1S2$  сигналов  $s_1, s_2$  с непересекающимися носителями как  $(s_1s_2)(z) = \max(s_1(z), s_2(z))$ .

Далее дается понятие речевой модели, отдельно для случая непрерыв-

ного и для случая дискретного времени.

Понятие речевой модели в дискретном случае формулируется в определении **1.15**. Упрощенно это определение выглядит так:

**Речевой моделью** называется восьмерка  $(B, V, S, C, d, \Gamma, \Pi, \gamma)$ , где

- $B$  — конечный алфавит звуков;
- $V \subseteq B^*$  — конечный словарь команд;
- $S$  — множество дискретных звуковых сигналов — произнесений команд;
- $C$  — конечный алфавит — кодовая книга описаний сигналов;
- $\Pi : B \rightarrow \mathcal{P}(S)$  — функция произнесения (здесь  $\mathcal{P}(S)$  — множество всех подмножеств множества  $S$ );  $\Pi(V)$  — множество всех произнесений команды;
- $d : S \rightarrow (R^M)^*$  — функция описания речевых сигналов, сопоставляющая речевому сигналу  $s$  матрицу  $M \times n$ , где  $M$  — некоторая константа (размерность пространства описаний), а  $n$  линейно зависит от длины речевого сигнала. Каждый столбец матрицы описания задает точку в пространстве описаний  $R^M$ ;

**и**

- $\Gamma : C \rightarrow B$  — кодирующее отображение, сопоставляющее каждому столбцу матрицы описания элемент кодовой книги  $C$ .  $\Gamma$  естественно обобщается на слова как  $\Gamma : (C^M)^* \rightarrow B^*$ ;
- $\gamma : B^* \rightarrow \mathcal{P}(S)$  — функция распознавания, определяемая как

$$\gamma(b) = \{s \in S \mid b \in \Gamma(\Pi(s))\}.$$

Обозначим через  $\Pi(\emptyset)$  множество сигналов, не являющихся произнесениями никаких звуков:

$$\Pi(\emptyset) = \bigcup_{b \in B} \Pi(b)$$

Будем считать, что функция произнесения  $\Pi$  удовлетворяет следующим свойствам:

1.  $Y/3 \in B \quad \Pi(/3)^0$ ;
2.  $V/3 \in B \quad A \wedge /3 \Rightarrow \text{ИОД} \Pi 11(/3) = 0$ ;
3.  $V/3 \in \mathcal{L} \quad \forall s \in \Pi(/3) \quad s \neq 0$ ;

А.  $UP \in B : /3 = \text{И}6_2 \dots 6_n$ , и  $n > 1 \quad \forall s \in \Pi(/3)$

$3s_0l, Si, S12, S2, \dots, S_n - \text{И} S_n S_n 0 \in 5: S = S_0 \text{И} Si Si S_2 \dots S_n - \text{И} S_n S_n 0$ , где

$Si \in 11(3/4)$  для  $i = 1, 2, \dots, n$ ,

$s_{ij} \in 11(6, -) \cup 1\%$   $\cup \Pi(\text{И}0)$  для  $i, j = 1, 2, \dots, n$ ,

$S_{\text{И}0}, s_{\text{И}0j} \in \Pi(\_) \cup \Pi(6\gg) \cup \Pi(6_0)$  для  $\gamma = 1, 2, \dots, n$ ,

$\_$  — специальный символ алфавита (символ паузы)

и все сигналы  $S_j, \mathbf{S}_{ij}$  имеют непересекающиеся носители.

Как правило, функция описания  $d$  строится в виде композиции оконного отображения, вырезающего с постоянным шагом  $L$  из сигнала временное окно в общем случае переменной длины, и функции локального описания сигнала, сопоставляющей каждому окну сигнала точку в пространстве описаний  $R^M$ .

**Образом**  $V(b)$  с  $\text{ИИ}^M$  фонемы  $b \in B$  как в непрерывном, так и в дискретном случае называется совокупность всех точек описания всех ее произнесений. Для дискретного случая

$$V(b) = \{p \in \text{ИИ}^M \setminus \exists y \in d(H(b)) : y = \|y_1\| y_2 \dots \|p\| \dots \|y_n\| \|y_n\|\},$$

где  $\|y_1\| y_2 \dots \|y_n\|$  обозначает матрицу, состоящую из столбцов  $y_1, y_2, \dots, y_n$ .

В **утверждении** 1.6. показано, что для любых допустимых значений периода дискретизации  $S$  и шага описания  $L$  дискретный и непрерывный образы одной и той же фонемы совпадают. Это, по сути, означает эквивалентность понятий образа фонемы в непрерывном и дискретном случае, что позволяет, ничего не теряя в точности решения задачи, использовать для распознавания речевую модель с дискретным временем.

При этом можно считать, что мы рассматриваем описания сигналов с шагом 1.

Задача распознавания речи во введенной речевой модели состоит в восстановлении функции распознавания  $\gamma$ . Содержательно, распознавание речи сводится к тому, чтобы по заданному кодовому слову  $y \in C^*$ , про которое известно, что оно является кодом описания произнесения одной из команд словаря  $B$ , найти в словаре эту команду.

Задачу распознавания можно решать методом сравнения с эталонами. Пусть каждой команде  $\beta \in B$  сопоставлен эталон  $e^*$  - элемент некоторого множества  $E = \{e^1, e^2, \dots, e^N\}$ . Содержательно эталон команды является компактным способом представления множества кодовых слов всех ее произнесений.

Пусть также на декартовом произведении  $C^* \times E$  введена функция **расстояния**  $p : C^* \times E \rightarrow R$  между кодовыми словами и эталонами. Алгоритм  $(E, p)$  распознавания речи методом сравнения с эталонами восстанавливает функцию распознавания  $\gamma$  с помощью формулы  $\gamma(y, p) = \text{Argmin}_{e \in E} \sum_{z=1}^N p(y_z, e_z)$ . Если  $\gamma(y, p) \neq \beta$ , будем говорить, что имеют место ошибки распознавания.

Качество  $q(E, p)$  словаря команд  $B$  для данного алгоритма распознавания  $(E, p)$  естественно определить, как  $1 - \text{рош.}$ , где  $\text{рош.}$  - среднее число ошибок распознавания команд из  $B$  алгоритмом  $(E, p)$ .

Важно заметить, что каждой команде  $\beta$  соответствует множество кодовых слов описаний всех ее произнесений  $\Gamma(\beta)$ , которое является конечным подязыком языка  $C^*$  всех кодовых слов. Автоматный подход к распознаванию речи состоит в том, чтобы представить язык кодовых слов каждой команды некоторым автоматом. В случае конечно-детерминированных автоматов мы погружаем язык кодовых слов команды в некоторый известный нам бесконечный регулярный язык, а в случае вероятностных автоматов каждому слову языка кодовых слов при-





писываем некоторое действительное число — его вероятность, что дает нам возможность рассматривать случай пересекающихся языков кодовых слов.

**Во второй главе** эталоны команд задаются конечно-детерминированными автоматами. Построен алгоритм распознавания речи на основе конечных детерминированных автоматов, где множество произнесений команды представляется некоторым бесконечным регулярным языком, приведена конструкция автоматического порождения указанных регулярных языков, введена функция качества словаря команд и получена нижняя оценка значения функции качества в детерминированно-автоматном случае.

Предложен следующий способ построения автоматов.

Пусть пространство описаний сигналов  $R^M$  разбито на  $m$  непересекающихся подмножеств

$$K^M = \bigcup_{i=1, m} [J_i, Q_i]$$

$C = \{c_1, c_2, \dots, c_m\}$ . Сопоставим каждому элементу  $c_i$  разбиения  $C$  символ  $C_i$  и выберем алфавит  $C = \{c_1, c_2, \dots, c_m\}$  в качестве кодовой книги нашей речевой модели.

Кодирующее отображение будем строить на основе разбиения  $C$ . А именно, для каждой точки  $p \in R^M$  определим  $\Gamma(p)$  как  $\Gamma(p) = \{c \in C : p \in c\}$ . Предположим, что при распознавании мы не умеем вычислять  $\Gamma(p)$  для всех точек  $p \in R^M$ . Пусть, тем не менее, для каждой пары элементов разбиения  $C$  мы умеем "отличать их друг от друга". Более строго,

**Отделяющей функцией** для двух множеств  $c_i, c_j \in C$  назовем функцию  $\phi_{c_i c_j} : \mathbb{R}^M \rightarrow C \cup \{?\}$ , такую что  $\forall p \in \mathbb{R}^M$

$$0^\circ \phi_{c_i c_j}(p) = c_i;$$

1°  $\phi_{CiCi}(p) = ipc_{jCi}(p)$  (симметричность);

3°  $p \in CI \Rightarrow \phi_{Ci} c M \in fa, ?$ ;

С помощью функции  $\phi$  можно построить функцию локального распознавания кода описания сигнала:

**Функцией локального распознавания  $\Phi$**  будем называть функцию  $\Phi : E^M \rightarrow 2^{\mathcal{C}}$ , определенную для каждой точки  $p \in R^M$  следующим образом:  $\Phi(p) = \{c \in C : \forall c' \in C, c' \neq c \Rightarrow \phi_{cc'} > \epsilon\}$ .

Множество значений функции  $\Phi$  порождает регулярный язык "близких" гипотез распознавания кода описания сигнала, описываемый регулярным выражением  $\Gamma(y) = \Gamma(\Phi(I/1)) \cap (\Phi(I/2)) \dots R(\wedge(y_n))$ , где матрица  $y = \|y_1 y_2 \dots y_n\|$  регулярное выражение  $R$  определено как  $R(\{a_1, a_2, \dots, a_m\}) = (a_1 a_2 \dots a_m)$ .

Обозначим через  $\{ \%$  регулярный язык, описываемый регулярным выражением  $I$ .

Справедливо **Утверждение 2.2.** о том, что множество  $\Phi(p)$  "близких" гипотез локального распознавания содержит всегда "правильную" гипотезу  $\Gamma(p)$ :

$$\forall p \in R^M \quad \Gamma(p) \in \Phi(p).$$

Как следствие получаем, что для описания  $y = d(s)$  любого сигнала  $s$  справедливо  $\Gamma(y) \in (\Gamma(y))$ -

Говорим, что  $c \in C$  и  $c' \in C$  **отделимы** с помощью отделяющей функции  $\phi$ , если  $\forall p \in C \cup C' \quad \phi(p) \in \{c, c'\}$ .

**Окрестностью  $O(C)$**  некоторого подалфавита  $C \subset C$  кодовой книги  $C$  будем называть множество всех букв  $c$  из  $C$ , не являющихся отделимыми хотя бы от одной буквы алфавита  $C$ .

Обозначим через  $c(B)$  множество  $c(B) = \{c' \in C : c' \Gamma(B) \neq 0\}$ .

Свойство 4 функции произнесения  $\Pi$  позволяет нам ввести эталоны

команд в виде регулярных языков, порождаемых регулярными выражениями специального вида, которые строятся непосредственно по текстам команд.

А именно, **эталонным кодом** команды  $(\beta = b_1 \cdot \dots \cdot K$  назовем регулярное выражение

$$f(\beta) = D(0(c(\_) \cup c(b_0) \cup \phi_i)))D(0(c(b_0)))D(0(c(b_1) \cup c(b_2) \cup \dots \cup c(b_n)))L(0(c(b_a))) \cdot \dots \cdot \\ \dots D(0(c(b_{n-1})))D(0(c(b_{n-2}) \cup c(b_{n-1}) \cup c(b_n)))D(0(\phi_{n-1}))D(0(c(b_n) \cup c(b_0) \cup c(b_h))).$$

Код  $r(d(s))$  любого произнесения  $s \in \Pi(\beta)$  команды  $\beta$  лежит в эталоне этой команды - регулярном языке  $(\Gamma(\beta))$ .

Справедлива более сильная **теорема 2.1**:

В речевой модели, в которой функция произнесения обладает свойствами 1-4, для любой звуковой команды  $\beta = bfa..b_n$  и описания  $y = d(s)$  ее произвольного произнесения  $s \in \Pi(\beta)$  имеет место вложение  $(\Gamma(y)) \subset (\Gamma(\beta))$ .

Далее во второй главе показано, что если звуки объединить в классы эквивалентности с помощью отношения отделимости, и в качестве разбиения  $S$  пространства описаний рассмотреть образы этих классов при отображении  $d$  в  $\Gamma$ , то каждый язык гипотез распознавания  $(\Gamma(y))$  будет состоять ровно из одного элемента - кодового слова  $\Gamma(y)$ .

Теперь автоматный алгоритм  $(E, p)$  распознавания речи в модели  $(B, V, S, C, \langle i, \Gamma, \Pi, \gamma \rangle)$  можно задать множеством эталонов  $E = \{e_1, e_2, \dots, e_m\} \subset I = (\Gamma(\beta))$ , и функцией расстояния  $p$ , определяемой для кодового слова  $\gamma = \Gamma(\phi)$ ,  $\gamma \in \Pi(A)$  как:

$$I \begin{cases} 0, & \text{если } \gamma \in E \\ \in I - 1, & \text{иначе} \end{cases}$$

Результат распознавания кодового слова  $y$  алгоритмом  $(E, p)$  есть множество команд  $\Gamma(y_{\gamma/\gamma})(\gamma) = \{P \in I \mid P \in (f(\beta))\}$ .

Справедливо **утверждение 2.4** о том, что в условиях модели для любого кодового слова  $\gamma$  имеет место вложение  $\Gamma(\gamma) \subset \Gamma(E, p)(I)$ -

Функцию качества словаря для введенного детерминированно-автоматного алгоритма распознавания  $(E, p)$  определим следующим образом. Если словарь  $B$  состоит из одного слова  $/3$ , то положим, что его качество равно 1 (естественно предположить, что в словаре из одной команды распознаватель не делает ошибок). Для каждого словаря  $B \in \{ /3, (3') \}$ , состоящего из двух команд, определим функцию качества как

$$(a \rightarrow a') \rightarrow 0 \quad \text{и} \quad \frac{7 \in \Gamma(\langle i(\Pi(/3)) \Pi((3')) \rangle)}{|\Gamma(\text{даоз}) \Pi(\text{поз}')|}$$

Пусть  $S$  — некоторое конечное множество,  $/ : S \times S \rightarrow R$ . Обозначим через  $@s(f)$  среднее значение функции  $/$  на  $S \times S$ :

$$w^s \wedge I \sim \frac{|S| - (15) - 1}{2}$$

Теперь для произвольного словаря команд  $B$  определим его качество как среднее качество по всем парам слов из  $B$ :

$$Я(E, p)(B) = e_B \{ q(E \rightarrow p) \} \quad \text{Введем расстояние } p(e, e')$$

между регулярными языками  $e$  и  $e'$  как

$\left\{ \begin{array}{l} 0, \quad \text{если } e \Pi e' \neq 0, 1, \\ \text{иначе} \end{array} \right.$  Справедлива **Теорема 2.2** о нижней оценке качества словаря команд при распознавании детерминированно-автоматным алгоритмом  $(E, p)$ :

$$\langle ?(Я, p)(\xi) \rangle \in C_B(/3).$$

**В третьей главе** эталоны задаются в виде автономных монотонных вероятностных автоматов.

Этот подход не является новым. В современных системах распознавания речи распознавание производится с помощью метода скрытых марковских моделей [18]. В основе этого метода лежит алгоритм, который можно перевести на язык вероятностных автоматов.

Формально, **автономный вероятностный автомат** — это тройка  $(C, Q, \nu)$ , где  $C$  - алфавит выходных символов;  $Q$  - конечный алфавит состояний автомата,  $\nu$  - функция, определенная на множестве состояний автомата  $Q$  и принимающая в качестве своих значений вероятностные меры на множестве  $C \times Q$ , такая, что она разлагается в произведение  $\nu = \gamma \cdot P$ , где  $\gamma$  действует на множестве  $Q$  и имеет значениями вероятностные меры на множестве  $Q$ , а  $P$  действует на том же множестве  $Q$  и имеет значениями вероятностные меры на множестве  $C$  [26]. Содержательно  $\nu\{c, q'/q\}$  интерпретируется как условная вероятность перейти в состояние  $q'$  и выдать символ  $c$  при условии, что в предыдущий момент времени автомат находился в состоянии  $q$ . Функцию  $\nu$  можно задать двумя матрицами - квадратной  $m \times m$ -матрицей  $\gamma = \|\gamma_{jj}\|$ :  $\gamma^i = \gamma(qj/qi), i, j = 1, m$  и прямоугольной  $m \times \kappa$ -матрицей  $P = \|\gamma^i\|$ :  $P_{ik} = P(ci/qi), i = 1, m, I = 1, \kappa, \kappa = |C|$ . Матрицы  $\gamma$  и  $P$  являются стохастическими (по строкам). Далее для вероятностных автоматов вида  $(C, Q, \gamma \cdot P)$  будем использовать обозначение  $(C, Q, \gamma, P)$ .

В каждый момент времени состояние автомата характеризуется вектором распределения вероятностей вида  $(p_1, p_2, \dots, p_m)$ , где  $m = |Q|$ ,  $p_i$  — вероятность находиться в состоянии  $i$ ,  $\sum_{i=1}^m p_i = 1$ . Будем называть

$$i=1, m$$

вероятностный автомат **инициальным**, если заданы распределения  $u^0$  и  $v^F$  для его начального и финального состояний соответственно (обозначение  $(C, Q, \gamma, P, u^0, v^F)$ )-

Можно рассмотреть обобщение понятия вероятностного автомата, положив, что матрицы  $\gamma$  и  $P$  не являются стохастическими, но состоят из неотрицательных элементов и сумма элементов в каждой их строке не превосходит 1 (такие матрицы будем называть **слабостохастическими**). Далее автоматы со слабостохастическими матрицами  $\gamma$  и  $P$  также будем называть вероятностными.

Функционирование инициального автомата определяется следующим



образом. Начальное состояние определяется распределением  $z/\circ$ . На каждом шаге, находясь в некотором состоянии  $q$ , автомат сначала выдает некоторую букву  $c$  алфавита  $C$ , исходя из распределения вероятностей  $P(q)$ , а затем переходит в следующее состояние, исходя из распределения вероятностей  $\Gamma(q)$ . Поскольку матрицы  $\Gamma$  и  $P$  не являются в общем случае стохастическими, будем считать, что на каждом шаге существует вероятность того, что автомат не выдал никакой буквы или не осуществил перехода. Если на некотором шаге автомат перешел в состояние  $q_i$ , для которого  $\sum_{c \in C} \Gamma(q_i, c) = 1$  считаем, что автомат завершил работу. В

этом случае буквы, которые автомат начиная с начального состояния последовательно выдавал на выход в процессе своей работы, образуют некоторое слово  $\gamma$  над алфавитом  $C$ . Говорим в таком случае, что автомат "выдал" слово  $\gamma$ .

Если некоторое слово  $\gamma = c_1 c_2 \dots c_n$  выдано автоматом  $(C, Q, \Gamma, P, z/\circ, u^F)$ , **вероятностью** этого слова назовем величину

$$p_A(\gamma) = u^G P(c_1) P(c_2) \dots P(c_n) (u^F)^T,$$

где  $^m$  означает транспонирование матрицы, а для каждого  $q \in Q$   $P(q)$  - квадратная  $m \times m$  матрица и  $P\{c\}_q = \Gamma(q, c) \cdot P_q$  — вероятность того, что находясь в состоянии  $q$ , автомат выдал букву  $c$  и перешел в состояние  $q'$

Инициальный автономный вероятностный автомат  $\mathcal{A} = (C, Q, \Gamma, P, z/\circ, u^F)$  будем называть **монотонным**, если выполнены условия:

1.  $\Gamma_{ij} = 0$  при  $i > j$ ;
2.  $z/\circ = (1, 0, 0, \dots, 0)$ ; 3.  $u^F = (0, 0, \dots, 0, 1)$ ;

**Теорема 3.1.** утверждает, что для любого монотонного автономного

вероятностного автомата  $\mathcal{L} = (C, Q, \Gamma, P, u^0, u^F)$ , для которого выполнено условие  $\sum_{i=1}^m p_{ti} < 1$  для всех  $i = 1, m$ , события, связанные с выдачей автоматом различных слов из  $C^*$ , являются несовместными.

Это позволяет сопоставить каждому автономному монотонному вероятностному автомату  $\mathcal{L}$  **стохастическую (слабостохастическую) словарную функцию**  $\text{rd} : C^* \rightarrow [0,1]$ , вычисляющую вероятность выдачи слов из  $C^*$  автоматом  $\mathcal{L}$  и обладающую свойством

где сумма берется по всем словам  $\gamma \in C^*$ , занумерованным в лексикографическом порядке. Каждую такую словарную функцию можно задать бесконечным вектором

$$(p_1, p_2, \dots, p_n, \dots) \in h(C^*),$$

где  $p_i = \text{rd}(\gamma_i)$ ,  $h = \{(p_1, p_2, \dots, p_n, \dots) : \sum_{i=1}^{\infty} p_i < \infty\}$ ,

$$h/2 = \{(p_1, p_2, \dots, p_n, \dots) : \sum_{i=1}^{\infty} p_i^2 < \infty\}$$

Методу скрытых марковских моделей соответствует в нашей формализации алгоритм  $(E, p)$  распознавания речи, в котором эталоны  $E = \{e_1, e_2, \dots, e_d\}$  команд  $\{z_1, z_2, \dots, z_n\}$  задаются в виде некоторых автономных монотонных вероятностных автоматов, синтезированных по примерам кодовых слов одним из известных методов [19-22], а расстояние между эталонами и кодовыми словами вводится как  $p(e^i, \gamma) = 1 - P(e^i | \gamma)$ .

Содержательно метод сравнения с эталонами в вероятностно-автоматном случае сводится к поиску команды, на которой достигается максимум вероятности выдачи распознаваемого кодового слова соответствующим этой команде эталонным вероятностным автоматом.

**Качество** словаря команд для алгоритма  $(E, p)$  будем также вводить следующим образом. Для  $B$  вида  $\{z\}$  положим  $Y(E, p)(\{z\}) = 1 - D^{\text{ЛЯ}}$  словаря из двух команд определим его качество как единица минус удвоен-





ная вероятность ошибки распознавания в этом словаре при условии, что распознаваемые кодовые слова выдаются автоматами, соответствующими каждой из команд, и автоматы включаются по очереди с одинаковой вероятностью  $1/2$ . Для словаря, состоящего из более чем двух команд, определим его качество как среднее качество всех его подсловарей из двух команд.

Далее в третьей главе решается задача оценки качества словарей команд для алгоритма  $(E, p)$ . С этой целью на множестве эталонов команд (автономных монотонных вероятностных автоматов) вводится метрика.

Если на множестве стохастических словарных функций метрика каким-либо образом введена, расстоянием между автономными монотонными вероятностными автоматами определяем как расстояние между соответствующими им словарными функциями.

Определения 3.8, 3.9, 3.11 задают три способа введения метрики на множестве стохастических словарных функций:

$$\bullet P_1(P_1, P_2) = \sum_{\gamma \in C^*} P_1(\gamma) - P_2(\gamma);$$

$$\bullet P_2(P_1, P_2) = 1 - \sum_{\gamma \in C^*} \min(P_1(\gamma), P_2(\gamma));$$

$$\bullet P_3(P_1, P_2) = \sqrt{\sum_{\gamma \in C^*} |P_1(\gamma) - P_2(\gamma)|^2}.$$

где  $\|\bullet\|$  — норма в  $\mathbb{R}^n$ .

Все эти три функции расстояния обладают свойствами метрики. Показано также, что на множестве стохастических словарных функций первые две метрики совпадают с точностью до постоянного множителя.

Связь метрики  $p_2$  и вероятности ошибки при распознавании иллюстрирует Лемма 3.4., утверждающая, что  $p_2(p_{el}, P_{e2}) = 1 - 2p_{0ш.}$ , где  $p_{0ш.}$  — вероятность ошибки распознавания в словаре  $\{1, 2, 3, 4\}$ .

Это позволяет выразить функцию качества словаря команд через расстояние между соответствующими командам эталонными вероятностными

ми автоматами. **Теорема 3.2** утверждает, что

$$Q(E,P)(B) = Q_B(P2)-$$

Далее доказывается, что метрика  $\frac{3}{4}$  эффективно вычисляется по автоматам.

Для этого в определении 3.12 вводится понятие **декартова произведения автономных автоматов**:

Декартовым произведением  $A = \prod_i A_i$  автономных вероятностных

автоматов  $A = (C, Q, \gamma, P)$

$A = (C, \gamma, Q, P)$  называем вероятностный автомат

$\gamma = (c^f, g_2 \times g_2, \gamma, i/\gamma, i/\gamma)$ ,

где  $\gamma$ -матрица,  $m = \gamma_{ii}$ ,  $m_i = \gamma_{ii}$ ,  $i = 1, 2$ ,

$$P(hj)k = (Pl)ik * (P2)jk'i$$

**Ай = to<sup>0</sup>**. - ( $\frac{3}{4}\frac{3}{4}$ ;  $Lu$ ) = ("А ■ ("А-

Приведенной матрицей переходов  $\gamma$  вероятностного автомата  $(C, Q, \gamma, P)$  назовем  $m \times m$ -матрицу,  $(\gamma, \gamma)$ -й элемент которой определен как

$$I=I$$

**Теорема 3-3** показывает связь между понятием декартова произведения  $A = \prod_i A_i$  автоматов  $A_i$  и скалярным произведением соответствующих этим автоматам словарных функций в  $\mathcal{L}$

(здесь  $E$  - единичная матрица).

Используя известную связь между нормой и скалярным произведением ( $\|v\| = \sqrt{v^T v}$ ) получаем, как следствие, формулу для вычисления расстояния  $r_z$  между автономными монотонными вероятностными авто-



матами  $L_i = \{C, Q_h, T_{g_2}, P_h, z/D, u\}$  и  $L_2 = (C, Q_2, T_{g_2}, P_2, \wedge^2, \wedge^2)$ :

$$P_3(L, L) \quad \wedge((E - P \cdot z \Gamma^T), \wedge \wedge \wedge)^{-1},$$

$m$   
2

(здесь  $m_i = |Q_i|$ ,  $m_2 = |Q_2|$ ,  $T_{g_2}$  — приведенная матрица переходов автомата  $L_i \times L_2$ ,  $T_{g_i}$  — приведенная матрица переходов автомата  $L_i \times L_i$ ,  $\wedge^2$  — приведенная матрица переходов автомата  $L_2 \times L_2$ )

Окончательно, качество словаря команд  $B$  с помощью метрики  $r_z$  определяется как

## Приложения.

На основе изложенных в работе подходов написаны прикладные программы распознавания речи. Система распознавания речи с ограниченным словарем ([a2]) работает в операционной системе Windows и позволяет в реальном времени распознавать речевые команды из словаря объемом до 100 команд. Система позволяет вводить команды на любом естественном языке, поскольку для обучения эталонов используется поэлементный метод [33]. В качестве моделей эталонов система позволяет одинаково эффективно работать как с методом ДДВ, так и с методом СММ. Система использует алгоритм распознавания речи, основанный на конечно-автоматном подходе, изложенном во второй главе. В качестве алфавита классов звуков выбран алфавит  $\{L, O, /, \_, X\}$  ((см. приложение 2, [a1]), где класс  $X$  включает шипящие согласные,  $\_$  — символы паузы и глухие смычки, а классы  $A, O, I$  соответствуют низкочастотным, среднечастотным и высокочастотным гласным звукам соответственно).

Конструктор систем распознавания речи в условиях сильных акустических шумов ([a4]) также работает с отдельно произносимыми командами и фразами из ограниченного словаря. В качестве устройства ввода

для конструктора могут подключаться как обычные микрофоны, так и системы из различных датчиков: акустических, фотодатчиков слежения за губами говорящего, датчиков дыхания и т.п. В отличие от описанной выше системы распознавания, конструктор позволяет динамически строить функции описания сигналов на основе сигналов с датчиков с помощью богатого семейства унарных и бинарных операций. Конструктор позволяет строить также устойчиво работающие модули выделения границ речевых сообщений на фоне шума, что повышает в сотни раз надежность распознавания речи в шуме по сравнению с традиционно используемыми методами распознавания. Работа конструктора основана на описанной в первой главе формальной речевой модели.

Система автоматического подравнивания произношения является дополнением к мультимедийным системам обучения иностранным языкам. После настройки на голос ученика электронный учитель позволяет определить количество, место и тип ошибок при произнесении учеником фраз на иностранном языке. При построении системы использовалась оригинальная методика построения модели «правильного произношения» для ученика в виде монотонного вероятностного автомата.

Метрика на множестве автономных вероятностных автоматов, введенная в третьей главе настоящей работы, была эффективно использована при решении задачи оптимального выбора фонетического алфавита при разработке системы распознавания русской речи в рамках гранта с фирмой Intel Corp., США ([35]). С помощью метрики  $p\delta$  была построена матрица попарных расстояний между фонемами русского языка (см. приложение 3), представленных в виде автономных вероятностных автоматов с 4 состояниями, которые были синтезированы на основе русской речевой базы данных. На основе проведенного эксперимента удалось показать, что алфавит из 150 фонемных символов для русского языка можно сократить без потенциальной потери точности при распознавании до

120 символов, что может значительно увеличить эффективность системы распознавания речи на русском языке.

Введенная в настоящей работе метрика на множестве автономных монотонных вероятностных автоматов имеет широкие возможности практического применения. Она может быть использована в тех областях, где в настоящее время требуется уметь эффективно вычислять расстояние между марковскими автоматами: при решении задач подбора фонетического алфавита ([36]) и словаря команд ([38]), распознавания и идентификации диктора, конверсии и синтеза речи ([37] - здесь метрика используется для решения задачи интерполяции голоса диктора), при программировании стратегических компьютерных игр ([40]) и моделировании стратегий поведения человека в реальных условиях ([39]). Во всех этих практических примерах в настоящее время используются статистические меры близости типа расстояний Махаланобиса, Кульбака-Лейбнера, метод Монте-Карло и им подобные. Поскольку эти меры близости зачастую не обладают всеми свойствами метрики (не выполняется неравенство треугольника), введение метрики на множестве скрытых марковских моделей может дать более эффективные результаты при решении данных практических задач.

Автор выражает благодарность своему научному руководителю доктору физико-математических наук, профессору Бабину Дмитрию Николаевичу за постановку задачи и помощь в работе над диссертацией.

# ГЛАВА 1. ФОРМАЛЬНАЯ МОДЕЛЬ РЕЧИ

## §1.1. Речь как математический объект. Основные определения

В данном разделе будет введена формальная модель речи для естественного языка, на котором произносятся слова. Без ограничения общности можно считать, что рассматривается русский язык, хотя построенная модель инвариантна относительно выбора языка.

Пусть  $A$  - некоторый конечный алфавит букв языка. Пусть  $A$  - фиксированное подмножество множества  $A^*$  слов в алфавите  $A$ . Назовем множество  $A$  *словарем команд*, а элементы множества  $A$  - *командами*.

Пусть  $B$  - некоторый алфавит, который будем в дальнейшем называть *алфавитом звуков языка*, или *фонетическим алфавитом* (фонетический алфавит для русского языка приведен в приложении 1). Будем считать, что алфавит звуков содержит некоторый специальный символ - символ паузы  $\_$ . Пусть  $\mathcal{B}$  - фиксированное подмножество  $B^*$  слов в алфавите  $B$ , такое что  $B \subset \mathcal{B}$ . Назовем множество  $\mathcal{B}$  *фонетическим словарем*, или *словарем звуковых команд*, а элементы этого множества - *фонетическими словами* {звуковыми командами}.

Через  $|Z|$  будем обозначать длину слова  $Z$  (количество букв в нем), а через  $R$ ,  $N$  и  $Z$  - множества действительных, натуральных и целых чисел соответственно.

Определение 1.1. *Звуковым сигналом* назовем функцию  $s: R \rightarrow R$ , непрерывную на некотором отрезке  $[t_0, t_1] \subset R$  и принимающую вне этого отрезка нулевые значения (см. рис 1.1).



Область определения звукового сигнала будем называть *временем*, а точки на ней — *отсчетами* {моментами времени}.

Началом  $t_0(s)$  звукового сигнала  $s$  назовем момент времени  $t_0(s) = \inf\{t: s(t) \neq 0\}$ , а концом  $t_i(s)$  сигнала  $s$  - момент времени  $t_i(s) = \sup\{t: s(t) \neq 0\}$ . Длиной  $|s|$  сигнала  $s$  будем называть величину  $|s| = t_i(s) - t_0(s)$ .

Через  $S$  обозначим множество всех звуковых сигналов.

С физической точки зрения каждый звуковой сигнал  $s$  задает изменение величины звукового давления на отрезке времени  $[t_0(s), t_i(s)]$ .

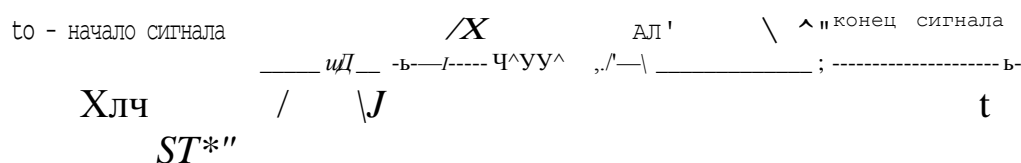


Рисунок 1.1. Звуковой сигнал

На множестве  $S$  введем отношение эквивалентности.

**Определение 1.2.** Сигналы  $s_1$  и  $s_2$  будем называть *эквивалентными*, если при всех  $t \in \mathbb{R}$  выполнено  $s_2(t) = s_1(t+k)$  для некоторого фиксированного  $k \in \mathbb{R}$ .

Эквивалентность сигналов  $s_1$  и  $s_2$  будем обозначать как  $s_1 \sim s_2$ .

Через  $S[u,v]$  будем обозначать *ограничение сигнала  $s$  на отрезок  $[u,v]$* , т.е.

звуковой сигнал, определенный следующим образом:  $S[u,v](t) = s(t)$

$\forall t \in [u,v] \quad S[u,v](t) = 0 \quad \forall t \notin [u,v]$

Напомним, что *преобразованием Фурье*  $\Phi_8$  функции  $s: \mathbb{R} \rightarrow \mathbb{R}$  называется функция  $O_s: \mathbb{R} \rightarrow \mathbb{C}$ ,  $\langle D_s(f) = \int_{-\infty}^{+\infty} s(t) e^{-2\pi i f t} dt \rangle$  [29].

Определение 1.3. *Речевой моделью* назовем десятку

$\langle A, D, B, \delta, \tau, \Pi, T, \alpha, \Gamma \rangle$ , где

- $A, B$  - конечные алфавиты букв и звуков соответственно;
- $T \subseteq A^*$  - конечный словарь команд;
- $S \subseteq B^*$  - конечный словарь звуковых команд;
- $S$  - множество звуковых сигналов;
- $\tau: A^* \rightarrow B^*$  - конечно-автоматное отображение *транскрибирования*,

и

- $\Pi: B \rightarrow 2$  - отображение *произнесения*,

обладающее следующими свойствами:

$$1^\circ \forall p \in \Pi((3)^0,$$

$$2^\circ p \in M, MP_2 \Rightarrow \Pi(p) \cap \Pi(P_2) = \emptyset;$$

3°  $\forall p \in B: |p| > 1 \quad \forall \text{sen}(p) \quad \Phi_8(i) > 0$  при  $|f| > F_{\max}$ , где  $F_{\max}$  - действительная константа (на практике  $F_{\max} = 7000$  Гц), которую мы будем называть *максимальной частотой речевых сигналов*

$$4^\circ \forall p \in B \quad \text{sen}(\Pi(p)), \text{Si} \sim s \Rightarrow \text{Si} \in \Pi(p)$$

$$5^\circ \exists \text{ константы } T_{\min}, T_{\max}: 0 < T_{\min} < T_{\max}: \forall p \in S \quad \forall \text{sell}(p), \forall t \in \mathbb{R}$$

$$\exists t \in [T_{\min}, T_{\max}], \text{ такое что звуковой сигнал } s_{[t, T/2, t+T/2]} \text{ удовле-}$$

$$\text{творяет свойству } O_{s_{[t, T/2, t+T/2]}}(f) \neq 0 \text{ при } |f| > F_{\max}. 6^\circ \forall p \in B \quad \text{sell}(p) \Rightarrow$$

$$s \neq 0 \text{ (т.е. } s \text{ - не тождественный } 0)$$

- $T: S \times \mathbb{R} \rightarrow \mathbb{R}$ , - функция, задающая размер *окна анализа*,  $T_s(t)$  - окно анализа сигнала  $s$  в момент времени  $t$ , определяемая как

$$T_s(t) = \min\{T \in [T_{\min}, T_{\max}]: O_{s_{[t-T/2, t+T/2]}}(f) \neq 0 \text{ при } |f| > F_{\max}\}.$$

Отрезок  $[t-T_s(t)/2, t+T_s(t)/2]$  будем называть *окном анализа*.

- $r: S \rightarrow 2^S$  - *распознающая функция*, определяемая как  $r(s) = \{p \in V: z \in \Pi(P)\}$ ;

- $\Pi: SXN \rightarrow R$  - *отображение фонемной разметки* звуковых сигналов, такое что  $Vp \in B$ ,  $p = b_1 b_2 \dots b_m$ ,  $V \text{sell}(p)$  выполнено  $0 = u(s, 1) < u(s, 2) < \dots < u(s, 2m+2) = |s|$ ,  $u(s, n) = 0$  при  $n > 2m+2$  и справедливы свойства:

$$1^\circ \text{ Si} \sim s \Rightarrow |u(s_i, n)| = |u(s, n)| \quad \forall n \in N.$$

$$2^\circ \forall j = 1, 2, \dots, m \quad \forall t \in [i(s, 2j), i(s, 2j+1)] \quad s_{[t-T(s)/w(s)/2, t+T(s)/w(s)/2]} \in (b_j)$$

$$3^\circ \quad \forall j = 1, 2, \dots, m+1 \quad \forall t \in [f_j(s, 2j-1), i(s, 2j)], \text{ выполнено либо } s_{[t-T(s)/2, t+T(s)/2]} \in (b_N) \cup (b_j), \text{ либо } s_{[t-T(s)/2, t+T(s)/2]} \in (b) \text{ при всех}$$

$b \in V$  (здесь  $b_0 = b_{m+1} = \_$  - символ паузы)

4° Если  $P = b$  - фонетическое слово из одной буквы, а  $\text{sell}(P)$  - произнесение этого слова, то  $f_i(s, 1) = u(s, 2) = 0$ ,  $\Pi(s, 3) = |u(s, 4)| = |s|$ .

Будем называть звуковой сигнал  $s$  *произнесением* фонетического слова  $p \in V$ , если  $b \in \Pi(p)$ .

Звуковые сигналы, являющиеся произнесением какого-то фонетического слова, будем называть *речевыми сигналами*. Далее, если противное не оговорено особо, символом  $S$  будем обозначать множество всех речевых сигналов.

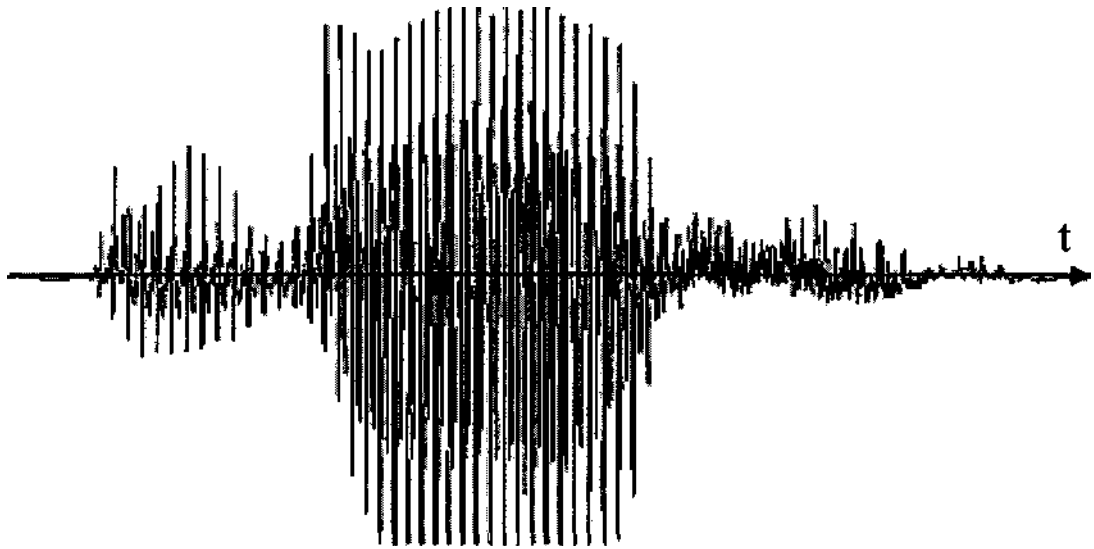


Рисунок 1.2 Пример речевого сигнала, являющегося произнесением слова «мыс»

Для каждого речевого сигнала  $s$ , являющегося произнесением некоторого фонетического слова  $(\mathbb{Z}e\mathbb{Z})$ , будем говорить, что временные отрезки вида  $[|Li(s, 2i)|, |Li(s, 2i+1)|]$  соответствуют произнесению  $i$ -й фонемы фонетического слова  $(\mathbb{Z})$ , а интервалы вида  $(|i(s, 2i+1)|, |j(s, 2i+2)|)$  соответствуют переходу от  $i$ -й к  $i+1$ -й фонеме фонетического слова  $p$  (интервал  $(|Li(s, 1)|, |Ll(s, 2)|)$  будем называть *переходом от паузы к первой фонеме слова*, а интервал  $(|Li(s, 2|P|+1)|, |Li(s, 2|p|+1)|)$  - *переходом от последней фонемы к паузе*). Величину  $|Li(s, 2i+1)| - |Li(s, 2i)|$  будем называть *длительностью произнесения  $i$ -й фонемы в речевом сигнале  $s$* , а величину  $|Li(s, 2i+2)| - p(s, 2i+1)|$  — *длительностью перехода от  $i$ -й фонемы к  $i+1$ -й*.



образом (у него не меняется тип голоса, акцент, словарный запас и т.п.), можно считать, что для каждого отдельно взятого человека и фиксированного словаря команд речевая модель единственна. Поэтому каж-

дое из отображений  $t$ ,  $\Pi$ ,  $g$  фиксировано для данной речевой модели и, в силу связи с реальными процессами произнесения слов и звуков и восприятия речи, обладает рядом описанных выше свойств, содержательный смысл которых состоит в следующем.

Отображение транскрибирования  $t$  задается правилами "озвучивания" отдельных слов естественного языка, поэтому для каждого типа акцента, диалекта и т.п. отображение  $t$  задано единственным образом. Для каждого естественного языка правила, задающие отображение  $t$ , известны. Пример таких правил для русского языка приведен в приложении 1.

Отображение произнесения  $\Pi$  задает различные варианты произнесения отдельных звуков, сочетаний звуков, фонетических слов данным носителем языка. Поэтому отображение  $\Pi$  единственно и обладает рядом естественных свойств.

Свойства  $1^\circ$ - $6^\circ$  содержательно означают, что у любого фонетического слова естественного языка есть хотя бы одно произнесение ( $1^\circ$ ); один и тот же речевой сигнал не может быть произнесением различных фонетических слов ( $2^\circ$ ); речевые сигналы содержат в себе только колебания с частотами меньше некоторой частоты  $F_{\max}$  ( $3^\circ$ ); если какой-то речевой сигнал является произнесением некоторого слова, то все эквивалентные ему речевые сигналы также являются произнесениями того же слова ( $4^\circ$ ); для любого речевого сигнала в любой момент времени существует окрестность этого момента времени, имеющая длину в заданных пределах, которая содержит в себе только колебания с частотами меньше частоты  $F_{\max}$  ( $4^\circ$ ); произнесение слова не может быть нулевым сигналом ( $5^\circ$ ).

Распознающая функция  $g$  моделирует процесс восприятия речи человеком — она восстанавливает фонемную транскрипцию слова по любому его произнесению.

Функция размера окна анализа  $T$  определяется как длина минимального окна, такого что сигнал, ограниченный на него, имеет ограниченный спектр.

Наконец, функция фонемной разметки ( $i$  позволяет в речевом сигнале выделить отрезки времени, соответствующие как произнесению отдельных фонем, так и переходу между фонемами. Эквивалентные сигналы имеют одинаковую фонемную разметку.

Зафиксируем речевую модель  $(A, A, B, B, S, x, Y, T, [i, r])$ .

**Определение 1.4.** Функцию  $0: S \times R \rightarrow S$ , определенную как

$$0(s, t)(i) = s(t + r), X \in \Pi T_s(t)$$

, будем называть *оконным отображением*, иначе

ем.

**Определение 1.5.** Отображение  $d: S \times R \rightarrow R^M$ , где  $d(s, t) = d'(0(s, t))$ , где  $d': S \rightarrow R^M$ , - некоторая функция, будем называть *функцией описания сигнала*, а отображение  $d'$  - *функцией описания окон*, если при любом фиксированном  $s$  функция  $d(s, t)$  непрерывна по  $t$  на отрезке  $[t_0(s), t_i(s)]$ .

При фиксированном  $s$  функцию времени  $d_s: [t_0(s), t_i(s)] \rightarrow R^M$ , где  $d_s(t) = d(s, t)$ , будем называть *описанием звукового сигнала  $s$* . Согласно определению 1.5,  $d_s(t)$  - непрерывная  $M$ -мерная действительностнозначная функция времени.

**Определение 1.6.** Функцию описания  $d = d' \circ 0$  будем называть *нормальной*, если отображение  $d'$  инъективно.





Описание речевого сигнала с помощью нормальной функции описания также будем называть *нормальным*.

Содержательно это означает, что нормальными мы будем называть функции описания, для которых по описанию речевого сигнала в данный момент времени однозначно восстанавливается звуковой сигнал на окне анализа.

Приведем пример функции описания звуковых сигналов, основанной на операторе преобразования Фурье.

**Пример 1.1.** В качестве примера функции описания возьмем отображение, сопоставляющее окну анализа вектор значений преобразования Фурье на равноотстоящих значениях частот  $7_M F_0, 2/M F_0, \dots, M/M F_0$ :  
 $c_{\Gamma T} = (\langle X \rangle_s(V_M \% \Phi_8 (2/M F_0), \dots, \Phi, (\% F_0)), \text{ где } F_0 > F_{\max}.$

В §1.2 будет показано, что введенная таким образом функция описания для достаточно больших  $M$  является нормальной (утверждение 1.4).

Пусть  $s$  - некоторый речевой сигнал. Пусть  $d_s$  - описание этого речевого сигнала, т.е. функция зависимости  $M$ -разрядного вектора действительных чисел от времени. Рассмотрим все возможные произнесения всех фонем.

**Определение 1.7.** Образом  $V(b)$  фонемы  $b \in B$  для данной функции описания  $d$  будем называть множество  $D(b) = \{d_s(t) \mid b \in \Pi(B), t \in [t_0(s), t_1(s)]\} \subset \mathbb{R}^M$

**Утверждение 1.1.** Образы различных фонем для нормальной функции описания не пересекаются. Доказательство.

Пусть  $d$  - нормальная функция описания сигналов. Предположим  
 противное. Тогда существуют две различные фонемы  $b_1$  и  $b_2$ , два  
 сигнала  $S_1$  и  $S_2$ , являющиеся произнесениями этих фонем, и  
 Два МоМента Времени  $t_1$  и  $t_2$ , Такие ЧТО  $t_1 \in [t_0(S_1), t_1(S_1)]$ ,  $t_2 \in [t_0(S_2), t_1(S_2)]$  и  
 $d(S_{b_1 t_1}) = d(S_{b_2 t_2})$ . Отсюда, учитывая инъективность отображения  $d'$ , полу-  
 чаем:  $z_f(u+x) = B_2(z_2+x) \quad \forall f \in [-T(t)/2, T(t)/2]$ , где  $T(t) = \min(T_{S_1}(t), T_{S_2}(t))$ ,  
 пусть, для определенности,  $T(t) = T_{S_1}(t)$ . По определению для любого ре-  
 чевого сигнала  $s$  и момента времени  $t$   $T_s(t) = \min\{T: T > T_{\min},$   
 $\langle E \rangle(s_{[t-T/2, t+T/2]})(f) = 0 \text{ при } f > F_{\max}\} \Rightarrow$  для сигнала  $s_2$   $T_{S_2}(t) < T(t) \Rightarrow T_{S_2}(t) = T_{S_1}(t)$ .

Однако по свойству 2° функции  $\wedge$  каждый из сигналов  $S_i(t)$ ,  $i=1,2$ ,  
 ограниченный на отрезок  $\left| \begin{array}{c} 2 \\ 2 \end{array} \right|$ , является произнесением со-  
 ответствующей фонемы  $b_i$ . Поскольку эти сигналы совпадают, по свой-  
 ству 2° функции  $\Pi$  должны совпадать и фонемы  $b_1$  и  $b_2$ . Мы пришли к  
 противоречию.

Утверждение доказано.

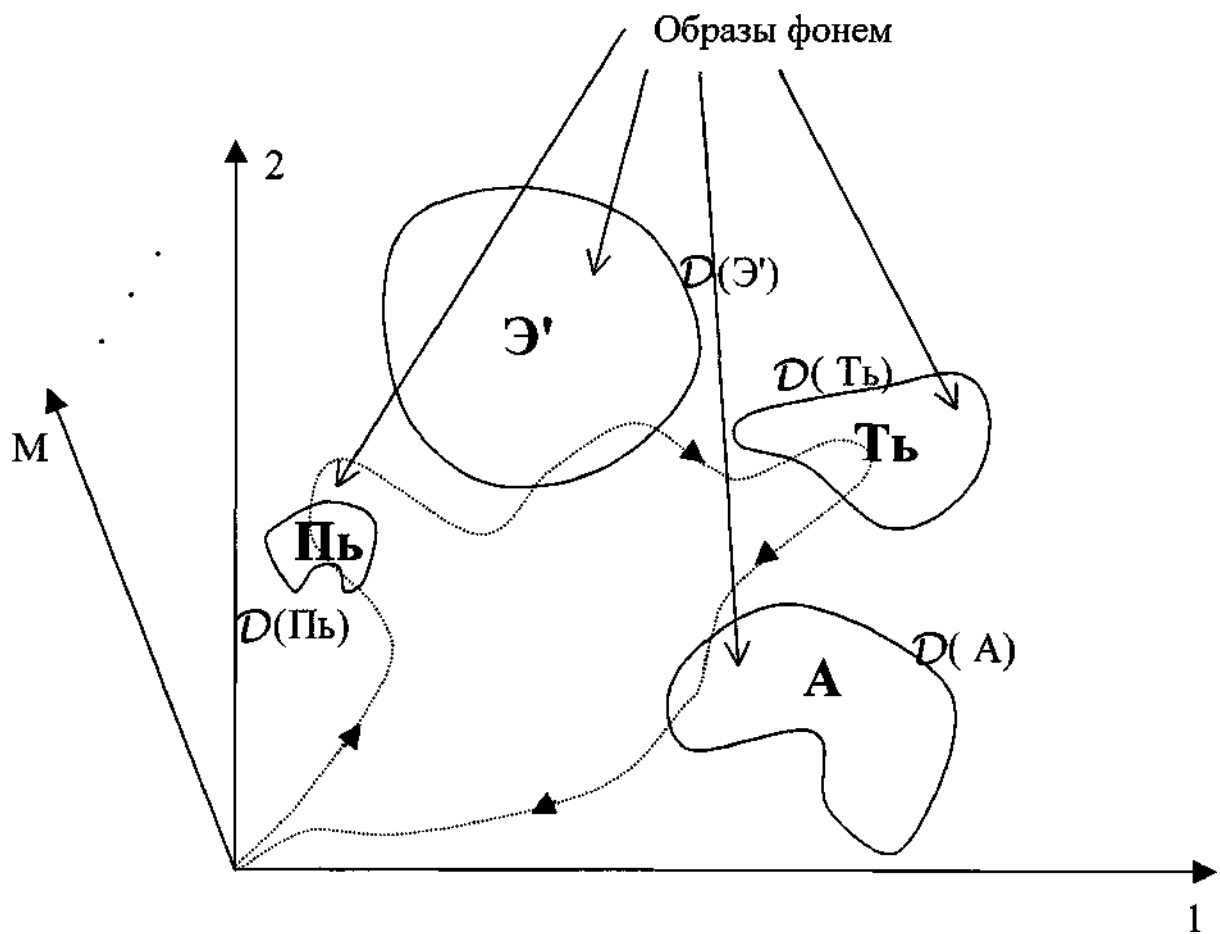


Рисунок 1.4. Описание произнесения слова ПЕТЯ и образы фонем из транскрипции этого слова для нормального оператора описания сигналов.

## §1.2. Дискретизация

Перейдем теперь от непрерывного времени к дискретному. Зафиксируем некоторую действительную величину  $\delta > 0$ , которую назовем периодом дискретизации. Выделим на оси времени точки  $\{n\delta + \phi, n \in \mathbb{Z}\}$ , где  $\phi$  - некоторая действительная константа (начальный сдвиг), т.е. точки, расположенные на расстоянии  $\delta$  друг от друга. Известная из теории связи теорема Котельникова [32] утверждает, что для того, чтобы по точкам, равноотстоящим друг от друга на величину  $\delta$ , можно было однозначно восстановить сигнал  $s(t)$ , необходимо, чтобы было выполнено условие  $D_s(f) = 0$  при  $|f| > 1/(2\delta)$ . Если  $\delta < 1/(2F_{\max})$ , то  $1/(2\delta) > F_{\max} \Rightarrow$  по свойству 3° функции произнесения  $D_s(f) = 0$  при  $|f| > 1/(2\delta)$ .

Следовательно, любой речевой сигнал может быть восстановлен единственным образом по последовательности равноотстоящих отсчетов  $\{s(n\delta + \phi), n \in \mathbb{Z}\}$ , если  $1/\delta > 2F_{\max}$ , т.е.  $\delta < 1/(2F_{\max})$ .

**Замечание 1.1.** Согласно теореме Котельникова [32], для восстановления звукового сигнала  $s$  по значениям  $s'(n)$  этого сигнала в точках  $\{n\delta + \phi, n \in \mathbb{Z}\}$  используют формулу:

$$s(t) = \sum_{n=-\infty}^{+\infty} s'(n) \frac{\sin \pi(n\delta - t - \phi)}{\sin \pi(n\delta - \phi)} \quad (1)$$

**Определение 1.8.** Дискретным звуковым сигналом будем называть любое отображение  $s': \mathbb{Z} \rightarrow \mathbb{R}$ , такое что  $\exists z_0, \exists \delta > 0 : s'(z) = 0 \quad \forall z \in [z_0, z_0 + \delta]$ .

Понятия начала  $t_0(s')$ , конца  $t_1(s')$  и длины  $|s'|$  дискретного звукового сигнала определяются аналогично соответствующим понятиям для непрерывных звуковых сигналов.

Обозначим через  $S'$  множество дискретных звуковых сигналов, а через  $R^+$  - множество неотрицательных действительных чисел.

**Определение 1.9.**  $(\delta, \phi)$ -дискретизацией назовем отображение

$D_{\delta, \phi}: S \times R^+ \times R \rightarrow S'$ , сопоставляющее звуковому сигналу  $s$  дискретный звуковой сигнал  $S_{\delta, \phi} = A_{\delta, \phi}(s): \forall n \in \mathbb{Z} \quad S_{\delta, \phi}(n) = s(n\delta + \phi)$ , где  $\delta \in R^+$  - период дискретизации,  $\phi \in R$  - начальный сдвиг дискретизации.

Дискретным речевым сигналом (для фиксированных  $\delta$  и  $\phi$ ) будем называть любую  $(\delta, \phi)$ -дискретизацию любого речевого сигнала. Будем считать, что период дискретизации  $\delta$  выбран так, что он гарантирует однозначность обратного восстановления непрерывного сигнала из дискретного (это возможно по свойству 3° функции произнесения). При этом по дискретному речевому сигналу  $s'$  можно восстановить непрерывный сигнал  $s$  при заданных  $\delta$  и  $\phi$  по формуле (1), т.е. существует отображение  $A'^{\delta, \phi}$ , обратное к  $A_{\delta, \phi}$ . Отображения  $A_{\delta, \phi}$  и  $D'^{\delta, \phi}$  при  $\phi=0$  будем обозначать через  $A_{\delta}$  и  $A'^{\delta}$  соответственно.

Таким образом, для каждого фиксированного  $\delta < 1/(2F_{\max})$  и фиксированного сдвига  $\phi$  между множеством речевых сигналов  $S$  и множеством дискретных речевых сигналов (будем его обозначать через  $S_{\delta, \phi}$ ) существует взаимнооднозначное соответствие, задаваемое отображением  $A_{\delta, \phi}$ .

**Утверждение 1.2.** Для любых  $\delta, \phi \quad S_{\delta, \phi} = S_{\delta}$ , е-

Доказательство.

Пусть  $s' \in S'$ . Тогда  $\exists s \in S: s = A_{\delta, \phi}(s')$ , такой что  $\forall k \in \mathbb{Z} \quad s(k\delta + \phi) = s'(k)$ . Рассмотрим сигнал  $Y: s(t) = s(t - Q + q) \quad \forall t \in R$ . Очевидно, что звуковой сигнал  $S'' \sim s$ , поэтому он также является речевым сигналом. Рассмотрим



$s'' = A_5, e(?)$ .  $\forall k \in \mathbb{Z} \quad s''(k) = s(k+9) = s(k-9+9) = s(k) = s'(k) \Rightarrow s' = s''$   
 $\Rightarrow s' \in S_{s,e} \Rightarrow S_{s,e} = S_{s,e}$ .

Аналогично доказывается, что  $S_{s,e} = S_{s,e}$ . Значит,  $S_{s,e} = S_{s,e}$ .

Утверждение доказано.

Следовательно, в обозначении множества дискретных речевых сигналов второй параметр - величину сдвига - можно опускать и обозначать это множество через  $S_s$ -

Далее, если обратное не оговорено особо, будем для упрощения опускать штрих в обозначении дискретных звуковых сигналов.

Через  $S_{[n,m]}$ ,  $n, m \in \mathbb{Z}$ , будем обозначать *ограничение дискретного сигнала  $s$  на отрезок  $[n,m]$* , т.е. дискретный звуковой сигнал, определенный следующим образом:

$$S_{[n,m]}(i) = s(i) \quad \forall i = n, n+1, \dots, m$$

$$S_{[n,m]}(i) = 0 \quad \forall i < n \text{ и } i > m$$

Введенную нами в §1.1 речевую модель (определение 1.3) можно переопределить для дискретного времени. Далее, если противное не оговорено особо, будем рассматривать везде дискретизацию с нулевым сдвигом ( $\phi=0$ ).

Итак, пусть для данного естественного языка заданы некоторая речевая модель  $(A, B, S, x, n, T, f, r)$  и период дискретизации  $\Delta < 1/(2F_{\max})$ .

**Определение 1.10.** *Речевой моделью с дискретным временем назовем декартову  $\langle A, B, S, T, n, T, |J_s, r \rangle$ , где*

- $A, B$  - конечные алфавиты букв и звуков соответственно;
- $D \subset A^*$  - словарь команд;



- $SaB^*$  - словарь звуковых команд;
- $Sg$  - множество дискретных речевых сигналов;
- $T:A^* \rightarrow B^*$  - отображение транскрибирования;
- $\Pi_5: S^2 \rightarrow S^2$ ,  $\Pi_5 = A^5 \circ \Pi$  - отображение произнесения;
- $T: Sg \times Z^N$ ,  $Ts_5(n) = [Ts(n_5)/5]$  - функция, задающая размер *окна анализа*.

- $rgiSs \rightarrow \dots, rg^{\wedge}roAs^{n1}$  - *распознающая функция*;
- $(o^{\wedge}SgxN - ^{\wedge}Z$  - *отображение фонемной разметки* звуковых сигналов, такое что  $V(3G3, (3=bib_2 \dots b_m, VSSGIT^P)$  выполнено  $0=M_6(S_6, 1) < Li_5(s_5, 2) < \dots < i_6(s_5, 2m+2) = |s_s|, i_6(s_5, n) = 0$  при  $n > 2m+2$ .

Отображения  $j$ , и  $(Li_5$  связаны соотношением  $V_i = 1, 2, \dots, 2m+2$   
 $jLi_5(s_5, i) = [(|u(A_5^{n1}(s)_5 i)/5]$ .

Свойства функции произнесения  $\Pi_5$  и отображения фонемной разметки  $\Pi_5$  для дискретного случая легко получаются из соответствующих свойств для непрерывного случая.

Свойства функции  $\Pi_5$ :

$$1^\circ V(3GB \Pi_5(p))^* 0,$$

$$2^\circ p \in P_2 \in B, P_1 * p_2 \Rightarrow \Pi_5(P_1) \cap \Pi_5(P_2) = \emptyset \quad 3^\circ, 5^\circ \text{ для}$$

дискретного случая теряют смысл  $4^\circ Vpe_8$

$$s_5 Gn_5(P)_5 s'_6 \sim s_s \Rightarrow s'_6 \in \Pi_6(P) \quad 6^\circ VpeSs_6 Gn_5(P) \Rightarrow s^{\wedge} 0$$

Свойства функции  $\backslash x \mathfrak{S} \backslash$

$$2^\circ Vj=1..m \quad Vn=(i_6(s_5, 2j)+1, \quad j i_5(s_5, 2j)+1, \dots \quad |i_5(s_5, 2j+1)-1$$

$$S8[n-[T(s)/2], n+[T(s)/2]]^c \Pi_5(B_{II})$$

$3^\circ \forall j=1..m+1 \forall n=|i|(s_5, 2j-1), n \in (s_5, 2j-1)+1, \dots |L_6(s_6, 2j)|$ , выполнено, что либо  $S_5[n, [T(s)/2], n, [T(s)/2]] \in n_5(b_j, i) \cup n_5(b_j)$ , либо  $\forall b \in B$   $S_5[n, [T(s)/2], n, [T(s)/2]] \wedge n_5(b)$  (здесь  $b_0 = b_{m+1} = \_$  - символ паузы)  $4^\circ$  Если  $p \in B$  - звуковая команда из одной буквы, а  $sg \in \mathcal{G}(p)$  - произнесение этой команды, то  $|Ag(sg, 1)| = 1115(85, 2) = 0$ ,  $(a_6(s_5, 3) = |A_6(s_5, 4)| = |j s_5|$ .

Как и для непрерывных речевых сигналов, для каждого дискретного речевого сигнала  $sg$ , являющегося произнесением некоторой звуковой команды (Зс73, будем говорить, что отрезки натурального ряда вида  $[|Hg(sg, 2i)|, |i(ss, 2i+1)|]$  соответствуют произнесению  $i$ -й фонемы звуковой команды  $p$ , а отрезки натурального ряда вида  $[|i_b(s_b, 2i+1)+1|, |Hg(sg, 2i+2)-1|]$  соответствуют переходу от  $i$ -й к  $i+1$ -й фонеме звуковой команды  $p$ . Величину  $p, g(sg, 2i+1) - |Ag(sg, 2i)+1|$  будем называть длительностью произнесения  $i$ -й фонемы в дискретном речевом сигнале  $Sg$ , а величину  $(i_5(ss, 2i+2) - |Li_5(sg, 2i+1)-1|$  — длительностью перехода от  $i$ -й фонемы к  $i+1$ -й.

Определение 1.11. Отображение  $d\$: S \times N^M \rightarrow R^M$  будем называть функцией описания дискретных речевых сигналов, если  $\forall s \in S, \forall i \in Z$  выполнено  $dg(s_6, i) = d(Ag''(sg), i_5)$ , где  $d: S \times R^M \rightarrow R^M$  - некоторая функция описания непрерывных речевых сигналов.

Содержательно это определение означает, что для описания дискретного речевого сигнала мы сначала восстанавливаем из него непрерывный речевой сигнал, а затем применяем какую-либо функцию описания к непрерывному сигналу.

**Замечание 1.2.** Как в непрерывном, так и в дискретном случае можно считать, что область значений функции описания является многомерным комплексным пространством. Для приведения комплексного случая к действительному нужно действительные и мнимые части значений функции описания рассматривать как отдельные координаты многомерного действительного пространства.

Поскольку по определению любая функция описания непрерывных речевых сигналов представляется в виде  $d = d^{\circ} 0_5$  то  $Vs \S GS5, ViGZ$  выполнено

$$de(S5, i) = \wedge(0(\wedge(85), 16)).$$

**Определение 1.12.** Функцию  $0 \S : S_5 \times Z \wedge S_5$ , определенную как

$$0_5(s_5, z)(z') = s_5(z + z'), \quad z' \in T \setminus T^* \quad \text{ЗД}$$

, будем называть дискретным

$O_5$  иначе

оконным отображением.

**Утверждение 1.3.**  $\forall s_5 \in S_5 \forall i \in Z \quad 0(A_5^{-1}(s_5), i) = A_5^0(s_5, i)$ .

Доказательство

Возьмем любое  $t \in R$ , такое что  $t = kS$ ,  $k \in \mathbb{Z}$ .

Если  $t \in [-r_{A_5}, (iS)/2, \Gamma_{A_5}, (iS)/2]$ , то по определению 1.12 имеем

$0(A_5^{-1}(s_5), i)(t) = A_5^{-1}(s_5)(t + iS) = s_5(k + i)$ , в то время как  $A_5^0(s_5, i)(t) = 0_5(s_5, i)(k) = s_5(i + k)$ , откуда  $0(A_5^{-1}(s_5), i)(t) = A_5^{-1}(0_5(s_5, i))(t)$  для  $t$  из этого отрезка.

Если же  $t \in [-r_{A_5}, (iS)/2, \Gamma_{A_5}, (iS)/2]$ , то по определению 1.4 оконного отображения  $O_5$  имеем  $0(A_5^{-1}(s_5), i)(t) = 0$ , в то время как  $A_5^{-1}(0_5(s_5, i))(t) = 0_5(s_5, i)(k) = 0$  по определению 1.12.



Звуковой сигнал  $A_8 \sim (O_5(s_5, i))$  удовлетворяет свойству  $\Phi_{d \leftarrow \langle \cdot, \cdot \rangle}^*$ ,  $(\cdot, \cdot) = 0$  при  $|\mathbf{f}| > F_{\max}$  по определению функции  $A_5^{-1}$ , а звуковой сигнал  $O(A_5^{-1}(s_5), i_5)$  удовлетворяет этому свойству в силу свойства 5° функции произнесения. По теореме Котельникова два сигнала с ограниченным спектром, совпадающие на множестве точек  $t = k_5$ ,  $k \in \mathbb{Z}$ , совпадают всюду.

Утверждение доказано.

**Следствие 1.1.** Любая функция описания  $d_8$  дискретных речевых сигналов может быть представлена в виде  $d_8 = d' \circ \theta_8$ , где

$d'_8: S \rightarrow \mathbb{R}^M$  - дискретная функция описания окон;

$\theta_8$ : дискретное оконное отображение.

Доказательство

По утверждению 1.3  $V_{s_8} \in S_8, V_{i \in \mathbb{Z}}$   
 $d_8(s_5, i) = d(A_8^{-1}(s_5), i_5) = d'(O(A_5^{-1}(s_5), i_5)) = d'(A_8^{-1}(O_5(s_8, i)))$ , откуда  $d_8 = d' \circ A_8^{-1} \circ O_5$ . Обозначим  $d'_8 = d' \circ A_5^{-1}$ . Отображения  $d'_8$  и  $\theta_8$  обладают, очевидно, искомыми свойствами.

Следствие доказано.

Содержательно следствие 1.1 означает, что для описания дискретного сигнала нет необходимости восстанавливать непрерывный сигнал из дискретного и применять непрерывную функцию описания. Существует дискретный аналог функции описания и дискретный аналог оконного отображения, которые позволяют определить функцию описания для дискретного случая аналогично непрерывному случаю как композицию оконного отображения и функции описания окон.



С другой стороны, по формуле обратного преобразования Фурье для дискретного времени ([32])

$$s_5(n) = f - \frac{1}{\Phi - C/Y^2} * \wedge$$

171 \

~2S

Приравнявая два выражения для  $s_5(n)$ , получим соотношение

$$\begin{matrix} 1 & +\infty & 1 \\ 0 & * - \langle \rangle & 0 \end{matrix}$$

Поскольку в нашем случае рассматриваются значения  $f$  из промежутка  $(0, F_0]$ ,  $F_{\max} < F_0 < 1/28$  и по свойству 3° отображения  $\Pi$   $\Phi(ii)=0$  при  $|f| > F_{\max}$ , то в сумме (\*) только одно слагаемое отлично от нуля (это слагаемое соответствует случаю  $k=0$ , т.к. для  $k>0$   $f+k/8 > 0+2F_{\max}=2F_{\max}$ , а для  $k<0$   $f+k/5 < F_0-2F_0=-F_0 < -F_{\max}$ )

$$\Rightarrow \Phi^L(f) = \frac{1}{d} \Phi(f), \text{ или, с индексами } s \text{ и } s\$, \Phi^L_{3\$}(:E) = \Phi_s(r)/8.$$

Последнее равенство справедливо для любых звуковых сигналов  $s$  и  $s_s = A_5(s)$ , для которых  $\Phi_5(r)=0$  при  $|f| > F_{\max}$ , в частности, для удовлетворяющих этому свойству по свойству 5° отображения  $\Pi$  звуковых сигналов  $C > 5(S_5)$  и  $0(s)$ .

Следовательно, выполнено равенство  $0 * 0g(S_5)(f)8 = \Phi \circ f \otimes$ , откуда  $dg(s_5, i) = d^f(0(s, i8))$ , что завершает доказательство утверждения.

Утверждение доказано.

**Определение 1.13.** Функцию описания дискретных речевых сигналов будем называть *нормальной*, если нормальной является соответствующая ей функция описания непрерывных речевых сигналов.

**Утверждение 1.5.** Функция описания непрерывного речевого сигнала, порожденная функцией описания окон  $d'(s) = (\langle D_s(\cdot/2M5) \rangle \Phi_8(hus), \blacksquare - , \Phi_5(M/2M6))$  (здесь  $F_0 = 1/2\$$ ) и соответствующая ей дискретная функция опи-





сания, порожденная функцией описания окон  $\wedge^{\wedge}(\Phi^{\wedge}3C^{\wedge}мб)^{\wedge}$ ,  $\Phi^{Д85}(^2/2M5)^5$ , ...,  $\Phi^{Д85}(^M/2Mб)^{\S}$ , являются нормальными функциями опи-сания при  $M > T_{\text{тах}}/8$ .

### Доказательство.

Рассмотрим описание произвольного речевого сигнала  $s: R \rightarrow R$  с помощью функции описания  $d=d'oO$  в произвольной точке  $t_o$ . Из свойства 4° функции произнесения  $\Pi$  следует, что можно не ограничивая общности считать, что  $t_o = 5 \cdot i_o$ , для некоторого  $i_o \in \mathbb{Z}$  (действительно, если это не так, то вместо сигнала  $s$  можно рассмотреть эквивалентный ему речевой сигнал, смещенный по оси времени на  $t_o - [t_o/5]8$ ).

Покажем, что отображение  $d^f$  инъективно.

Обозначим через  $s'$  речевой сигнал  $s$  после применения оконного преобразования  $O$ :  $s'(t) = O(s, t_o)(t)$ , а через  $s'_s$  речевой сигнал  $sg$  после применения оконного преобразования  $Og$ :  $s'_s(n) = O_s(s, j_o)(n)$ .

$$\text{Тогда} \quad \Phi 8^{<T/2Mб)/5} = \Phi^{Д8.5}(^T/2M5) - \frac{-2mn}{2M} *$$

$$\frac{T_{S_s}(fr)}{t 4(*K^{2J})}$$

$$H = T_s * O_o$$

$$)e^{LM} - (\text{поскольку } M > T_{\text{тах}}/8, \text{ откуда } M > T_{S_s}(i_o)/2)$$

$$5 \gg e^{TM} = Y, s_s^f(n-M)e^{n \sim M} \quad n=0$$

$$e \sim 5 >; (и-ло e^{2M} \quad \partial y u \quad 1Y \pm J \&)$$

Применяя формулу обратного дискретного преобразования Фурье ([32]), получим:

$$-mm \quad 1 \quad \backslash \quad \text{Ч} \quad \kappa \quad \kappa \quad \frac{i2\pi k(n+M)}{2M}$$

$k=M, M+1, \dots, 2M-1$  положим  $\Phi^* \# \left( \frac{A}{2MS} \right) = \Phi^* \text{Ч} \left( \frac{2M-1-A}{2M8} \right)$  (из свойств пре-

образования Фурье для действительного сигнала, [32]).

Возвращаясь снова к непрерывному сигналу,  $s'(t) = O(s, t_0)(t) =$

$O(A V(^{\wedge})_5 t_0)(O =$  (по доказанному ранее утверждению)  $=$

$A \cdot j(O_d(s)_{9i_0})(t) = A V(s'_s)(O$ . Итак,  $s'(t) = A^{-1}(s'_s)(t)_9$  в то время как

$$\blacksquare \quad 1 \quad 2\wedge \quad \text{Л} \quad \frac{I\wedge k(n+M)}{L2M}$$

Таким образом, мы показали, что сигнал  $s'$  однозначно восстанавливается по значению  $d'(<v')$ , т.е. что отображение  $d^f$  - инъективно. Следовательно, порожденная этим отображением функция описания нормальна.

Утверждение доказано.

Пусть  $X$ -прямоугольная  $m \times n$  матрица,  $X^i$ ,  $i=1, 2, \dots, n$  - ее столбцы. Будем в таком случае писать, что  $X = \|X^1 | X^2 | \dots | X^n\|$ .

Матрицу  $d_5(s_6) = \|d_5(s_5, \text{to}(s_5)) | d_5(s_6, \text{to}(s_5)+1) | \dots | d_6(s_5, \text{ti}(s_8))\|$  размерности  $M \times |s_5|$  будем называть *описанием дискретного речевого сигнала  $S_5$* .

С целью экономии памяти в практических приложениях удобнее вместо матрицы  $d_g(s_5)$  описания речевого сигнала  $s_5$  использовать матрицу, составленную из некоторого подмножества множества столбцов матрицы  $d_6(s_5)$ .

Более строго, *описанием дискретного речевого сигнала  $s_5$  с шагом  $L$* ,  $L \in \mathbb{N}$ , будем называть матрицу размерности  $M \times \lfloor |s_5|/L \rfloor$ , столбцы кото-

рой соответствуют описанию сигнала в моменты времени, кратные  $L$ , т.е. являются подмножеством множества столбцов матрицы  $d_s(ss)$ :

$$d_{5,L}(S_6) = \{d_5(s_5, \frac{\Phi_6}{L} + 1, Z) | d_5(s_5, \frac{\Phi_6}{L} + 2, \dots, Z) | \dots | d_6(s_5, \frac{\Phi^*}{L}, L)\}$$

Очевидно, что  $d_5(s_5) = d_{5,L}(s_5)$ .

Пусть  $sg$  - некоторый дискретный речевой сигнал. Пусть  $d_{5,L}(s_5)$  - описание этого речевого сигнала с шагом  $L$ . Столбцы этой матрицы, рассматриваемые как координаты точек пространства  $R^M$ , задают некоторое множество точек в этом пространстве.

Рассмотрим теперь все возможные произнесения всех фонетических слов для данного естественного языка.

**Определение 1.14.** Дискретным образом  $\mathcal{L}(B)$  фонемы  $B \in V$  для данной функции описания  $d_s$  и шага описания  $L$  будем называть следующее подмножество пространства  $R^M$ :

$$t_{5,L}(b) = \{d_5(s_5)(z) | 8_6 \in \Pi_5(B), ZG[t_0(s_5), \wedge^{(3/4)}]\}.$$

Содержательно данное определение означает, что дискретный образ фонемы состоит из тех и только тех точек, соответствующих дискретным описаниям с шагом  $L$  всех произнесений этой фонемы, которые задаются дискретными моментами времени, кратными шагу  $L$ .

**Утверждение 1.6.** Для любых допустимых значений периода дискретизации  $5$  и шага описания  $L$  дискретный и непрерывный образы одной и той же фонемы совпадают.

Доказательство.

Очевидно, что каждый дискретный образ фонемы вложен в соответствующий непрерывный, т.е.  $D_{s,L}(b) \subset D(b)$ . Это следует из того, что



для моментов времени вида  $t = Ln5$  соответствующая индексу  $Ln$  точка описания дискретного сигнала принадлежит описанию соответствующего непрерывного сигнала.

Обратное утверждение следует из свойства 4° функции произнесения  $\Pi$ . Пусть  $b$  - произвольная фонема,  $x \in TXb$ ). Тогда существует непрерывный звуковой сигнал  $SGil(b)$  и момент времени  $t'$ :  $d(s, t') = x$ . Рассмотрим сигнал  $s'$ :  $s'(t) = s(t - t') \quad \forall t \in R$ . Очевидно, что  $s' \sim s \Rightarrow b' \in \Pi(b)$  (свойство 4° функции  $\Pi$ ) и  $fi(s, t) = fi(s', t) \quad \forall t \in [t_0(s), t_j(s)]$  (свойство 1° функции  $\Pi$ ).

$d(s', 0) = d(s, t') = x$ . Если обозначить через  $s'_8$  дискретный звуковой сигнал  $s'_8 \in As^8$ , то по определению  $d_s$  имеем  $d_s(s'_8, 0) = d(s', 0) = x$ . Момент времени 0 имеет вид  $k8L$ ,  $k \in Z$  (здесь  $k=0$ )  $\Rightarrow$  по определению  $D^8 Xb$  имеем  $x \in 0_{5, b}(b) \Rightarrow D^8 b \in 0_{5, b}(b)$ .

Итак,  $0_{5, b}(b) \subset D^8 b$  и  $D(b) \subset D^8 L(b)$ , следовательно,  $TXb = V_{b, L}(b)$ .

Утверждение доказано.

Содержательный смысл утверждения 1.6 следующий: как бы сильно мы не "прореживали" сигнал при его описании, образы фонем будут в точности теми же, поскольку мы имеем свободу при выборе начальной точки дискретизации и прореживания.

Поскольку  $1Xb = 0_{5, b}(b) \cap V$  допустимых  $5$  и  $L$ , индексы  $5$  и  $L$  можно опустить и в дальнейшем и в дискретном случае обозначать образы фонемы  $b$  через  $TXb$ ).

При дискретизации речевого сигнала при достаточно большом значении шага описания  $L$  может случиться, что множество точек описания сигнала, соответствующих дискретному описанию сигнала с шагом  $L$ , не будут пересекаться с каким-либо дискретным образом фонемы, несмотря на то, что соответствующая этому множеству фонема (символ алфавита

В) входит в произнесенное фонемное слово. Будем говорить в таком случае, что в описании данного (дискретного) произнесения данного фонемного слова указанная фонема "пропущена" (см. рисунок 1.5).

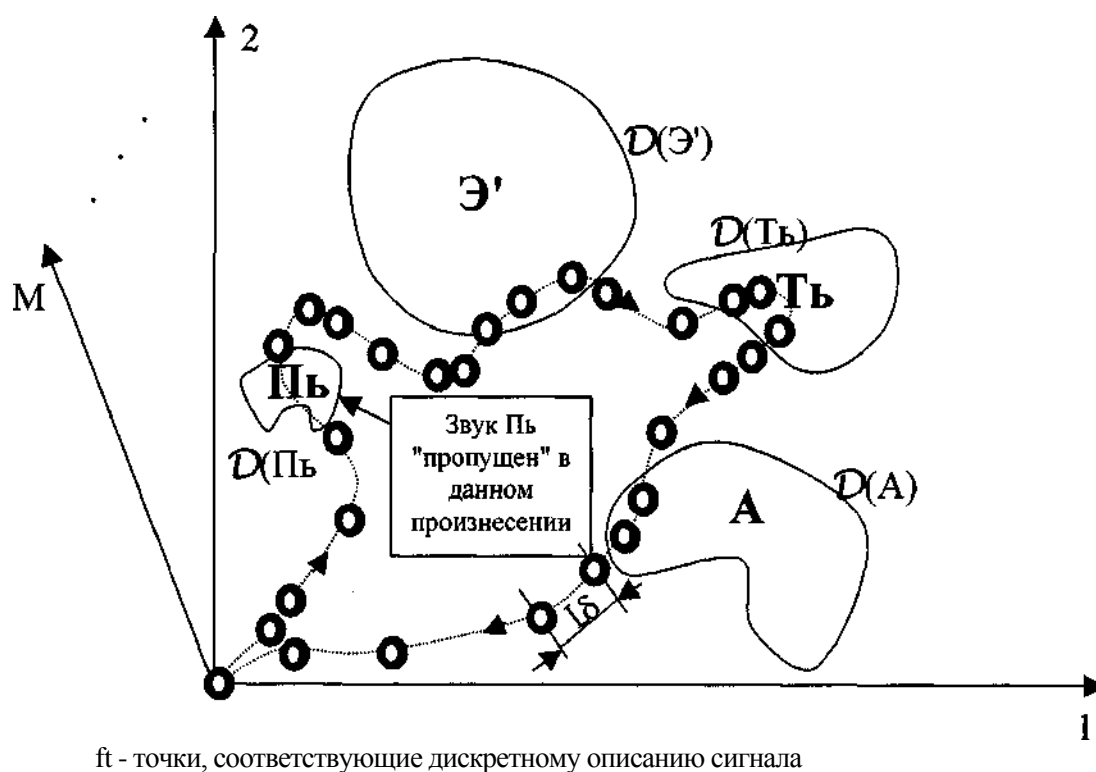


Рисунок 1.5. Описание с шагом  $L$  дискретного речевого сигнала, являющегося произнесением слова «Пе'тя»

### **§1.3 Формализация задачи распознавания речи.**

#### **Функция качества словаря команд.**

В предыдущем параграфе мы ввели понятие дискретного речевого сигнала и показали, что если период дискретизации  $\Delta$  выбран правильно, множества дискретных и непрерывных речевых сигналов находятся во взаимнооднозначном соответствии.

Понятие речевой модели было сформулировано как для случая непрерывного (определение 1.3), так и для случая дискретного времени (определение 1.10). Было показано, что операторы описания в непрерывном и дискретном случае устроены одинаково (следствие 1.1).

Мы показали также, что несмотря на то, что дискретное описание речевого сигнала берется с некоторым временным шагом, образы фонем и в дискретном, и в непрерывном случае в точности одни и те же (утверждение 1.6).

Таким образом, не ограничивая общности, в дальнейшем можно считать, что мы имеем дело с дискретными речевыми сигналами и рассматриваем их описания с шагом 1. Для упрощения записи индекс  $\Delta$  при обозначении элементов речевой модели с дискретным временем будем в дальнейшем опускать.

Поскольку распознающая функция  $\gamma$  восстанавливает по речевому сигналу речевые команды из  $\mathcal{C}$ , а не тексты команд из  $\mathcal{A}$ , на практике используются такие словари команд, транскрипции которых попарно различны (т.е. в словаре нету т.н. "омонимов" - одинаково звучащих слов). В таком случае без ограничения общности можно считать, что команды - это слова в алфавите  $\mathcal{B}$ , и исключить из нашей речевой модели алфавит букв  $\mathcal{A}$  и отображение транскрибирования  $\tau$ .

На практике системы распознавания речи работают не со звуковыми сигналами, а с их описаниями как с исходным материалом для обучения и распознавания речевых образов. Если через  $(R^M)$  обозначить множество прямоугольных матриц с  $M$  строками и произвольным количеством столбцов с элементами из  $R$  (иначе говоря, множество слов над континуальным алфавитом  $R^M$ ), то описание любого речевого сигнала является элементом множества  $(R^M)^*$ .

Часто бывает удобным разбить пространство  $R^M$  на конечное число непересекающихся подмножеств:  $R^M = \bigcup_{i=1,2,\dots,J_C} C_i$ ,  $C = \{c_i, i = 1, 2, \dots, J_C\}$ , которым можно сопоставить символы некоторого конечного алфавита  $S = \{s_1, s_2, \dots, s_{J_C}\}$  (алфавит  $S$  называется обычно *кодовой книгой*). В этом случае каждая матрица из  $(R^M)$  может быть закодирована словом в алфавите  $S$ , если каждому столбцу этой матрицы сопоставить символ из  $S$ , соответствующий элементу разбиения  $C$ , содержащему точку в  $R^M$ , координаты которой задает этот столбец.

Таким образом, понятие речевой модели, с учетом включения в модель функции описания, кодовой книги  $S$  и кодирующего отображения  $\Gamma$ , может быть, не ограничивая возможностей ее применимости, сформулировано следующим образом:

**Определение 1.15.** *Речевой моделью* называется совокупность 11 элементов  $\langle V, W, S, (R^M)^*, C, n, T, |i, d, r, r\rangle$ , где

- $V$  - конечный алфавит звуков;
- $W$  - словарь команд,  $W \subset V^*$ ;
- $S$  - множество дискретных речевых сигналов;
- $(R^M)^*$  - множество прямоугольных матриц над  $R$  с  $M$  строками;
- $C$  - кодовая книга (конечный алфавит);



- $\Gamma: \mathbb{R}^M \rightarrow C$  - кодирующее отображение, которое естественно обобщается на слова как  $\Gamma: (\mathbb{R}^M)^* \rightarrow C^*$ , если сопоставить каждой матрице  $y \in (\mathbb{R}^M)^* \setminus y = \|y^1 y^2 \dots y^n\|$  слово  $\Gamma(y) = \Gamma y^1 \Gamma y^2 \dots \Gamma y^n$  - последовательный код  $y$ .

- $d: S \times N \rightarrow R^M$  - функция описания речевых сигналов, имеющая вид  $d = d' \circ O$ , где  $d': S \rightarrow R^M$  - функция описания окон, а  $O: S \times Z \rightarrow S$  - оконное отображение;  $d$  также обобщается на сигналы как  $d: S \rightarrow K R^M$ , если сопоставить сигналу  $s$  его описание  $d_s = d(s) = \|d(s, 1) d(s, 2) \dots d(s, n)\|$ ;

- $I: \mathcal{B}^2$  - функция произнесения, обладающая свойствами:

1°VJ3eB Π((3)\*0,

$$2^\circ \text{ } p_1, p_2 \in \mathbb{N}, p_1 \wedge p_2 \Rightarrow \Pi(p_1) \cap \Pi(p_2) = \emptyset$$

4°VPeB sen<sup>s</sup>sen<sup>S</sup>CP)

$$6^\circ \text{VpGBsGn(P)} \Rightarrow s^0$$

- $T: S \times Z \rightarrow N$ , - функция, задающая размер *окна анализа*, такая что  $T_{\min} < T_s(n) < T_{\max}$ , где  $T_{\min}, T_{\max} \in M$  - некоторые натуральные числа,  $T_{\min} < T_{\max}$  (здесь  $T_s(n)$  - окно анализа сигнала  $s$  в момент времени  $n$ )

- ц:  $S_x N \rightarrow Z$  - отображение фонемной разметки звуковых сигналов, такое что  $UP \in B$ ,  $P = bib_2 \dots b_m$ ,  $VSGII(P)$  выполнено  $0 = \wedge (s, l) < p, (s, 2)$

$\langle \dots \langle a(s, 2m+2) = |s|, u(s, n) = 0 \text{ при } n > 2m+2 \text{ и справедливы свойства: } 1^\circ \text{ Si} \sim s$

$$\Rightarrow |\text{Li}(\text{si}, n) = |\text{i}(\text{s}, n) \quad \forall n \in \mathbb{N}. 2^\circ \quad \forall j=1..m \quad \forall n=\wedge(\text{s}, 2j)+1,$$

$$K_{s,2j)+1}, \quad \dots, \quad |a_{(s,2j+1)-1}$$

$$S[n-T(s)/2, n+T(s)/2]^c \cap \mathbf{b} \mathbf{j}$$

3°  $\forall j=1..m+1 \forall n=j(s, 2j-1), p, (s, 2j-1)+1, \dots \} J.(s, 2j)$ , выполнено, что

либо  $s_{[n.T(sy2),n+iT(sy2)]}e\Pi(bj.i)uri(bj)$ , либо  $VbeB$

$$S[n - \lceil T(s)/2 \rceil, n + \lceil T(s)/2 \rceil] \not\subseteq \Pi(B)$$

(здесь  $b_0 = b_{m+1} = \_$  - СИМВОЛ паузы)

4° Если  $(3=b$  - звуковая команда из одной буквы, а  $sell(P)$ , то  $|i(s,1)=|i(s,2)=1, |i(s,3)=|i(s,4)=|s|$ . •  $\Gamma: \Gamma \rightarrow 2^B$  - *распознающая функция*, определенная как  $\Gamma(y) = \{P \in B: y \in \Gamma(c_i(\Pi(P)))\}$ .

**Замечание 1.2.** Если образы фонемных множеств при отображении  $\Gamma$  не пересекаются, то  $\forall y \in \Gamma(c_i(\Pi(B))) | \Gamma(y)|=1$ , т.е. речевая команда однозначно восстанавливается по коду описания любого ее произнесения.

**Замечание 1.3.** Необходимым условием для того, чтобы  $\forall y \in \Gamma(c_i(\Pi(B)))$  выполнялось  $| \Gamma(y)|=1$ , является условие нормальности функции описания  $d$ .

Под *задачей распознавания речи* в речевой модели  $(B, B, S, (R^M)^{C, n, T, |u, d, r, r})$  будем понимать задачу восстановления функции  $\Gamma$  на основе конечного числа примеров произнесения каждой из звуковых команд из  $B$ .

Обычно задача распознавания речи решается с помощью *метода сравнения с эталонами*. Опишем этот метод более строго.

Итак, пусть словарь команд  $B = \{P_1, P_2, \dots, P_N\}$  - Сопоставим каждой команде  $P_j$  из  $B$  *эталон*  $e_j$  - элемент некоторого абстрактного множества  $E = \{e_1, e_2, \dots, e_N\}$ .

Содержательно эталон  $e_j$  команды  $P_j$  является компактным способом представления множества кодовых слов всех ее произнесений и

строится обычно на основе набора примеров ее произнесения - некоторого непустого множества кодовых слов  $\{y_1, y_2, \dots, y_{i-1}(i)\} \in \Gamma(\Pi(\mathbb{Z}))$ .

На декартовом произведении  $S^* \times E$  введем функцию *расстояния*  $p$ :  $S^* \times E \rightarrow \mathbb{R}$  между кодовыми словами и эталонами. Распознавание речи методом сравнения с эталонами производится по формуле  $\Gamma(E, p)(y) = \text{Argmin}(p\{y, e_i\})$ . Если  $\Gamma(E, p)$ , будем говорить, что имеют место

ошибки распознавания. Алгоритм распознавания речи, основанный на сравнении кодовых слов с эталонами из  $E$  с помощью расстояния  $p$ , будем обозначать через  $(E, p)$ .

Содержательно метод сравнения с эталонами позволяет находить эталоны, ближайшие к распознаваемому коду.

Качество  $q(E, p)$  словаря команд  $V$  для данного алгоритма распознавания  $(E, p)$  естественно определить, как  $1 - 2p_{\text{ош.}}$ , где  $p_{\text{ош.}}$  - среднее число ошибок распознавания команд из  $V$  алгоритмом  $(E, p)$  (более строгое определение функции качества будет дано для конкретных примеров автоматных алгоритмов распознавания речи в следующих главах).

Важно заметить, что каждой команде  $p$  соответствует множество кодовых слов описаний всех ее произнесений  $\Gamma(\text{cl}(\Pi(p)))$ , которое является конечным подязыком языка  $S$  всех кодовых слов. Автоматный подход к распознаванию речи состоит в том, чтобы представить язык кодовых слов каждой команды некоторым автоматом. В случае конечно-детерминированных автоматов (глава 2) мы погружаем язык кодовых слов команды в некоторый известный нам бесконечный регулярный язык, а в случае вероятностных автоматов (глава 3) - в множество всех слов  $V^*$  над алфавитом  $V$ , приписав при этом каждому слову языка кодо-



вых слов действительное число - его вероятность. Последний подход дает нам возможность рассматривать случай пересекающихся языков кодовых слов.

## ГЛАВА 2. ДЕТЕРМИНИРОВАННЫЕ АВТОМАТНЫЕ МОДЕЛИ §2.1

### Решение задачи локального распознавания речи.

Пусть задана речевая модель  $(B, V, S, (R^M)^*, C, \Pi, T, j, d, r, r)$  и решается задача распознавания речи в этой модели. Предположим, что разбиение  $C$  и соответствующая ей кодовая книга  $C$  таковы, что существует алгоритм распознавания речи  $(E, p)$ , восстанавливающий  $\gamma(y)$  без ошибок при условии, что для каждого  $p \in K^M$  мы умеем вычислять  $\Gamma(p)$  (далее в этой главе будет показано, что такое  $C$  существует). К сожалению, на практике вычисление  $\Gamma(p)$  для таких  $C$  оказывается сложной задачей. Существование "идеального" алгоритма  $(E, p)$  распознавания кодовых слов сводит, таким образом, задачу распознавания речи к задаче восстановления кодирующего отображения  $\Gamma$ . Поскольку  $\Gamma$  вычисляется отдельно для каждого элемента описания сигнала, задачу восстановления  $\Gamma$  назовем *задачей локального распознавания речи*.

Будем говорить, что задача локального распознавания речи решена, если каждой точке  $p \in R^M$  сопоставлено некоторое множество  ${}^x F(p) \subseteq C$ , такое что  $\Gamma(p) \in {}^x F(p)$ .

Для конструктивного задания функции локального распознавания  ${}^x F$  предположим, что задана функция  $\|/\|$ , которая умеет "отличать" элементы разбиения  $C$  попарно. Перейдем к формальным определениям.

**Определение 2.1.** Функцию  $\|/\|: C \times C \times R^M \rightarrow \mathcal{C}u\{?\}$  будем называть *отделяющей*, если для  $\forall p \in R^M$   $(0^\circ) \quad \|/\|_{cc}(p) \supseteq c$

$$(1^\circ) \quad \|/\|_{icj}(p) = \|/\|_{cji}(p) \text{ (СИММЕТРИЧНОСТЬ)}$$

$$(2^\circ) \quad \|/\|_{cm}(p, M_{ci, cj}, ?).$$

$$(3^\circ) p \in C_i^m / c_i \in C_j \text{Mci, ?}$$

Смысл функции  $i/c_i j$  следующий - она выдает ответ на вопрос, "какому из множеств  $c_i$ ,  $c_j$  принадлежит точка  $p$  при условии, что она лежит в объединении  $c_i \cup c_j$ . Ответы интерпретируются следующим образом:  $\forall i/c_i j(p) = C_i \Rightarrow p \in C_i$ ,  $\forall i/c_i j(p) = C_j \Rightarrow p \in C_j$ , '?' - Неизвестно.

**Определение 2.2.** Функцией локального распознавания  $\Upsilon$  будем называть функцию  $\Upsilon^*: R^M \rightarrow 2$  (здесь  $2^e$  - множество всех подмножеств алфавита  $C$ ), определенную для каждой точки  $p \in R^M$  следующим образом:

$$\Upsilon(p) = \{c \in C: \forall c' \in C, c' \neq c \Rightarrow i/c_c j(p) \in \{c, ?\}\}.$$

Смысл функции  $\Upsilon$  такой:  $c \in \Upsilon(p)$ , если гипотезу " $p \in c$ " невозможно опровергнуть при помощи отделяющей функции  $i/c_j$ .

**Определение 2.3.** Будем говорить, что функция  $W$  отделяет множества  $c_i$  и  $c_j$ ,  $W_{ij}$  (иначе говоря, множества  $c_i$  и  $c_j$  отделимы с помощью функции  $\Upsilon$ ), если для любой точки  $p \in R^M$  выполнены условия:

$$(1^\circ) p \in c_i \Rightarrow c_j \in \Upsilon(p)$$

$$(2^\circ) p \in c_j \Rightarrow c_i \notin \Upsilon(p)$$

Можно дать графовую интерпретацию вычисления функций  $i/c_j$  и  $\Upsilon$  для данной точки  $p$  из множества  $R^M$ . Итак, пусть точка  $p$  фиксирована. Построим ориентированный граф, вершины которого будут помечены символами алфавита  $C$ . Для каждой пары  $(c_i, c_j)$ :

если  $i/c_i j(p) = C_j$ , то проводим направленное ребро из вершины, помеченной символом  $C_i$  в вершину, помеченную символом  $C_j$ ;

если  $v/c_j(p) = C_b$  то проводим направленное ребро из вершины, помеченной символом  $C_j$  в вершину, помеченную символом  $C_b$ ; если  $v/c_j(p) = ?$ , никаких ребер не проводим.

**Утверждение 2.1.**  $s \in p) \Leftrightarrow$  из вершины, помеченной символом  $s$ , не исходит ни одного ребра.

Доказательство.

$s \in p) \Leftrightarrow \exists s' \in C, s' \wedge s$  выполнено: либо  $i/c_s(p) = ?$ , либо  $j/c_s(p) = b$   
 $\Leftrightarrow$  У вершины, помеченной символом  $s' \in C$ ,  $s' \wedge s$  выполнено: либо вершины, помеченные символами  $s$  и  $s'$ , не соединены ребром, либо в графе есть ребро, направленное от вершины, помеченной символом  $s'$  к вершине, помеченной символом  $s$ .

Утверждение доказано.

**Следствие 2.1.** Множество  $T(p)$  совпадает с множеством меток вершин графа, из которых не исходит ни одного ребра.

**Пример 2.1.** Пусть в речевой модели  $\langle B, B, S, (R^M)^*, C, n, T, (i, d, r, r) \rangle$  функция описания  $d$  - нормальная, а разбиение  $C$  пространства  $R^M$  на подмножества имеет вид  $R^M = \bigcup_{i=0,2..|D|} J V(b_i)$ , где  $\mathbb{N}_i = t(b_i)$  при  $i=1..|B|$  - образы

фонем из  $B = \{b_1, b_2, \dots, b_{|B|}\}$ , а  $c_0 = 0(b_0) = R^M \setminus \bigcup_{i=1,2..|D|} J V(b_i)$  - «межфонемное про-

странство». В этом случае можно считать, что  $C = B \cup \{b_0\}$ . Как будет показано ниже, для этого примера существует детерминированно-автоматный алгоритм распознавания речи, восстанавливающий речевые команды по кодовым словам описаний их произнесений.





**Пример 2.2.** Пусть задана речевая модель из примера 2.1,  $V=\{П,Р,И,М,Е\}$  (алфавит для слова ПРИМЕР),  $C=Bub_0$  и функция  $\|/$  в точке  $p$  задана следующей таблицей:

	П	Р	И	М	Е
П	П	9	И	М	Е
р	?	Р	И	М	?
И	И	И	И	?	И
М	М	М	?	М	?
Е	Е	?	И	?	Е

Тогда граф, соответствующий функции  $\|/$  в приведенной выше графической интерпретации, будет выглядеть так (вершины, соответствующие  $\forall(p)$ , обведены дважды):

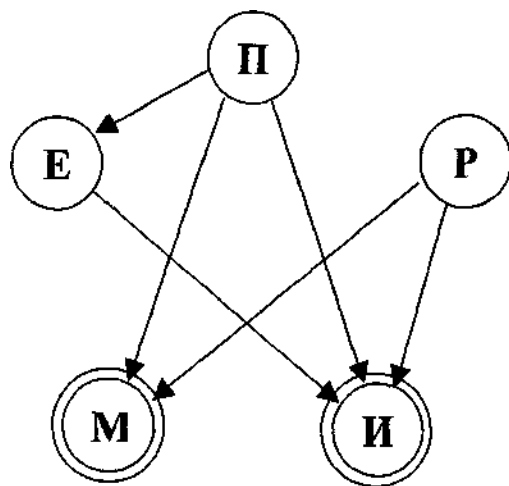


Рисунок 2.1. Графовая интерпретация функции  $\forall$ .

В данном примере  ${}^xP(p)=\{М,И\}$ .

**Утверждение 2.2.**  $U p \in E^M \Gamma(p) \in T(p)$

Доказательство.

Если  $p \in c_0$ , то  $\Gamma(p)=c_0 \Rightarrow$  по определению функции  $\|/ \forall c' \in C, c' \wedge c_0$  выполнено  $U c_0 C' C \{c_0, ?\} \Rightarrow$  по определению  $\forall c_0 \in T(p) \Rightarrow \Gamma(p) \in {}^xP(p)$ .

Если же  $pg\ c_0$ , то  $\exists c: \Gamma(p)=c$  и  $pc\ c \Rightarrow$  по определению функции  $\forall i$   
 $\forall c' \in C, \quad c \forall c$  выполнено  $i/\_{cc'} \in \{c, ?\} \Rightarrow$  по определению  $4^y$   $cc^\wedge(p)$   
 $\blacksquare \Rightarrow$

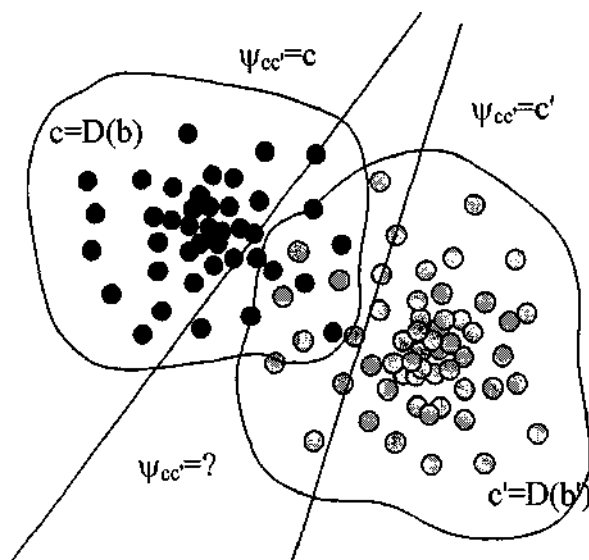
$\Gamma(p) \in^\wedge(p)$ .

Утверждение доказано.

Содержательно утверждение 2.2 означает, что множество "неопровергнутых гипотез локального распознавания"  $\Gamma(p)$  обязательно содержит "правильную гипотезу"  $\Gamma(p)$ .

На практике отделяющая функция  $\forall_{cc'}$  может быть введена с помощью построения поверхности (например, состоящей из фрагментов гиперплоскостей), разделяющей образы элементов разбиения, в нашем примере - двух фонем  $V\{h\}$  и  $0(B')$ .

Если такая поверхность на основе обучающей выборки может быть построена, данная пара образов фонем отделяется с помощью функции  $\forall_j$ ,



2.2. Построение отделяющей функции  $\psi$  для случая  $M=2$

являющейся характеристической для данной поверхности (если положить  $\forall_{cc}\{p\}=c$ , если точка  $p$  и множество  $0(b)$  лежат с одной и той же стороны относительно разделяющей поверхности, а иначе определить, что  $\forall_{cc}(p)=c'$ ).

противном случае, если Рисунок

так  $Y^{10}$

разделяющую поверхность построить не удастся

(например, когда  $\Pi(B) \cap O(B') \neq \emptyset$ ), зачастую можно построить две непересекающиеся разделяющие поверхности, которые разбивают все множест-



во  $TXB) \cup TXB')$  на три подмножества. Точки первого подмножества гарантированно принадлежат  $TXB)$ , второго подмножества - объединению  $D(b) \cup D(b')$  и третьего подмножества - множеству  $D(b')$ .

## **§2.2 Модель распознавания последовательного кода команды с помощью детерминированных автоматов.**

С помощью функции локального распознавания  $U^*$  невозможно однозначно восстановить код произнесенной команды, поскольку для каждого момента времени  $i$  функция  $W$  дает вместо однозначного ответа распознавания (кода  $\Gamma(y')$  точки  $y^1$  описания  $y$  сигнала в этот момент времени) множество гипотез распознавания, которое содержит в силу доказанного выше утверждения "наилучшую" гипотезу  $\Gamma(y')$ , но может содержать и другие буквы кодовой книги  $C$ .

Множество гипотез локального распознавания для каждого момента времени порождает некоторый язык слов в алфавите  $C$  возможных кодов сигнала, если рассматривать сигнал в целом.

Введем некоторые обозначения.

Пусть дано произвольное конечное множество  $A = \{a_1, a_2, \dots, a_r\}$ . Обозначим через  $R(A)$  регулярное выражение ([34]), описывающее регулярный язык однобуквенных слов в алфавите  $A$ :  $R(A) = (a_1 | a_2 | \dots | a_m)$ , где  $|$  означает символ «или».

Пусть  $a$  - регулярное выражение. Обозначим через  $(a)$  регулярный язык, которое оно описывает.

Пусть  $s$  - речевой сигнал,  $y = \|y^1 | y^2 | \dots | y^n\| \in (K^M)^*$  - его описание. Рассмотрим язык "гипотез локального распознавания" матрицы  $y$ :



$\{y \mid y = c_1 c_2 \dots c_n, c_i \in \Sigma^1\}$ . Этот язык, очевидно, описывается регулярным выражением  $f(y) = R(F(y^1))R(F(y^2)) \dots R(F(y^n))$ , причем  $\Gamma(y) \in \langle \Gamma(y) \rangle$ .

**Определение 2.4.** Пусть  $\Sigma^1$  - подалфавит кодовой книги. Окрестностью  $O(\Sigma^1)$  алфавита  $\Sigma$  назовем следующий подалфавит  $\Sigma$ :

$O(\Sigma^1) = \{c \in \Sigma : \exists c' \in \Sigma^1 : c \text{ и } c' \text{ не являются separable с помощью функции } T\}$ .

**Лемма 2.1** Для любой точки  $p \in K, \Sigma^1(p) \subset O(\Gamma(p))$ .

Доказательство.

Пусть  $\Gamma(p) = c \Rightarrow$  по утверждению 2.2  $c \in \Sigma^1(p)$ . Возьмем любую другую букву  $c' \in \Sigma^1(p)$ ,  $c \neq c'$ . Множества  $c'$  и  $c$  не являются separable с помощью функции  $T$ , так как противное противоречило бы определению separability. Следовательно,  $c' \in O(c) = O(\Gamma(p))$ .

Лемма доказана.

Объединив утверждение 2.2 и доказанную лемму, получим, что имеет место вложение:  $\Sigma^1(p) \subset O(\Gamma(p))$ .

л

Обозначим через  $c(b)$  множество элементов разбиения  $\Sigma$ , пересекающихся с образом  $D(b)$  фонемы  $b$ :

$$c(b) = \{c \in \Sigma : c \cap D(b) \neq \emptyset\}.$$

Через  $c(b)$  будем, соответственно, обозначать множество букв кодовой книги, соответствующих элементам разбиения из  $c(b)$ . Очевидно, что для любого звука  $b \in B$   $\Gamma(D(b)) \subset c(b)$ .

Покажем, что множество  $(f(y))$  гипотез локального распознавания описаний у произнесений любой речевой команды можно описать регулярным выражением, которое строится непосредственно по тексту этой команды.

**Определение 2.5.** Пусть  $p \in B$  - некоторая команда,  $p = b_1 b_2 \dots b_m$ .

Эталонным кодом  $\Gamma((3))$  команды  $|3|$  будем называть регулярное выражение

$$f(P) = R(O(cO^c(b_1)uc_0)) * R(O(c(b_1))) * R(O(c(b_1)uc(b_2)uc_0)) * R(O(c(b_2))) * \dots * R(O(c(b_m))) * R(O(c(b_m)uc(Juc_0)))$$

(здесь  $*$  - символ итерации).

<p><b>Теорема 2.1.</b></p>	<p>Пусть задана речевая модель <math>\langle B, B, S, (R^M)^*, C, n, T,  a, d, r, r\rangle</math>, отделяющая функция <math>u</math> и функция локального распознавания <math>\forall</math>. Тогда для любой звуковой команды <math>(3 = b_1 b_2 \dots b_m)</math> выполнено <math>(f(y))c:(f((3)))</math>, где <math>y = d(s)</math> - описание произвольного произнесения <math>sen(P)</math> команды <math>p</math>.</p>
----------------------------	--

Доказательство.

Поскольку  $\Gamma(y) \in (\Gamma(y))$

Пусть  $\Pi y U_i \dots I y C H a d$  и  $m(s) = M s, 1), |i(s, 2), \dots, f_i(s, 2m+2))$  - фонемная разметка сигнала  $s$ .

По свойству 2° отображения фонемной разметки  $\forall j = 1..m$

$\forall k = f_i(s, 2j) + 1, K s, 2j) + 1, \dots, |u(s, 2j+1) - 1 s_{[k \cdot T(s)/2; k+T(s)/2]} \cap n(b_j) \Rightarrow d(s, k) e$

$\wedge b_j) c c(b_j) \Rightarrow r(d(s, k)) e c(b_j)$  и по лемме 2.1

$\forall (d(s, k)) c O(r(d(s, k))) c O(c(b_j)) r >$

$\langle \Gamma(P(8, \wedge, 2i) + 1) I(8, \sqcup(8, 2i) + 2) | \dots I(8, \sqcup(8, 2i+1) - 1 ||) \rangle^{\wedge(0(c(B^\wedge)))^*}$ .

По свойству 3° отображения фонемной разметки  $\forall j = 1..m+1$

$\forall k = p(s, 2j-1), u(s, 2j-1) + 1, \dots |i(s, 2j)$ , выполнено, что либо



$S[k-[T(s)/2], k+[T(s)/2]]en(b_j.i)un(b_j)$ , либо не существует звука  $b \in B$ :

$S[k-[T(s)/2], k+[T(s)/2]]en(b) \wedge d(s, k)e^{\wedge b_j}uD(b_{j+1})u'D(b_0)cc(b_j)uc(b_j.i)uc_0 \Rightarrow$

$r(d(s, k))Gc(b_j)uc(b_j.i)uc(b_0)$  и по лемме 2.1

$\wedge(d(s, k))cO(r(d(s, k)))cO(c(b_j)uc(b_{j+1})uc(b_0)) \Rightarrow$

$\langle f(|d(s, Ks, 2j-1)|d(s, \wedge(s, 2j-1)+1)|\dots|d(s^{\wedge(s, 2j)}))|c$

$e \langle R(0(\{c^{\wedge}, c(6_{yi}), Co\}))^* \rangle.$

Следовательно,  $(f(d(s))) = (f(|d(s, 1)|d(s, 2)|\dots|d(s, n)|))> = =$

$\langle f(|d(s, rtM)|d^{\wedge} d(s, n < s, 3))|\dots|M < s, 2m+2)|))> = = \langle f$

$(|dCs, Ks, 1)|d < :s, H < s, 1)+1| . - |d(s, n(s, 2))|)f(|dCs, M < s, 2 > 4-1|.. |d(s, n(s, 3)-$

$1)|)\dots f(|d(s, Ks, 2m+1)|\dots|\wedge(s_5 2m+2)|))c$

$c(R(O(cO^{\wedge}c(b_1)uc_0))*R(O(c(b_1)))^{\text{!!!}}R(O(c(b_1)uc(b_2)uc_0))*$

$R(O(c(b_2)))^*\dots R(O(c(b_n)))^*R(0(c(b_n)ucO^{\wedge}c_0))^* \rangle = f((3),$  что и требовалось

доказать.

Теорема доказана.

Теорема 2.1 позволяет нам решить задачу локального распознавания речи в общем случае. Однако в связи с тем, что отделяющая функция  $\|/$  для большинства пар фонем устроена плохо - фонемы плохо отделимы друг от друга для используемых функций описания сигналов, - практическое применение этой теоремы ограничено. Дело в том, что для различных пар команд  $\{3$  и  $\{3'$  языки  $\Gamma(p)$  и  $\Gamma(p')$  часто пересекаются, что не позволяет использовать теорему 2.1. для решения задачи распознавания речи. Выход заключается в том, чтобы объединить некоторые фонемы в классы фонем и вместо разбиения пространства описаний сигналов на образы фонем использовать укрупнение этого разбиения - разбиение пространства  $BL^M$  на образы фонемных классов.

## §2.3 Модель и алгоритм распознавания кода команды для случая классов звуков. Нижняя оценка качества словаря команд.

Пусть теперь  $B = B^0 \cup B^1 \cup \dots \cup B^k$  - разбиение множества фонем на непересекающиеся классы,  $B^i \cap B^j = \emptyset$   $\forall i \neq j$ . Рассмотрим

разбиение

$R^M$  на подмножества  $(DB^i), i=1..k \cup P(B^0)$ , где  $c_i = Я(B^i) = \{ /Щ \}$  - Образы классов фонем, а через  $c_0 = O(B^0)$  обозначено межфонемное пространство  $R^M \setminus \bigcup_{i=1}^k DB^i$  (откуда  $DB^0 = DB_0$ ). Сопоставим каждому элементу

$c_i$  этого разбиения букву  $c_i$  кодовой книги  $C = \{c_i, i=0..k\}$ . В

этом случае справедливо следующее утверждение:

**Утверждение 2.3.** Если для всех  $C_j, C_j \in C$  функция  $*F$  отделяет  $q$  и  $c_j$  то

$$1) \forall p \in E^M \forall (p) = \{\Gamma(p)\}$$

$$2) \forall p \in S \forall y \in D(n(b))(f(y)) = \{r(y)\}$$

Доказательство.

1) По лемме 2.1 имеет место вложение  $\Gamma(p) \in \Gamma^x(p) \subset O(\Gamma(p))$ . Но по условию для кодовой книги  $C$   $O(c) = c$  для любой буквы  $c \in C$ . Значит,  $\forall (p) = \{\Gamma(p)\}$ .

$$2) f(y) = Rmy^1))Rmy^2)) \dots RWy^n)) = (B \text{ силу 1) } = (\Gamma(y^1))(\Gamma(y^2)) \dots (\Gamma(y^n)) \Rightarrow (\Gamma(y)) = \{\Gamma(y)\}$$

Утверждение доказано.

Содержательный смысл этого утверждения следующий.

Если разбиение множества фонем на фонемные классы произведено так, что удастся отделить друг от друга с помощью функции  $\Gamma^7$  образы этих классов фонем, то

- 1) Локальное распознавание каждой точки пространства описаний дает только один вариант распознавания - код этой точки
- 2) Код любого произнесения восстанавливается с помощью функции локального распознавания  $\psi^*$  однозначно

Замечание. Для случая, когда разбиение задано с помощью классов звуков, эталонный код  $\Gamma((3))$  любой команды ( $3 \in 3$ ,  $p = b_1 b_2 \dots b_m$  будет равен

$$f(P) = R(cQuc(b_1)uco) * R(c(b_1)) * R(c(b_1)uc(b_2)uco) * R(c(b_2)) * \dots * R(c(b_m)) * R(c(b_m)uc(J)uco) *$$

(здесь  $*$  - символ итерации).

Предложим теперь механизм разбиения множества фонем на фонемные классы, позволяющего применять утверждение 2.3.

Пусть заданы речевая модель, функция отделимости  $\|/\|$  и функция локального распознавания  $\psi$ . Введем на буквах кодовой книги  $S$  бинарное *отношение отделимости*, обобщив определение отделимости, приведенное ранее. Определение отделимости букв  $s$  и  $e'$  (обозначение  $s \sim e'$ ) будем вводить индуктивно.

### Определение 2.6.

1. Если  $s$  и  $e'$  отделимы с помощью  $\|/\|$  в смысле определения 2.5, то  $s \sim e'$ .
2. Если  $s \sim e'$ ,  $e' \sim e''$ , то  $s \sim e''$ .

Очевидно, что отношение отделимости  $\sim$  является отношением эквивалентности. Классы эквивалентности, на которые это отношение разбивает элементы кодовой книги из примера 2.1, являются примером удачного разбиения множества фонем на фонемные классы, позволяющего использовать образы этих классов для построения алгоритма распознавания речи.

Пусть словарь команд  $V$  речевой модели  $\langle V, V, S, (R^M), C, П, T, ц, c1, Г, г \rangle$  имеет вид  $V = \{(3^b p_2 v P_N)\}$  - Пусть заданы функции отделимости  $\backslash\backslash$  и локального распознавания  $Ч^*$ . Пусть кодовая книга  $C$  и соответствующее ей разбиение  $C$  заданы на основе образов классов эквивалентности введенного отношения отделимости  $\sim$ , а кодирующее отображение  $\Gamma$  определено с помощью  $\mathbb{Y}$  на основе утверждения 2.3 (пункт 1). Определим алгоритм распознавания речи, сопоставив каждой команде  $P_i$  из  $V$  эталон  $\epsilon_i = (\Gamma(P_i))$  и введя на декартовом произведении  $C^* \times E$ ,  $E = \{e_1, e_2, \dots, e_k\}$  функцию расстояния  $P(\cdot, Y >^e) - \begin{cases} 0, & \text{если } y \in e \\ 1, & \text{иначе} \end{cases}$

Справедливо

**Утверждение 2.4.** Для алгоритма

распознавания речи  $(E, p)$  для любого кодового слова  $y \in \Gamma(c_i(\Pi(V)))$  имеет место вложение  $\gamma(y) \subset \Gamma_{(E, p)}(y)$

Доказательство.

Пусть  $p \in \gamma(y)$ . Тогда по определению  $\gamma$   $y \in \Gamma(\check{e}(\Pi((3))))$ . По утверждению 2.3  $y \in (f(y)) \subset (f(p)) = e \Rightarrow$  по определению  $p$   $p(y, e) = 0 \Rightarrow P \in \underset{\text{ДеЯ}}{\text{Argmm}}(p(\cdot, e_i)) = \Gamma_{(E, p)}(y)$ . Итак, мы доказали, что  $\gamma(y) \subset \Gamma_{(E, p)}(y)$ .

Утверждение доказано.

Содержательный смысл доказанного утверждения состоит в том, что построенный на основе разбиения звуков на классы детерминирован-но-автоматный алгоритм распознавания речи каждый раз включает в результат распознавания все "правильные ответы".

**Утверждение 2.5.** Пусть задан введенный выше алгоритм распознавания речи  $(E, p)$ . Если словарь команд  $V$  нашей речевой модели та-

ков, что для любой пары команд  $p \in V$  соответствующие им эталоны  $e_j$  и  $e_j$  не пересекаются, то  $\Gamma = \Gamma(E, p)$ .

### Доказательство.

По утверждению 2.5 для любого  $y \in \Gamma(E, p)$ . Но если эталоны всех команд словаря попарно не пересекаются, то для любого  $y$  множество  $\Gamma(E, p)(y) = \{e \in E \mid \text{ШП}(e, y) > 0\}$

по построению функции расстояния состоит

*MB*

ровно из одного слова. Поскольку  $\Gamma(E, p)$  не пусто (оно содержит  $P: y \in \Gamma(E, p)$ ), получаем, что  $\Gamma(E, p) = \{P\}$ . Утверждение доказано.

Утверждение 2.5 имеет следующее практическое значение. Если на этапе построения системы распознавания речи мы имеем возможность выбирать команды произвольно, то можно подобрать команды, эталоны которых не пересекаются. В этом случае алгоритм распознавания будет работать идеально.

Возникает вопрос, как, не проводя экспериментов, сравнить среднее число ошибок алгоритма  $(E, p)$  для двух различных словарей команд, каждый из которых содержит пары пересекающихся эталонов. С этой целью введем понятие качества словаря команд.

Качество словаря команд будем определять индуктивно. Положим, что качество словаря, состоящее из одной команды, равно 1. Для словарей из двух команд определим качество как:

**Определение 2.7.** Качественным словарем команд  $V = \{p, P'\}$  для введенного алгоритма распознавания речи  $(E, p)$  назовем величину

гда<sup>^</sup>ипот)

Содержательно качество  $q(E,p)(\{P>P'\})$  словаря  $\{(3,p')\}$  связано со средним числом ошибок при распознавании алгоритмом  $(E,p)$ .

Пусть  $S$  - некоторое конечное множество,  $f: S \times S \rightarrow R$ . Обозначим через  $\odot_s(f)$  среднее значение функции/на парах различных элементов из  $S \times S$ :

$$\odot_s(f) = \frac{\sum_{a,a' \in S, a \neq a'} f(a,a')}{S(S-1)}$$

## 2 Определение 2.8. Качеством

словаря команд  $B$ , состоящего из более чем двух команд, назовем величину  $q(E,p)(S) = \odot_s(q(E,p))$ .

**Лемма 2.2.** Для любого словаря команд  $B$   $0 < q(E,p)(S) < 1$ .

Доказательство.

Для каждого  $y \in \Gamma(E,p)(y)$  СОСТОИТ либо из одной команды (и в этом случае  $\Gamma(E,p)(y) = \{y\}$ ), либо из двух команд. Следовательно, по определению 2.7 качество  $q(E,p)(\{P,P'\})$  любого словаря из двух команд  $\{(3,(3'))\}$  лежит на отрезке  $[0,1]$ . При усреднении по всем парам команд из словаря  $B$  неравенство  $0 < q(E,p)(S) < 1$  сохранится.

Лемма доказана.

Определим расстояние  $p(e,e')$  между множествами  $e$  и  $e'$  как

$$p(e,e') = \begin{cases} \Gamma O, & \text{если } e \setminus e' \neq \emptyset \\ 1, & \text{иначе} \end{cases}$$

Справедлива теорема о нижней оценке качества словаря команд для детерминированно-автоматного алгоритма  $(E, p)$ :

**Теорема 2.2.** Для введенного алгоритма распознавания речи  $(E, p)$  для любого словаря команд  $\mathcal{B}$  справедлива оценка:

$$q_{(E,p)}(S) \geq 0_6(\mathcal{B}).$$

Доказательство.

Случай словаря из одной команды для задачи распознавания речи лишен смысла.

Если словарь состоит из двух команд, то по определению

$$n \quad (R / ?'П - ? \quad \frac{r \Gamma (<*(\Pi(i) \Pi(\mathcal{B}))}{\Gamma_{\text{да}}(y) \Gamma \Pi(\mathcal{B})) \Gamma} \\ ЯшШР*)-*$$

Если эталоны команд  $p$  и  $p'$  пересекаются, то  $0(p) = 0$ , откуда, используя лемму 2.2, получаем требуемую оценку для этого случая.

Если же эталоны команд  $P$  и  $p^1$  не пересекаются, то  $0(P) = 1$ .

Однако в этом случае для каждого  $y \in \Gamma_{(E,p)}(y)$  состоит ровно из одной команды, откуда и  $q_{(E,p)}(S) = 1$ .

Усреднение по всем парам команд из словаря не изменит этой оценки, откуда следует справедливость теоремы для словарей с произвольным числом команд.

Теорема доказана.

## ГЛАВА 3. ВЕРОЯТНОСТНЫЕ АВТОМАТНЫЕ МОДЕЛИ.

### §3.1. Понятие монотонного автономного вероятностного автомата.

#### Модель распознавания речи с помощью вероятностных автоматов.

Наиболее распространенным математическим объектом, использующимся на практике при построении алгоритмов распознавания речи, является вероятностный автомат. Обусловлено это тем, что применение конечных детерминированных автоматов, как это показано в предыдущей главе, позволяет решать задачу распознавания речи только для специально подобранных словарей команд. При создании эффективных алгоритмов распознавания речи в общем случае используются вероятностные автоматы.

Вероятностный автомат хорош тем, что он может применяться в качестве "эквивалента" множества кодов описаний всех речевых сигналов, являющихся произнесениями какого-либо слова словаря или звука речи. Вероятностный автомат в нашем случае используется в качестве генератора этих кодов.

Обучение вероятностных автоматов происходит по правилам, так что коды описаний произнесений «родных» команд порождаются ими с большой вероятностью, и эта вероятность выше вероятности порождения этих кодов «чужими» автоматами. Задача обучения параметров вероятностных автоматов находится за рамками настоящего исследования.

Пусть  $(B, S, (R^M), C, P, T, |i, (i, \Gamma, \gamma))$  - речевая модель, задан словарь команд  $S = \{p_i, \{3\gamma, \dots, P_N\}$  и решается задача распознавания кодов описаний произнесений этих команд, т.е. кодовых слов  $y \in T(d(nCB))$ .

**Определение 3.1.** *Вероятностным автоматом* называют четверку

$(B, C, Q, v)$ , где

$B$  - конечный алфавит входных символов;



$S$  - конечный алфавит выходных символов;  $Q$  - алфавит (в нашем случае - конечный) состояний автомата,  $v$  - функция, определенная на декартовом произведении  $B \times Q$  и принимающая в качестве своих значений вероятностные меры на декартовом произведении  $B \times Q$  ([26, 34]).

В содержательной интерпретации вероятностного автомата как преобразователя информации числа  $v(b, q)(c, q') \in [0, 1]$  означают условную вероятность перехода автомата в состояние  $q'$  и выдачи символа  $c$ , при условии, что автомат находился в состоянии  $q$  и получил на вход символ  $b$ . Будем обозначать поэтому  $v(b, q)(c, q') := v(c, q'/b, q)$ . Числа  $v(c, q'/b, q)$  в силу этого обладают свойствами  $0 < v(c, q'/b, q) < 1$  и  $\sum v(c, q'/b, q) = 1 \forall b, q$ .

Вероятностный автомат функционирует подобно детерминированному автомату. Находясь с некоторой вероятностью в некотором начальном состоянии  $q$ , автомат работает по тактам и в каждый такт времени получает на вход букву  $b$  алфавита  $B$ , выдает на выход букву  $c$  и переходит в состояние  $q'$  с вероятностью  $v(c, q'/b, q)$ , и т.д.

*Инициальным вероятностным автоматом* называется пятерка  $(B, S, Q, v, v_0)$ , где  $v_0$  - вероятностная мера на множестве состояний  $Q$ , задающая распределение для начального состояния автомата.

Нам будет интересно рассмотреть частный случай вероятностных автоматов - вероятностный автомат Мура без входа. Именно такие автоматы, обладающие дополнительно свойством монотонности (о нем будет рассказано ниже), используются при построении алгоритмов распознавания речи.

**Определение 3.2.** *Вероятностный автомат Мура без входа* (далее - *автономный вероятностный автомат*) - это тройка  $(S, Q, v)$ , где  $S$  - алфа-

вит выходных символов;  $Q$  - конечный алфавит состояний автомата,  $v$  - функция, определенная на множестве состояний автомата  $Q$  и принимающая в качестве своих значений вероятностные меры на множестве  $C \times Q$ , такая, что она разлагается в произведение  $v = \pi \circ P$ , где  $\pi$  действует на множестве  $Q$  и имеет значениями вероятностные меры на множестве  $Q$ , а  $P$  действует на том же множестве  $Q$  и имеет значениями вероятностные меры на множестве  $C$  ([26]).

Содержательно  $v(c, q \rightarrow q')$  интерпретируется как условная вероятность перехода в состояние  $q'$  и выдачи символа  $c$  при условии нахождения в предыдущий момент времени в состоянии  $q$ . Условие представимости вероятностной меры  $v$  в виде произведения вероятностных мер  $\pi$  и  $P$  означает содержательно, что два события - переход в следующее состояние и выдача некоторого символа на выход автомата - несовместны.

Прямоугольную матрицу будем называть *стохастической*, если все ее элементы неотрицательны и сумма элементов которой в каждой строке равна 1.

Как в общем случае, так и для автономных автоматов, функция  $v$  полностью определяется некоторой конечной системой матриц.

Так, для автономного автомата функция  $v$  задается стохастической матрицей  $\pi$  размерности  $m \times m$  (где  $|Q|=m$ ) вероятностей переходов и стохастической матрицей  $P$  размерности  $m \times k$  (где  $|C|=k$ ),  $i$ -я строка которой задает распределение вероятностей выходных символов для состояния  $q_i$ . Матрицы  $\pi$  и  $P$  являются стохастическими.

Определение автономного вероятностного автомата можно обобщить на случай, когда матрицы  $\pi$  и  $P$  не являются стохастическими, но обладают

тем свойством, что сумма элементов в каждой их строке не превосходит 1 (будем такие матрицы называть *слабо-стохастическими*). Добавив дополнительное (поглощающее) состояние и дополнительную (пустую) букву, можно доопределить данный автомат естественным образом до автономного вероятностного автомата со стохастическими матрицами переходов  $\kappa'$  и выходов  $P'$  следующим образом:

$$\kappa' = \begin{array}{|c|c|} \hline \tau\Gamma & \lambda \\ \hline 0 & 1 \\ \hline \end{array}, \quad P' = \begin{array}{|c|c|} \hline P & P \\ \hline 0 & 1 \\ \hline \end{array}$$

(вектор-столбец  $\kappa$  в матрице  $\tau\kappa'$  состоит из чисел, дополняющих каждую из строчек матрицы  $\tau\kappa$  до 1, а вектор-столбец  $P$ , соответственно, дополняет до 1 распределение выходных букв в каждом состоянии, т.е. строчки матрицы  $P$ )

Благодаря возможности такого обобщения многие утверждения, справедливые для автономных вероятностных автоматов, верны и для автоматов со слабо-стохастическими матрицами. Если это не оговорено особо, мы будем и такие автоматы называть вероятностными.

Будем считать, что выделено некоторое подмножество состояний автомата, которые будем называть финальными. Множество финальных состояний автомата однозначно задается вектором-столбцом  $j_{\text{ф}}$ , в котором в  $i$ -й позиции стоит единица тогда и только тогда, когда состояние  $q_i$  - финальное, а иначе там стоит 0.

Будем считать, что вероятностный автомат, в котором заданы финальные состояния, функционирует таким образом, что при попадании в финальное состояние он прекращает свою работу, или, что эквивалентно, в финальном состоянии автомата (вне зависимости от наличия у автомата входа и того, какой символ поступил на вход, если у автомата вход есть) вероятность перехода в себя и в любое другое состояние равна 0. Вероятностные автоматы, в которых выделены финальные состояния, будем обозначать, включая вектор финальных состояний в число объектов, определяющих автомат, на-

пример,  $(B, C, Q, v, v_0, v_F)$ . Заметим, что матрица вероятностей переходов  $t_c$  в таком случае уже не будет стохастической, т.к. финальным состояниям в ней будут соответствовать нулевые строки.

Для заданного автономного инициального автомата  $A = (C, Q, t_c, P, V_0, V_F)$  можно рассматривать вероятность  $p(y)$  того, что автомат перешел в финальное состояние, выдав при этом на выход слово  $y \in C^*$ . Известно ([26]), что  $p(y)$  можно посчитать с помощью матричной формулы

$$p(y) = \pi_0 M(y) \pi^T \text{ где } M(y) = M(c_1 c_2 \dots c_n) = M(c_1) M(c_2) \dots M(c_n) \quad (1)$$

(здесь через  $M(c/)$  обозначена матрица вероятностей переходов автомата и одновременной выдачи символа  $C/$ ,  $M(c_i)y = t_{iy} P_{iy}$ ).

Более подробно эта формула выглядит так:

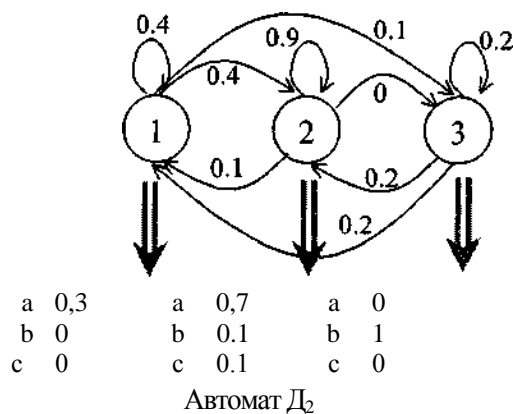
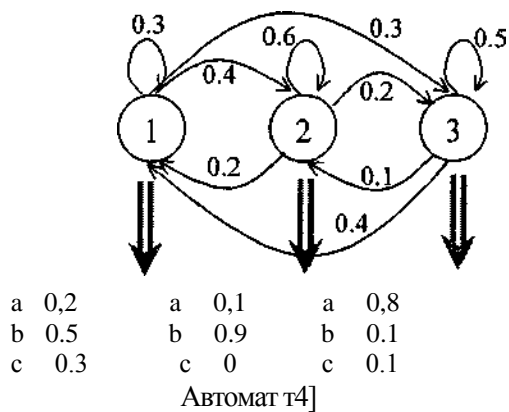
$$P(y) = \pi_0 P(\Gamma \setminus \Gamma_0 = \wedge (c_1 O^{\wedge} K A^{\wedge} K, \dots, W^{\wedge},$$

В дальнейшем мы будем рассматривать только автономные вероятностные автоматы, поэтому слово автономный для простоты будем иногда опускать.

### Пример 3.1.

Пусть вероятностные автоматы  $A_i = (\{a, b, c\}, \{1, 2, 3\}, t_{c_i}, P_i)$  и

$A_2 = (\{a, b, c\}, \{1, 2, 3\}, t_{c_2}, P_2)$  представлены диаграммами:



У автомата  $A_4$  матрицы переходов и выходов являются стохастическими:

$$P_1 = \begin{pmatrix} 0,2 & 0,5 & 0,3 \\ 0,8 & 0,9 & 0 \\ 0,4 & 0,2 & 0,5 \end{pmatrix}, P_i = \begin{pmatrix} 0,3 & 0,3 \\ 0,2 & 0,2 \\ 0,4 & 0,5 \end{pmatrix}$$

У автомата  $A_j$ , напротив, ни матрица переходов, ни матрица выходов стохастическими не являются:

$$P_2 = \begin{pmatrix} 0,3 & 0 & 0 \\ 0,7 & 0,9 & 0 \\ 0 & 1 & 0 \end{pmatrix}, P_2 = \begin{pmatrix} 0,3 & 0,4 & 0,1 \\ 0,2 & 0,9 & 0 \\ 0,2 & 0,2 & 0,2 \end{pmatrix}$$

Матрицы  $M_j(a)$ ,  $M_j(b)$ ,  $M_j(c)$ ,  $i=1,2$ , участвующие в формуле (1), для автоматов  $A_4$  и  $A_i$  имеют следующий вид:

$$M_i(a) = \begin{pmatrix} 0,06 & 0,08 & 0,06 \\ 0,02 & 0,06 & 0,02 \\ 0,32 & & 0,08 \\ 0,4 & & \end{pmatrix}, M_i(b) = \begin{pmatrix} 0,15 & 0,2 & 0,15 \\ 0,18 & 0,54 & 0,18 \\ 0,04 & & 0,01 \\ 0,05 & & \end{pmatrix}, M_i(c) = \begin{pmatrix} 0,09 & 0,12 \\ 0,09 & \\ 0 & 0 \\ 0 & 0,04 & 0,01 \end{pmatrix}$$

0,05

$$M_2(a)=\begin{pmatrix} 0,09 & 0,12 \\ 0,03 & 0,07 \\ 0,63 & 0\ 0 \\ 0 & 0 \end{pmatrix}, M_2(b)=\begin{pmatrix} 0 & 0 \\ 0,01 & 0,09 \\ 0,2 & 0,2 \end{pmatrix}, M_2(c)=\begin{pmatrix} 0 & 0 \\ 0 & 0 \\ 0,01 & 0,09 \\ 0 & 0 \\ 0 & 0 \end{pmatrix}$$

Каждому автономному вероятностному автомату естественным образом ставится в соответствие *словарная функция*  $p: C^* \rightarrow [0,1]$ , сопоставляющая каждому слову  $u$  из  $C^*$  вероятность  $p(u)$ . Эту функцию можно рассматривать как счетное множество пар вида  $(u, p(u))$ , где пары следуют согласно лексикографическому порядку, заданному на словах  $u$ .

**Определение 3.3.** Два инициальных вероятностных автомата называются *эквивалентными*, если тождественно равны их словарные функции.

**Определение 3.4.** Для автономного инициального автомата

$T_4 = \{C, Q, T, P, V_0, V_F\}$  квадратную матрицу  $A$ , где  $A_{ij} = \sum_{a \in C} p_{ij}(a)$ , будем называть *приведенной матрицей переходов* (здесь  $k = |C|$  - мощность выходного алфавита).

( $i, j$ )-й элемент приведенной матрицы переходов имеет смысл вероятности того, что автомат перешел из  $i$ -го состояния в  $j$ -ое и подал какую-то букву на выход.

**Лемма 3.1.** Для любого автономного вероятностного автомата  $T_4 = \{C, Q, T, P, V_0, V_F\}$  выполнено равенство  $\sum_{r \in C^*: |r|=n} p(r) = V_0 A^n V_F$  где  $A$  - приведенная матрица переходов автомата  $A$ .

Доказательство.

$$2>(r) = 2 \quad 2 v_0(f_{fi})^* > (c_{\pi} k_{i(\pi)} (c_{\pi} x_2, \dots \cdot p_{la} \{c_{jn} k_{\pi}$$

$$y \in C : |\gamma| = n \quad r = c_{h_1} c_{h_2} \dots c_{j_m} Y^{\wedge 2} Y^{\wedge} Y_{m+z}$$

$$L \rightarrow i \quad I \rightarrow i \quad UVJ, / \quad /, V \quad y, \bullet \quad /i/_2 \quad '2^V \quad \wedge_2^Y \quad '2'3 \quad 'x^4 \quad y_{\pi}{}^7 \quad \forall n$$

$$Z^v_o(Q_i)Y_{\alpha}{}^{\pi}{}_{\beta}{}^6 - 'Y_i I \quad Y^{\wedge}(cI)P_t(c, ) \dots - \mathcal{L}(c, )$$

$$r^{\circ} J \tilde{r}^{\circ} J 2^{-c} J,$$

$$Z \quad vofe, )^{3/4 3/4} * - \blacksquare^{\wedge}_{\pi+1} Z^{\wedge}, (<v, ) Z^{\wedge}_2 ({}^c_{\pi}) - IX ({}^c_{\pi})$$

$${}^v_F (< //_{\hat{a}+} ) H$$

$$4 = ?, y/2 - < 7, j_1, \ddot{r}$$

$$M^{?/_{\ast, +1}} = I$$

Лемма доказана.

**Следствие 3.1.** Для любого автономного вероятностного автомата  $A = \{C, Q, 7C, P, V_o, v_F\}$  со стохастической матрицей  $P$  выполнено  $\hat{p}(y) = v_B n^n v_F$ .  
 $y \in C^*: \forall n$

Доказательство

Для этого автомата матрицы  $n$  и  $n$  совпадают, откуда и следует иско-  
мое равенство.

Лемма доказана.

**Лемма 3.2.** Если  $R$  - квадратная матрица (такая что матрица  $E - R$  - невырожденная) и матричный ряд  $\prod_{l=1}^{+\infty} \Gamma^d$  сходится (т.е. сходящимися последо-  
вательностями являются все соответствующие элементы  $\prod_{l=1}^N$  матриц  ${}^{\wedge}R''$ ,

$N = 1, 2, \dots$ ), то сумма этого ряда равна  $(E - R)^{''l} - E$ .

Доказательство.

$$N \quad N \quad N+1 \quad XT_{L1}$$

$$l=0 \quad l=0$$

$$N \quad J \quad \Pi'' = (E - R^{N+l})(E - RY^l)^{l=1} \quad E$$

$$=> \mathcal{L}V \quad iN+b$$

$$= (E - R^{N+l})(E - Ry^l - E$$

$$\ll=0$$



+00

тождество. Лемма доказана.

$0 \ 1 \cdot$	$i^{\bullet\bullet} t$	$0 \ 0$	$\bullet\bullet \ 0$
$0 \ 0$	$\cdot \ 1$	$\cdot$	$0$
$0 \ 0$	$\blacksquare^{\bullet}$	$1 \ 0$	$\bullet\bullet \ 0$
$0 \ 0$	$\blacksquare^{\bullet}$	$0 \ 1$	$\bullet\bullet \ 0$
$\cdot$	$0$	$\cdot$	$0$
$0 \ 0$	$\cdot \ 0 \ 0$	$\bullet\bullet \ 1 \ 0$	
$\cdot$	$0$	$\cdot$	
$0 \ 0$	$\cdot \ 0 \ 0$	$\bullet\bullet \ 1$	
$\cdot$	$0$	$0$	
$0 \ 0$	$\cdot \ 0 \ 0$	$\bullet\bullet \ 0$	
$\cdot$	$0$	$\cdot$	
$0 \ 0$	$\cdot \ 0 \ 0$	$\bullet\bullet \ 0$	

Рисунок 3.1. t-й косой ряд.

Заметим, что  $I_0$  - единичная матрица, а  $I_t$  при  $t > m-1$  - нулевая матрица. Кроме того, 1-й косой ряд обладает свойством:  $\Pi^t = I_t$ .

**Определение 3.6.** Жордановой клеткой  $J(X)$  размерности  $m$  называется квадратная  $m \times m$  матрица  $J(A, \epsilon) = A\epsilon + I_i$ .

80

**Лемма 3.3.** Для  $X: 0 < \kappa < 1$  матричный ряд  $\sum_{n=0}^{\infty} J(X)^n$  сходится.

Доказательство.

В следующих выкладках равенство матричных рядов понимается в том смысле, что они либо оба сходятся к одной и той же матрице, или оба расходятся.

$$\sum_{n=0}^{\infty} J(X)^n = \sum_{n=0}^{\infty} (K + \gamma^n) = \sum_{n=0}^{\infty} \gamma^n \sum_{n=0}^{\infty} J(X)^n = \sum_{n=0}^{\infty} \gamma^n \sum_{n=0}^{\infty} J(X)^n$$

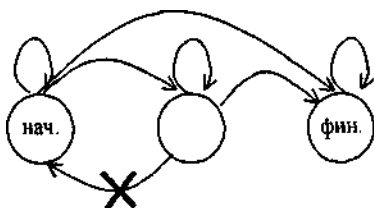
Таким образом, все элементы на  $k$ -й диагонали суммы данного матричного ряда равны  $\sum_{n=0}^{\infty} \gamma^n (i + n - 1) \dots (i + \gamma - 1)$ . Ряд, стоящий справа, при  $0 < \gamma < 1$  является сходящимся, поэтому сходится и данный в условии матричный ряд.

Лемма доказана.

Перейдем теперь к рассмотрению частного случая автономных вероятностных автоматов. А именно:

**Определение 3.7.** Инициальный автономный вероятностный автомат  $\mathcal{A} = (C, Q, \gamma, P, V_0, v_F)$  будем называть *монотонным*, если выполнены условия:

- 1)  $\gamma^i_j = 0$  при  $i > j$
- 2)  $v_0 = (1, 0, 0, \dots, 0)$
- 3)  $v_F^T = (0, 0, \dots, 0, 1)$



**Рисунок 3.2.** Монотонный автомат.



Пусть  $X = (E - ft)^x$ . Тогда из верхнетреугольности матрицы  $ft$  следует, что и матрица  $X$  имеет также верхнетреугольный вид. Причем из условия  $(E - ft)X = E$  следует равенство  $(E - ft)x = (0, 0, \dots, 0, 1)^T$ , где  $x$  - последний столбец матрицы  $X$ , откуда:

$$1) (1 - \sum_{i=1}^n ft_{ii}) X_{nn} = 1 \Rightarrow X_{nn} = \frac{1}{1 - \sum_{i=1}^n ft_{ii}}, \text{ в частности, } X_{nn} = 1,$$

$$2) (1 - \sum_{i=1}^m ft_{mi}) X_{mi} + (-\sum_{i=1}^m ft_{mi}) = 0, \text{ откуда } X_{mi} = \frac{\sum_{i=1}^m ft_{mi}}{1 - \sum_{i=1}^m ft_{mi}}. \text{ Из условия}$$

$\sum_{i=1}^m ft_{mi} < 1$  СЛЕДУЕТ, ЧТО  $ft_{mi} < 1 - \sum_{i=1}^m ft_{mi}$  СЛЕДОВАТЕЛЬНО,  $X_{mi} < 1$ .

3)  $(1 - \sum_{i=1}^n ft_{ni}) X_{ni} + \sum_{i=1}^m ft_{ni} X_{mi} + (-\sum_{i=1}^n ft_{ni}) = 0$ , откуда  $(1 - \sum_{i=1}^n ft_{ni}) X_{ni} = \sum_{i=1}^m ft_{ni} X_{mi} + \sum_{i=1}^n ft_{ni}$ , откуда  $X_{ni} = \frac{\sum_{i=1}^m ft_{ni} X_{mi} + \sum_{i=1}^n ft_{ni}}{1 - \sum_{i=1}^n ft_{ni}} < \frac{\sum_{i=1}^m ft_{ni} + \sum_{i=1}^n ft_{ni}}{1 - \sum_{i=1}^n ft_{ni}} < 1$ , что и требовалось доказать.

Теорема доказана.

**Следствие 3.2.** Для любого монотонного автомата  $A = (C, Q, \Gamma, P, v_0, v_F)$ , такого, что  $\forall i \in \Gamma, P(i) = ! \bullet$ , выполнено

Доказательство.

Доказательство для этого частного случая автоматов отличается от доказательства доказанной выше теоремы тем, что в этом случае  $ft$  -  $n$  - стохастическая матрица по всем строчкам, кроме последней.

Поэтому если  $X = (E - ft)^l$ , то из условия  $(E - ft)x = (0, 0, \dots, 0, 1)^T$ , где  $x$  - последний столбец матрицы  $X$ , получаем:

$$1) (1 - \sum_{i=1}^n ft_{ii}) X_{nn} = 1 \Rightarrow X_{nn} = \frac{1}{1 - \sum_{i=1}^n ft_{ii}}, \text{ в частности, } X_{nn} = 1,$$

2)  $(1 - \sum_{i=1}^m ft_{mi}) X_{mi} + (-\sum_{i=1}^m ft_{mi}) = 0$ . В нашем случае  $\sum_{i=1}^m ft_{mi} + \sum_{i=1}^n ft_{mi} = 1$ , откуда следует, что  $\sum_{i=1}^m ft_{mi} = 1 - \sum_{i=1}^n ft_{mi}$  следовательно,  $X_{mi} = 1$ .

$$m) (1 - \gamma_{i1})X_{1T} + (-\gamma_{i2} X_{2T}) + (-\gamma_{i3} X_{3T}) + \dots + (-\gamma_{im} X_{mT}) = 0,$$
откуда  $(1 - \gamma_{i1})X_{i1} = \gamma_{i2} X_{2T} + \gamma_{i3} X_{3T} + \dots + \gamma_{im} X_{mT} = \gamma_{i2} + \gamma_{i3} + \dots + \gamma_{im} = 1 - \gamma_{i1}$ , откуда  $X_{i1} = 1$ , что и требовалось доказать.

Следствие доказано.

Доказанное следствие имеет следующий смысл: события "монотонный автомат завершил работу и выдал слово  $y$ " для разных  $y$  несовместны.

Из доказательства теоремы 3.1 следует, что для монотонного автомата  $T_4 = (C, Q, \gamma_C, P, LIO, LIF)$  сумма матричного ряда  $\sum_{k=1}^{\infty} \gamma^k$  равна матрице  $(E - \gamma)^{-1}$ , которую будем называть *предельной приведенной матрицей переходов*.

$(y)_{ij}$ -й элемент предельной приведенной матрицы переходов - это вероятность того, что автомат перешел из  $i$ -го состояния в  $j$ -ое за конечное число шагов.

Введем теперь вероятностно-автоматный *алгоритм распознавания речи*. Пусть одним из известных методов ([18-22]) на основе набора примеров-кодовых слов каждой команде  $\{C\}$  словаря команд  $B$  сопоставлен вероятностный автомат  $A$ . Назовем эталоном команды  $P$ ; стохастическую словарную функцию  $\epsilon_i = p_{Ai}$  - Расстояние между кодовыми словами и эталонами команд определим как  $p(y, \epsilon_i) = 1 - p_{Ai}(y)$ .

Алгоритм распознавания речи  $(E, p)$  называется традиционно методом скрытых марковских процессов. Результат  $\Gamma_{(E, p)}(y)$  распознавания кодового слова  $y$  этим методом - это команда  $p$ , максимизирующая вероятность  $P_{Ai}(y)$  наблюдения слова  $y$  на выходе вероятностных автоматов  $D$ .

### §3.2 Метрика $\rho_i$ на множестве стохастических словарных функций.

Словарные функции, удовлетворяющие условию  $\sum_y f(y) = 1$ , будем называть *стохастическими*. Множество стохастических словарных функций будем обозначать через  $S$ . Если  $\sum_y f(y) < 1$ , словарные функции будем называть *слабостохастическими*. Множество слабостохастических словарных функций будем обозначать через  $S'$ .

Если на множестве стохастических (слабостохастических) словарных функций введена метрика, то она порождает метрику на вероятностных автоматах, где *расстоянием между монотонными автоматами* будем называть расстояние между соответствующими им стохастическими словарными функциями.

Каждую стохастическую словарную функцию  $\{(y, p(y))\}$  можно представить в виде вектора  $p = (p(y_1), p(y_2), \dots, p(y_0, \dots)) \in D^e$  слова  $y_1, y_2, \dots, y_0, \dots$  занумерованы в лексикографическом порядке. Будем обозначать множество таких стохастических векторов также символом  $S$  (или  $S'$  для слабостохастических словарных функций).

Метрику на множестве стохастических (слабостохастических) словарных функций можно ввести различными способами.

Поскольку  $\sum_y f(y) = 1$ , то  $p \in S$  где  $i = \{(X_1, X_2, \dots, X_n, \dots)\}$ :  
 $\sum_{y \in C^*} |f(x_i) - g(x_i)| < \infty$   $n = 1, 2, \dots$

([29]). Введем на множестве  $S$  метрику из пространства  $l^1$ .

**Определение 3.8.** Расстоянием  $\rho_i$  между стохастическими (слабо-стохастическими) словарными функциями  $P_1 = \{(y, P_1(y)), y \in C\}$  и  $P_2 = \{(y, P_2(y)), y \in C\}$  назовем функцию

$$\rho_1(\rho_2) = \inf_{y \in \rho_2} \rho(y, \rho_1)$$

**Утверждение 3.1.**  $\rho_1$  - метрика в  $\mathcal{P}_s$  и в  $\mathcal{P}'_s$ .

Доказательство.

Проверим корректность определения. Ряд сходится, т.к.  $\rho_1(\rho_n, \rho)$

$y$

$y$

$\rho$

Выполнимость свойств данной метрики из пространства  $\mathcal{P}$  проверяется просто:

1) Симметричность очевидна

2)  $\rho_1(\rho, \rho) = 0$  - очевидно

3) Неравенство треугольника:  $\rho_1(\rho(y), \rho_z(y)) \leq$

$$\rho_1(\rho(y), \rho_z(y)) \leq \rho_1(\rho(y), \rho(y)) + \rho_1(\rho(y), \rho_z(y)) =$$

$$\rho(y, \rho(y)) + \rho(y, \rho_z(y)) = \rho(y, \rho_z(y))$$

Утверждение доказано.

### §3.3 Метрика $p_2$ на множестве стохастических словарных функций. Функция качества словаря команд.

Более естественной метрикой на множестве монотонных автоматов является метрика  $p_2$ . Она имеет прямую связь с тем, как человек воспринимает различные слова и звуки и оценивает их похожесть.

**Определение 3.9.**  $p_2(P, Q) = 1 - 2^{-\min_{r \in C'} (A(r) \wedge_2(r))}$

**Утверждение 3.2.**  $p_2$  является метрикой на множестве  $/_s$ .

Доказательство.

Проверим корректность определения. Ряд сходится, т.к.

$y \in C'$

$y \in C'$

Проверим выполнимость свойств метрики:

1) Симметричность очевидна

$$2) p_2(p, P) = 1 - \bigwedge_{r \in C'} \min(p(r), P(r)) = 1 - \bigvee_{r \in C''} (\neg(r)) = 1 - 1 = 0$$

3) Неравенство треугольника.

Необходимо доказать, что для любой тройки словарных функций  $p, P, Q$  выполнено условие:

$$p_2(p, P) \leq p_2(p, Q) + p_2(Q, P), \text{ что эквивалентно}$$

$$1 - 2^{-\min_{y \in C'} (\bigwedge_{r \in C'} p(r) \wedge P(r))} \leq 1 - 2^{-\min_{y \in C'} (\bigwedge_{r \in C'} p(r) \wedge Q(r))} + 1 - 2^{-\min_{y \in C'} (\bigwedge_{r \in C'} Q(r) \wedge P(r))} \wedge 0,$$

или

$$2^{-\min_{y \in C'} (\bigwedge_{r \in C'} p(r) \wedge Q(r))} + 2^{-\min_{y \in C'} (\bigwedge_{r \in C'} Q(r) \wedge P(r))} \leq 2^{-\min_{y \in C'} (\bigwedge_{r \in C'} p(r) \wedge P(r))} < 1$$

Рассмотрим следующие случаи:

а)  $y: p_i(y) \leq p_2(y) \wedge p_3(y)$ , обозначим множество таких  $y$  через  $y_{123}$



- б)  $y: p_1(y) < p_3(y) \prec P_2(y)$ ? обозначим множество таких  $y$  через  $y_{123}$
- в)  $y: p_2(y) \prec p_1(y) - P_3(y)$ ? обозначим множество таких  $y$  через  $y_{213}$
- г)  $y: P_2(y) \wedge p_3(y) \prec P_1(y)$ ? обозначим множество таких  $y$  через  $y_{231}$
- д)  $y: P_3(y) \prec p_1(y) \wedge P_2(y)$  обозначим множество таких  $y$  через  $y_{312}$
- е)  $y: P_3(y) \prec P_2(y) \prec P_1(y)$  обозначим множество таких  $y$  через  $y_{321}$

Для каждого из случаев а)-е) посчитаем, чему равно  $f(y) =$

$$\min(p_1(y), p_2(y)) + \max(p_2(y), p_3(y)) - \max(p_1(y), p_3(y))$$

а)  $f(y) = P_1(y) + P_2(y) - p_1(y) = P_2(y) \quad \forall y \in y_{123}$

б)  $f(y) = P_1(y) + P_3(y) - P_1(y) = P_3(y) \quad \forall y \in y_{132}$

в)  $f(y) = P_1(y) + P_2(y) - p_2(y) = P_1(y) \quad \forall y \in y_{213}$

г)  $f(y) = P_3(y) + P_2(y) - p_2(y) = P_3(y) \quad \forall y \in y_{231}$

д)  $f(y) = P_3(y) + P_3(y) - p_1(y) = 2p_3(y) - p_1(y) \quad \forall y \in y_{312}$

е)  $f(y) = P_3(y) + P_3(y) - p_2(y) = 2p_3(y) - p_2(y) \quad \forall y \in y_{321}$

Нам необходимо доказать, что  $\bigcup_y f(y) \neq \emptyset$

Пользуясь тем, что все множества  $y_r$  попарно не пересекаются, преобразуем левую часть неравенства:

$$\begin{aligned} & \bigcup_y f(y) = \bigcup_{y \in y_{123}} P_2(y) + \bigcup_{y \in y_{132}} P_3(y) + \bigcup_{y \in y_{213}} P_1(y) + \bigcup_{y \in y_{231}} P_3(y) + \bigcup_{y \in y_{312}} (2p_3(y) - p_1(y)) + \bigcup_{y \in y_{321}} (2p_3(y) - p_2(y)) \\ & = 2 \cdot \bigcup_{y \in y_{123}} P_2(y) + \bigcup_{y \in y_{132}} P_3(y) + \bigcup_{y \in y_{213}} P_1(y) + \bigcup_{y \in y_{231}} P_3(y) + \bigcup_{y \in y_{312}} (2p_3(y) - p_1(y)) + \bigcup_{y \in y_{321}} (2p_3(y) - p_2(y)) \\ & = 2 \cdot \bigcup_{y \in y_{123}} P_2(y) + \bigcup_{y \in y_{132}} P_3(y) + \bigcup_{y \in y_{213}} P_1(y) + \bigcup_{y \in y_{231}} P_3(y) + \bigcup_{y \in y_{312}} (2p_3(y) - p_1(y)) + \bigcup_{y \in y_{321}} (2p_3(y) - p_2(y)) \\ & = 2 \cdot \bigcup_{y \in y_{123}} P_2(y) + \bigcup_{y \in y_{132}} P_3(y) + \bigcup_{y \in y_{213}} P_1(y) + \bigcup_{y \in y_{231}} P_3(y) + \bigcup_{y \in y_{312}} (2p_3(y) - p_1(y)) + \bigcup_{y \in y_{321}} (2p_3(y) - p_2(y)) \end{aligned}$$

Далее, поскольку  $p_3(y) - p_2(y) < 0$  при  $y \in \mathcal{Y}_2$ , то  $\int p_3(y) d\mathbb{P} - \int p_2(y) d\mathbb{P} < 0$ .

Аналогично  $\int (p_3(y) - p_2(y)) d\mathbb{P} < 0$ .

Так как  $p_2(y) < p_3(y)$  при  $y \in \mathcal{Y}_2$ , то  $\int p_2(y) d\mathbb{P} < \int p_3(y) d\mathbb{P}$ . Аналогично

$\int p_3(y) d\mathbb{P} < \int p_2(y) d\mathbb{P}$ .

Следовательно,  $\int p_2(y) d\mathbb{P} < \int p_3(y) d\mathbb{P}$  и  $\int p_3(y) d\mathbb{P} < \int p_2(y) d\mathbb{P}$ .

$\int p_2(y) d\mathbb{P} < \int p_3(y) d\mathbb{P}$  и  $\int p_3(y) d\mathbb{P} < \int p_2(y) d\mathbb{P}$  и требовалось доказать.

Утверждение доказано.

**Замечание 3.1.** Функция  $p_2$  не является метрикой на множестве  $\mathcal{Y}_2$ , т.к. свойство 2 метрики для слабостохастических словарных функций не выполнимо.

Будем называть стохастические словарные функции  $p_1$  и  $p_2$  «абсолютно разными», если выполнены условия:

1. Если  $p_1(y) > 0$ , то  $p_2(y) = 0$
2. Если  $p_2(y) > 0$ , то  $p_1(y) = 0$

**Утверждение 3.3.** На множестве стохастических словарных функций 4 метрики  $p_1$  и  $p_2$  совпадают с точностью до постоянного множителя.

Доказательство.





Поскольку ряд  $\sum_{i=0}^{\infty} (P_i(y), p_2(y))$  сходится, можно вычислять рас-

стояние между автоматами со сколь угодно большой точностью  $\epsilon$ , задавшись достаточно большим натуральным  $N$ , по формуле:

$$d(A, B) = 1 - \min_{i \leq N} (P_i(z), P_2(z))$$

Доказанное утверждение позволяет оценить число матричных умножений при вычислении расстояния между автоматами для векторов  $UQ$  и  $i_F$  специального вида. Сформулируем это более строго в виде следствия:

**Следствие 3.3.** Пусть даны два монотонных автомата  $A_1 = \{C, Q_1, T_1, P_1\}$  и  $A_2 = \{C, C_2, T_2, P_2, Q_2, P_2\}$ , такие, что  $|C_1| = |C_2| = T$ ,  $Q_2 = (1, \dots, 0)$ ,  $P_2 = (0, 0, \dots, 1)$ .

Пусть задано натуральное  $N$ . Тогда для приближенного вычисления расстояния между автоматами по формуле  $d(A_1, A_2) = 1 - \min_{i \leq N} (M_i(y), M_2(y))$

требуется порядка  $2m^N$  умножений матриц размерности  $m \times m$ .

#### Доказательство

Для расчета  $M_i(y)$  требуется  $|y|$  матричных умножений. Поскольку для вычисления  $M_i(y)$  можно воспользоваться индуктивной формулой  $M_i(y) = M_i(y')M_i(c)$ , где  $y = y'c$ ,  $c \in C^*$ ,  $c \in C$ , т.е. достаточно одного дополнительного матричного умножения, то общее число умножений матриц не превышает  $\sum_{j=0}^{N-1} 2^j m \sim 2^N m$ .

Следствие доказано.

Заметим, что если не использовать зависимость  $M_i(y) = M_i(y')M_i(c)$ , то для приближенного вычисления расстояния между автоматами требуется



уже порядка  $2Nm^{N+1}$  умножений матриц, так как общее число умножений матриц в этом случае равно  $2^{k-1} 2km \approx 2^{k-1} \frac{Nm^{N+1} - (N+1)m^{N+1}}{m-1} \approx 2Nm^{N+1}$ .

Пусть решается задача распознавания речи алгоритмом  $(E, p)$  в словаре  $B$  из двух команд  $B = \{B_1, B_2\}$ . Эту задачу можно рассматривать как задачу классификации ([35]) слов, выдаваемых соответствующими командами вероятностными автоматами  $A_1$  и  $A_2$ . А именно, пусть в каждый момент времени срабатывает только один автомат, автоматы равновероятны, и задача состоит в том, чтобы определить, какой из двух автоматов выдал заданное кодовое слово. Данная задача решается с некоторой ошибкой - *ошибкой классификации*, вероятностью которой будет вероятность такого события, когда слово было выдано одним автоматом, но отнесено нашим алгоритмом распознавания  $(E, p)$  к выходу другого. Ошибку классификации можно называть *ошибкой распознавания*, если выдача кодового слова у автоматом  $A_1$  происходит точно тогда, когда  $y \in \Gamma(\Pi(p_1))$ .

Вероятность ошибки классификации будет задаваться по формуле:

$$P_{\text{ош}} = P\{E(p)(y) = B_1 | y, A_1\}P(A_1) + P\{E(p)(y) = B_2 | y, A_2\}P(A_2) = z^*?$$

$$= P\{E(p)(y) = B_1 | y, A_1\}P(A_1) + P\{E(p)(y) = B_2 | y, A_2\}P(A_2)$$

Здесь  $P(D_i) = 1/2$  - вероятность "срабатывания"  $i$ -го автомата,  $P\{E(p)(y) = B_2 | y, A_1\}$  - вероятность того, что классификатор определил, что "сработал" 2-й автомат, а на самом деле у выдано первым автоматом,  $P^T(E, p)(y) = \text{fix} \setminus \Gamma \Pi$  - наоборот.

Справедлива

Доказательство.

$$\frac{1-2p_{\text{out}}}{2^{\min(A(r), A(r))}} = \prod_{y \in C^*: p_+(y) < p_-(y)} A(\Gamma_y) \cdot \prod_{y \in C^*: p_-(y) < p_+(y)} 2^{-A(\Gamma_y)}$$

Лемма доказана.

Содержательно данная лемма связывает расстояние между двумя автоматами как эталонами речевых объектов (слов, звуков и т.п.) с вероятностью ошибки распознавания в словаре из этих двух объектов, т.е. с вероятностью "спутать" произносимые слова (звуки и т.п.) при условии, что обучение (синтез) эталонов команд произведено идеально.





Введем теперь *функцию качества словаря команд* для вероятностно-автоматного алгоритма  $(E, p)$  распознавания речи.

Для словаря  $B$  из двух команд  $S = \{p_1, p_2\}$  определим  $Q(E, p)(CB)$  как  $Q(E, p)(CB) = 1 - 2p_{\text{ош.}}$ , где  $p_{\text{ош.}}$  - вероятность ошибки распознавания в этом словаре команд, о которой говорилось выше. Для словаря с произвольным числом команд, как и раньше, положим  $Q(E, p)(CB) = 0_B(Q(E, p))$ .

**Теорема 3.2.** Пусть  $B = \{p_1, p_2, \dots, p_n\}$  - словарь команд,  $A = \{A_1, A_2, \dots, A_n\}$  - соответствующие этим командам вероятностные автоматы. Тогда  $Q(E, p)(CB) = 0_A(p_2)$ .

Доказательство.

Для доказательства теоремы достаточно применить лемму 3.4. Для каждой пары  $\{p_1, p_2\}$  команд из словаря  $B$  имеем  $Q(E, p)(CB) = 1 - 2p_{\text{ош.}} = p_2(A_1, A_1)$ . Усредняя по всем парам команд, получаем искомое утверждение теоремы.

Теорема доказана.

Теорема 3.2 дает нам способ точной оценки качества произвольного словаря команд для алгоритма распознавания речи методом скрытых марковских моделей. Качество определяется с помощью метрики, введенной на множестве эталонов команд.

### §3.4 Метрика $\rho_3$ на множестве стохастических словарных функций.

Третий вариант введения метрики на множестве стохастических словарных функций состоит в следующем.

Если представить стохастическую словарную функцию  $\{(y, p(y))\}$  в виде стохастического вектора  $p = (p(i) \in PC^*, \dots, p(y_0 \in \dots) \in C^*)$  то  $\hat{p}(y) = I_{y \in C^*}$

$\Rightarrow \sum_{y \in C^*} \hat{p}(y) = 1$ , значит,  $p \in \mathcal{L}_2$ , где  $\mathcal{L}_2 = \{(x_1, x_2, \dots, x_n, \dots) : \sum_{i=1}^{\infty} x_i^2 < \infty\}$  ([29]).

$\mathcal{L}_2$  является евклидовым пространством со скалярным произведением  $(x, y) = \sum_{i=1}^{\infty} x_i y_i$ , нормой  $\|x\| = \sqrt{(x, x)}$  и метрикой  $\rho(x, y) = \|x - y\|$  ([29]).

Можно ввести в  $\mathcal{L}_2$  отношение эквивалентности следующим образом:  $x \sim y$ , если  $\exists k \in \mathbb{R}, k \neq 0: x = ky$ . Это отношение разбивает  $\mathcal{L}_2$  на классы эквивалентности, причем в каждом классе (кроме класса, состоящего из нулевого вектора) существует ровно один элемент, принадлежащий множеству  $\mathcal{L}$  стохастических векторов. Кроме того, каждый такой ненулевой класс эквивалентности содержит ровно один элемент из  $\mathcal{L}_2$ , норма которого равна 1.

Введем расстояние между стохастическими векторами следующим образом: расстоянием между  $p_1$  и  $p_2$  назовем расстояние между представителями классов эквивалентности, содержащими  $p_1$  и  $p_2$ , норма которых равна 1. Более строго:

**Определение 3.11.** Расстоянием  $\rho_3$  между стохастическими функциями  $p_1$  и  $p_2$  назовем величину  $\rho_3(p_1, p_2) = \inf_{\substack{A \in \Pi_1 \\ B \in \Pi_2}} \|A - B\|$

**Утверждение 3.5.**  $\rho_3$  - метрика на множествах  $\mathcal{L}_s$  и  $\mathcal{L}'_s$ .

### Доказательство.

- 1)  $\rho_3(p, p) = 0$  - очевидно
- 2)  $\rho_3(p_1, p_2) = \rho_3(p_1, p_2)$  - очевидно
- 3) Неравенство треугольника для метрики  $\rho_3$  и трех произвольных словарных функций  $p_1, p_2, p_3$  следует из справедливости неравенства треугольника для  $A \gg A_p$  и  $D_{\Gamma}$  для стандартной метрики в  $\mathbb{R}^n$ .

### Утверждение доказано.

Как известно, косинусом угла между векторами  $p_1$  и  $p_2$  евклидова пространства называется величина  $\cos Z(p_1, p_2) = \frac{(p_1, p_2)}{\|p_1\| \|p_2\|}$  ([29]).

**Утверждение 3.6.**  $\rho_3(p_1, p_2) = \sqrt{1 - \cos Z(p_1, p_2)}$ .

### Доказательство.

$$2\rho_3(p_1, p_2)^2 = \|p_1 - p_2\|^2 = (p_1 - p_2, p_1 - p_2) = \|p_1\|^2 + \|p_2\|^2 - 2(p_1, p_2) = 2 - 2\cos Z(p_1, p_2),$$

откуда  $\rho_3(p_1, p_2) = \sqrt{1 - \cos Z(p_1, p_2)}$

### Утверждение доказано.

### **Следствие 3.4.**

- 1) Для любых двух стохастических словарных функций  $p_1, p_2$  верно, что  $0 \leq \rho_3(p_1, p_2) \leq 1$
- 2)  $\rho_3(p_1, p_2) = 1$  точно тогда, когда  $\cos Z(p_1, p_2) = -1$
- 3)  $\rho_3(p_1, p_2) = 0$  точно тогда, когда  $p_1 = p_2$ .

Доказательство следует непосредственно из формулы  $\rho_3(p_1, p_2) = \sqrt{1 - \cos Z(p_1, p_2)}$ .



### Определение 3.12. Декартовым произведением $A \times A_2$ автономных

вероятностных

автоматов

$$A \sim \{C \& \setminus, \Pi, P \setminus, \mathcal{U} \setminus o, \mathcal{U} \setminus \mathbb{N}\}$$

и

$$\Lambda^2 = (0, 02, 712^{\wedge 2} 0^{\wedge})$$

назовем

вероятностный

автомат

$74 = (C, Q_i \times Q_2, 7C, P, v_0, (iF)$ , где  $\kappa$  имеет размерность  $m = \min_2 (m_j = |Q_j|, i=1,2)$  и

$$\% \text{ЛШ, } \text{Л})^{\wedge}, \text{ }^{\text{ш}} * u >^P (Y) \kappa = P_{i,k} - P_{2,j} b V_0^{(i)} = V! o' - V_{i_0}^j, V_F^{(i)} = v_i p^1 - V_{,F}^j.$$

**Теорема 3.3.** Пусть  $p_1$  и  $p_2$  - стохастические словарные функции,

порожденные

монотонными

автоматами

$$4i = \{C, Q_i, 7tbP_i, V_{i_0}, v_{iF}\}$$

и

$$\Lambda^2 = (C, Q_2, 7I_2, P_2, v_2O, V_2F), A = A_i \times A_2 = \{C, Q_i \times Q_2, n, ?, Vo, v_{\mathbb{N}}\} \text{ -декартово произ-}$$

ведение автоматов  $A \setminus$  и  $A_2$ ,  $k = |C|$ ,  $m = |Q_i \times Q_2|$ . Тогда  $(p_1 p_2 \geq ((E - m \varepsilon y^I \setminus_m$

$\kappa$

где  $n \setminus y_y = y^{\wedge} P_a$  -приведенная матрица переходов автомата  $A$ .

Доказательство.

$$(P_1 P_2) = \sum_{y \in C} E A(r) A(r) = \sum_{y \in C} \sum_{v \in \mathcal{U}} \sum_{\mathcal{U}} 2^{\wedge 0^{\wedge}} \wedge 0^{\wedge})$$

$f$

$$\sum_{y \in C} E \text{ ЦР}^{\wedge} \Gamma' Y M I U' Y I) = \sum_{r \in \setminus \mathcal{U}_i, \mathcal{U}_i} \left( \text{ЦР}(\Gamma^{\wedge} X \mathcal{U}_2) \right) = \sum_{y \in C} \text{£P00-}$$

Здесь  $q_x$  - путь из начального состояния в финальное в автомате  $A \setminus$ ,  $q_2$  - путь из начального состояния в финальное в автомате  $A_2$ ,  $q_x \times q_2$  - путь из начального состояния в финальное в автомате  $A = A \setminus \times A_2$ ,  $p(y, q)$ - совместная вероятность выдачи слова  $y$  и прохождения по пути  $q$ ,

$$P(\Gamma, I O^{\wedge} P(Y \setminus I O P(\mathcal{U})-$$

По теореме 3.1  $\text{£p00} = ((E - m \varepsilon y^I \setminus_m \geq (p_1 p_2) = ((\text{£} - \text{£})^{\text{''}0\%}$ , что и тре-

$y \in C$

бовалось доказать.

Теорема доказана.

Содержательно теорема 3.3 означает, что скалярное произведение двух стохастических словарных функций эффективно вычисляется как вероятность дойти из начального состояния до финального в декартовом произведении любых двух автоматов, реализующих эти словарные функции. В свою очередь, вероятность дойти из начального состояния до финального вычисляется путем нахождения правого верхнего углового элемента в предельной приведенной матрице переходов  $(E - \gamma)^{-1}$ .

**Следствие 3.5.** Расстояние между стохастическими словарными функциями  $p_1$  и  $p_2$ , реализующими автоматы  $A_1$  и  $A_2$ , эффективно вычисляется по формуле

$$P_3(p_1, p_2) = \sqrt{1 - \cos(\angle)} = 1 - \sqrt{1 - \cos(\angle)}$$

приведенная матрица переходов декартова произведения  $A = A_1 \times A_2 = \{C, Q_1 \times Q_2, P, \gamma, L_1 \cup L_2\}$  автоматов  $A_1$  и  $A_2$ ,  $L_1$  - приведенная матрица переходов декартова произведения  $D_1 \times L_1$ ,  $L_2$  - приведенная матрица переходов декартова произведения  $A_1 \times A_2$ .

Доказательство получается непосредственной подстановкой формулы  $(p_1, p_2) = (E - \gamma)^{-1}$  в формулу  $P_3(p_1, p_2) = \sqrt{1 - \cos(\angle)}$  для вычисления расстояния  $p_3$ .

Главным итогом этого параграфа является существования способа эффективного введения метрики на множестве автономных монотонных вероятностных автоматов. Для вычисления метрики  $p_3$  оказывается достаточным применить конечное число операций типа умножения матриц, нахождения обратной матрицы и определителя, извлечения квадратного корня, умножения и деления действительных чисел.



По аналогии с тем, как мы вводили функцию качества с помощью метрики  $p_2$ , можно дать следующее определение:

Пусть  $\delta = \{p_1, p_2, \dots, p_n\}$  - словарь команд,  $A'' = \{A_1, A_2, \dots, A_n\}$  - соответствующие этим командам вероятностные автоматы. Тогда *качеством* этого словаря назовем величину  $\chi'(E, p)CB = \theta_d(p_3)$ .

### §3.5 Связь между метрикой и вероятностью.

Пусть слово  $y_0 \in C^*$ .

**Определение 3.13.** *Характеристической словарной функцией слова  $U$  назовем следующую словарную функцию:*

$[0, \text{ иначе}]$  Пусть  $y = c_1 c_2 \dots c_n$ . Тогда монотонный вероятностный автомат, соответствующий функции  $pp$ , можно задать автоматом, приведенном на рис. 3.2:

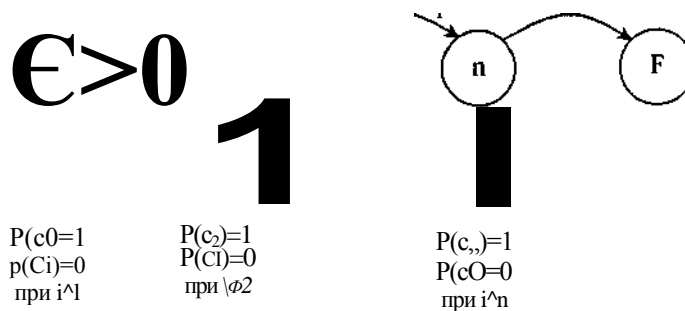


Рисунок 3.2.

Используя характеристическую словарную функцию для кодового слова, можно ввести расстояние между кодовыми словами и словарными функциями:

**Определение 3.14.** Пусть  $p$  - метрика на множестве  $/_s$  (или  $/_D$ )  $u \in C^*$ ,  $p \in /_s(/_D)$  Расстоянием между словом  $u$  и словарной функцией  $p$  назовем величину  $p(u, p) = p(p_Y, p)$ .

Посмотрим, чему равно расстояние между словом и словарной функцией для введенных ранее метрик  $p_1$ ,  $p_2$  и  $p_3$ .

**Лемма 3.5.** Пусть  $y \in C^*$ ,  $p \in \mathcal{P}_s$ . Тогда

$$(1) \quad p_1(y, p) = 2(1 - p(y)),$$

$$(2) \quad p_2(y, p) = 1 - p(y)$$

$$(3) \quad P_3(T, P) = \sum_{i=1}^n p_i, \text{ - норма словарной функции } p \text{ в } \mathcal{P}_s,$$

Доказательство.

$$(1) \quad p_1(y, p) = p_1(p_Y, p) = \sum_{y \in \Phi_Y} |P_Y(\Gamma') - P(\Gamma')| = \sum_{y \in \Phi_Y} |P_A(\Gamma') - P(\Gamma')| + 1 - P(Y)$$

$$2 \sum_{y \in \Phi_Y} (1 - p(y)) + 1 - P(Y) = 2 \sum_{y \in \Phi_Y} (1 - p(y)) + 1 - P(Y) = 1 - p(y) + 1 - p(y) = 2(1 - p(y)).$$

$$(2) \quad p_2(y, p) = p_2(p_Y, p) = 1 - \sum_{y \in C} mn(p_Y(y), p(y)) = 1 - p(y)$$

$$(3) \quad p_3(y, p) = P_3(p_Y, P) = \sum_{i=1}^n p_i$$

$$\sum_{i=1}^n p_i = \sum_{i=1}^n P_i(p_Y) = 1, (P_Y > P) = \sum_{i=1}^n P_i(p_Y) = P(Y) > \text{откуда } p_3(y, p)$$

$U \quad P(r)$   
 $V \quad \Pi$

Лемма доказана.

Из леммы 3.5 следует, что введенная в §3.1 при задании алгоритма распознавания  $(E, p)$  функция  $p$  расстояния между словарными функциями и кодовыми словами совпадает в смысле определения 3.14 с метрикой

**Следствие 3.6.** Если нам дано расстояние от слова  $y \in C^*$  до стохастической функции  $p \in \mathcal{P}_s$ , то можно рассчитать вероятность  $p(y)$  по одной из следующих формул:

$$(1) \quad p(Y) = \sum_{i=1}^n (1 - p_i(Y, p));$$

$$(2) \quad p(y) = 1 - p_2(y, p);$$

$$(3)P(Y)=|P\Pi(1-P3^2(Y,P)).$$

Метрика на вероятностных автоматах может также применяться для решения задачи **сокращения фонетического алфавита**. Коротко постановка и практический смысл этой задачи состоит в следующем. При распознавании речи вероятностные автоматы, выступающие в качестве эталонов слов словаря, обычно синтезируются путем последовательного соединения вероятностных автоматов, являющихся эталонами фонем, составляющих фонетическую транскрипцию слова.

Поскольку качество обучения (синтеза) автоматов зависит прямо пропорционально как от объема обучающей выборки, так и от количества автоматов, которые необходимо обучить, то одной из задач создания распознавателя речи является задача сокращения фонетического алфавита, имеющая важное практическое значение.

Задачу эту можно решать путем выбора в качестве начального приближения фонетического алфавита максимального алфавита, включающего в себя все возможные символы транскрипции с максимальной детализацией, а затем путем разбиения этого алфавита на классы эквивалентности и использования этих классов в качестве нового фонетического алфавита на следующем, заключительном шаге.

В один класс эквивалентности целесообразно помещать фонемы, которые близки друг к другу в смысле некоторой меры близости. Отсюда напрямую вытекает необходимость построения метрики на множестве фонем данного языка (фонетическом алфавите).

Поскольку фонемы моделируются вероятностными автоматами, возникает необходимость введения *метрики* на множестве вероятностных автоматов без входа.

Метрика  $r_3$  была приведена на практике для решения задачи сокращения фонетического алфавита для русского языка. В результате



проведенного исследования было показано, что можно сократить алфавит фонем для русского языка со 143 до 132 элементов.

## ПРИЛОЖЕНИЕ 1. ПРИМЕР АЛФАВИТА БУКВ И ЗВУКОВ, ПРАВИЛ ТРАНСКРИБИРОВАНИЯ ДЛЯ РУССКОГО ЯЗЫКА.

Алфавит букв русского языка: { а, б, в, г, д, е, ё, ж, з, и, й, к, л, м, н, о, п, р, с, т, у, ф, х, ц, ч, ш, щ, ь, ы, ь, э, ю, я, ' }.

Алфавит звуков русского языка (согласно [41], [42])

Символом " \* " обозначается мягкость звука слева, справа или с двух сторон (в зависимости от того, с какой стороны стоит символ " \* ").

### ГЛАСНЫЕ

Символом "+" обозначается гласный под ударением.

#### ГЛАСНЫЕ ПОД УДАРЕНИЕМ

Обозначение фонемы	Буква	Примеры	Транслитерация
a+	а	арка, буржуа	al
a*+	а	рвань	ail
*a+	а	чан	ial
	я	яма	
*a*+	а	часть	iail
	я	бязь	
o+	о	дом, наотмашь	ol
o*+	о	озеро, моль, наощупь	oil
*o+	о	трущоба,	iol
	е	черт, пес, ружье	
*o*+	о	на Печоре, йодистый	ioil
	е	песик	
y+	у	угол, жук, какаду, наука, недоумок	
y*+	у	пуля, усики, научит, баульчик	uil
*y+	у	чудо, свищу	iul
	ю	дюны, юг	
*y*+	у	чучело	iuil
	ю	бюстик	
i+	и	шина, цирк	yi
	ы	огурцы, пыл, сады	
i*+	и	истина, ширь, циник	yil
*i+	и	бинт	il

*и*+	и	линь	П1
	э	этот, пэр	el
e+	е	кафе, шест	
	э	эти, о пэре	eil
e*+	е	дельта, шесть	
	е	тесто, ель	iel
*e+	е	сельский, мель, свистеть	ieil

## БЕЗУДАРНЫЕ ГЛАСНЫЕ

**СИЛЬНЫЕ БЕЗУДАРНЫЕ** в 1-ом предударном слоге или после паузы в начале слова или после согласной перед паузой или последняя гласная в группах гласных.

Обозначение фонемы		Примеры	Транслитерация
Буква			
	О	оса. пока	а
	а	аванс, сажа	
	о	жокей	ai
	я	катя пять	ia
		атак	
4' Q 4'		пять осин	iai
У2у*2	У	указ	u
*у2	У	пустяк	ui
*у*2	ю	юла	iu
	У	чудить	III
И	ю	тюлень, юлить	
	и	широк	У
и* *и	ы	мытарить	UI
* JJ*	и	зима	i
	я	пятак	
	и	мирить	ii
	э	экран	e
ip nф	э	электростанция	ei
	е	поле	ie
	е	еще	iei



## СЛАБЫЕ БЕЗУДАРНЫЕ -

только после согласных и  
не в 1-ом предударном слоге и  
не в конце слова перед паузой.

Обозначение фонемы	Буква	Примеры	Транслитерация
Ъ*	а	балагур, шаловлив	ах
	ы	пасынок	
	о	богатырь, шоколад	
	е	сторожем	
	о	вымочить	ахі
уі	а	выкатить	
	У		w
	У	путевой, сулема	wі
у*і			
*уі	ю	полос, вынюхать, кюрасо	іw іwі
*Ъ	а	часовой	і
	е	камнем	
*Ъ*	я	впятером	іі

В ИНОСТРАННЫХ СЛОВАХ гласные следует обозначать в соответствии с тем, как они слышатся.

## СОГЛАСНЫЕ

Символом ":" обозначается длительность согласного (например, в случае удвоенного согласного).

### СОНОРНЫЕ ЗУБНЫЕ

Обозначение фонемы	Буква	Примеры	Транслитерация
Р	р	рама, пар, тигр	г
Р:р*	рр	суррогат	гг
р*:	р	рюмка, январь	ГІ
н	рр	территория	ггі
н:	н	нам, сын, канва	п
н*	нн	ванна, колоннада	пп
н*:	н	няня, дрянь, синька, бантик	пі
л	нн	теннис	ШІІ
л:	л	лавка, стол, цикл	І
л*	лл	балласт, мулла	И
л*:	л	лягу, моль, льстить	ІІ
	лл	капилляр	Ш

## ГУБНЫЕ

Обозначение фонемы	Буква	Примеры	
		Транслитерация	
	м	мыловар, дом, реноме	m
м			
м:	мм	грамматика	mm
м*	м	мял, семь, познакомьте	mi
м*:	мм	комментатор	mmi
б	б	бак, объект	b
б:	бб	аббат	bb
б*	б	бязь, бязь	bi
б*:	бб	баббит	bbi
п	п	пора, суп	P
	б	зуб, герб, трубка	
п:	пп	аппарат, группа	PP
п*	п	пятый, сыпь	Pi
	б	дробь, приспособьте	
п*:	пп	аппетит	ppi
	в	вуз, враг	v
	ф	Афганистан	
	в	ввоз	w
	в	вялый, выюга, въехать, вбить	vi
в .	в	ввел	wi
Ф	ф	факт, шеф	f
	в	остров, вкус	
ф:	фф	диффузный	ff
ф*	ф	тюфяк, верфь, потрафьте fi	
	в	кровь, приготовьте, впишет	
ф*:			ffi

## ЗУБНЫЕ ШУМНЫЕ

Обозначение фонемы	Примеры		Транслитерация
Буква			
	з	зал, призрак, язва	z
	с	сбор, сгубить	
з:	зз	беззаботный	zz
	сз	сзади	
	з	зябну, резьба, изъездить, уздечка zi	
	с	сбить	
	зз	безземельный	zzi
	с	сало, нос, стол	s
	з	глаз, смазка	
	ее	компрессор, касса, ссадина	
	с	сс	
		снесся	



<b>Д</b>  Д: Д* Д <sup>н</sup>	З	грызся	
	С	сяду, весь, моська	si
	З	мазь, лезть	
	ее	рассердиться	ssi
	Д	дать, подъем	d
	Т	отбить	
	ДД	аддуктор	dd
	Д	дятел, свадьба, две	di
	ДД	поддержка	ddi
	Т	там, пот, отъезд	t
	Д	яд, подкинуть	
	ТТ	атташе, оттащить, гетто	tt
	Т	тяга, петъ, отняли	ti
	Д	сельдь, обладьте, складчина	
гр*К•	ТТ	аттестат, аттический, оттесывать	tti

### ЗАДНЕБНЫЕ

Обозначение фонемы Буква		Примеры	Транслитерация
<b>Г:</b>	Г	газ, когда	
	К	анекдот	
	Г	гиоель	gl ggi
	К	кадка, так, краб	k
<b>К:</b> К* К*	Г	ногти, берег	
	КК	мокко	kk
	К	кинуть	ki kki
	Х	хата, задохся, слух	h
<b>Х:</b>	Г	мягко	hh
	Х	хитрый	hi
	Г	мягкие	hhi

### ШИПЯЩИЕ И [Ц]

Обозначение фонемы Буква		Примеры	Транслитерация
<b>Ж</b>	Ж	жалъ, вежливый, ружье	zh
	ЗЖ	разжать, езжу	
	ЕЖ	сжать	
	Ж	ложбина	



ж: ж*	жж	жужжать	zzh
ж*: ш			zhi
			zzhi
ш:	ш	шаг, шью, шла, грош, мышь	sh
	ж	нож, рожь, ложка	
щ*	с	сшить	ssh
	з	лезший	
щ*:	щ	пощада, вещью, плющ	sch
	сч	счастье, разносчик	
	сщ	расщелина	ssch
ч	34	навязчивый	
ч:	жч	мужчина	ch
			cch
ч*:			
ц	ч	часть, чью, чрево, врач	chi
ц:			cchi
	ц	цоколь, матрац, клетки	ts
			tts
ЙОТ			
й*	я	яблоко	j

### ОЗВОНЧЕННЫЕ (на стыке слов)

Обозначение фонемы	Буква	Примеры	Транслитерация
(хг)	хг	господи, белых гусей	x
(дз)	цб	отец бы	dz
(дж*)	чб	дочь бы	dxh

### СЛУЖЕБНЫЕ ФОНЕМЫ

СМЫЧКИ обозначаются соответствующим согласным, заключенным в квадратные скобки:

[т], [т\*], [д], [д\*], [п], [п\*], [б], [б\*], [р], [р\*], tcl ticl del did pel picl bcl bid rcl ricl  
 [к], [к\*], [г], [г\*], [ц], [дз], [ч], [ч\*], [дж\*], kcl kicl gel gicl tscl dzcl chel dzhel  
 [т:], [т\*:], [д:], [д\*:], [п:], [п\*:], [б:], [б\*:], [р:], [р\*:], ttcl tticl did ddicl picl ppicl  
 bbel bbicl pel ppicl  
 [к:], [к\*:], [г:], [г\*:], [ц:], [ч:], [ч\*:], [дж\*] kkcl kkicl ggcl ggicl tscl echel  
 dzhel

ПАУЗА в начале фразы и в конце фразы - %%  
 ПАУЗА между словами - %  
 ПАУЗА вставочная (эпентетическая) - []

**Транскрибирование производится по следующим правилам:**  
(согласно [31], [41], [42])

**Переход "буква-фонема".**

Осуществляет такие операции над орфографической записью, как устранение орфографических фикций, обработка особых случаев произнесения стечений согласных, устранение твердых и мягких знаков, обработка йотированных и мягких букв с соответствующей интерпретацией твердости-мягкости соседних согласных и введением йота.

*Пример правил этого блока:*

- 1) здн—»з н, лнц->н ц, стн-^с н, рдц—»р ц, и т.п. (пропуск непроезжих согласных)
- 2) тся-»ц ц а, тья-»ц ц а, что->ш т о, его->е в о й т.п. (замена сочетаний согласных в окончаниях)
- 3) сш-»ш ш, сщ-»щ, стч-»щ, тц-»ц ц, и т.п. (ассимиляция по способу образования)
- 4) сть—»сь ть и т.п. (ассимиляция гомоорганых согласных по твердости-мягкости)
- 5) ль->ль (смягчение согласных перед мягким знаком и гласная)
- 6) бы'-»бь й и' в слове воробы и т.п. (правила вставления й в ударном и безударном слоге)
- 7) не-»нь е и т.п. (смягчение согласных перед йотированными гласными)
- 8) сс-»с, нн-»н, мм->м и т.п. (замена двух одинаковых согласных одним в некоторых позициях)

**Переход "фонема-звукотип".**

Включает правила, обрабатывающие случаи позиционного озвончения, оглушения согласных и редукции гласных (с учетом двух степеней редукции).

*Примеры правил озвончения-оглушения:*

1) б-»п в слове столб и т.п. (оглушение согласных в конце слов)

2) бс->пс в слове обстрел и т.п. (ассимиляция согласных по звонкости-глухости)

### *Правила редукции гласных:*

Ударный гласный — наиболее полногласный, т. е. наиболее мощный по длительности и интенсивности (обозначен цифрой 3).

1-ая степень редукции — на одну ступень ниже по интенсивности и длительности, т. е. (3 - 1) -> 2.

2-ая степень редукции ещё ниже, т. е. (3 - 2) -> 1.

Формула Потебни:

1	2	3	1	1
2-ой предуд. слог	1-ый предуд. слог	ударный слог	1-ый заударн. слог	2-ой заударн. слог

или:

#2	2	3	1	2#
2-ой предуд. неприкр. слог	1-ый предуд. слог	ударный слог	1-ый заударн. слог	2-ой заударн. открыт, слог

Это означает следующее: 1-ая степень редукции наблюдается в первом предударном слоге, в абсолютном начале и в абсолютном конце слова, т. е. после или перед пробелом (знаком "#"). 2-ая степень редукции наблюдается во всех закрытых сзади согласными заударных слогах и во втором, третьем и др. прикрытых спереди согласными предударных слогах.

Например: "повтори+ть" -> 12 3

"перестро+йка" -> 12 3 2

"за+поведь" -> 3 11

"анахрони+зм" -> 2 12 3

Знак "+" после гласного означает ударность гласного. Правила замен фонем на редуцированные звукотипы для каждой степени редукции:



### 1-ая степень редукции

фонема	позиционно обусловленные звукотипы					
	#+Г	Г+Г	С+Г	с+Г	Г+#	с__Г
0	а	а*	а*	и	а*	а
а	а	а	а*	и	а	а
э	э	э	ы*	и	э	э
ы	ы	—	ы	—	ы	ы
и	и	и	—	и	и	и
у	у2	у2	у2	у2	у2	у2

\* здесь означает наличие исключений, которые задаются списками

### 2-ая степень редукции

фонема	позиционно обусловленные звукотипы	
	с+Г	с+Г
0	ь	ь
а	ь	ь
э	ь	ь
ы	ь	—
и	—	ь
у	yi	yi

## ПРИЛОЖЕНИЕ 2. ПРИМЕР ФОНЕМНЫХ КЛАССОВ ДЛЯ РУССКОГО ЯЗЫКА

В работе [а5] использовались следующие классы фонем для русского языка:

**класс \_** ("паузы"): пауза, [п], [п\*], [к], [л\*], [т], [т\*], [ч], [ц]

**класс X** ("шипящие"):

с, с\* , ц, ч, ш, щ, в, в\* , д, д\* , ж, з, з\* , й, л, л\* , п, п\* , т, т\* , ф, ф\* , х, х\*

**класс А** ("низкочастотные гласные"): а, а+, э, э+ **класс О**

("среднечастотные гласные"): о, о+, у1, у2, у+ **класс I**

("высокочастотные гласные"):

и, и+, \*и, \*и+, и\*, и\*+, \*и\*, \*и\*+, а\*+, о\*+ и т.п., л, л\*, м, м\*, н, н\*, ЗВОНКИЕ

СМЫЧКИ

### ПРИЛОЖЕНИЕ 3. МЕТРИКА НА МНОЖЕСТВЕ ЗВУКОВ РУССКОГО ЯЗЫКА.

Для примера взяты 10 звуков русского языка из набора, приведенного в  
приложении 1 (реально матрица имеет размерность 150x150)

	[п]	<b>М</b>	*а*+	У+	и	м	л*	ш	с	ж
[п]	0	14.657	143.974	134.180	136.092	116.847	138.666	96.287	48.550	114.615
<b>М</b>	14.657	0	144.745	139.320	136.029	117.216	139.747	99.394	48.370	116.181
*а*+	143.974	144.745	0	116.897	125.756	105.587	119.652	143.963	144.392	135.984
У+	134.180	139.320	116.897	0	126.364	113.008	118.581	150.204	139.529	120.719
и	136.092	136.029	125.756	126.364	0	127.115	105.292	122.289	147.063	115.253
м	116.847	117.216	105.587	113.008	127.115	0	108.844	134.464	117.647	103.921
л*	138.666	139.747	119.652	118.581	105.292	108.844	0	137.965	148.690	90.546
ш	96.287	99.394	143.963	150.204	122.289	134.464	137.965	0	100.939	110.394
с	48.550	48.370	144.392	139.529	147.063	117.647	148.690	100.939	0	109.848
ж	114.615	116.181	135.984	120.719	115.253	103.921	90.546	110.394	109.848	0

## ЛИТЕРАТУРА

- [I] СЕ. Левинсон. Структурные методы автоматического распознавания речи. ТИИЭР, т.73, №11, ноябрь 1985 г., с. 100-128.
- [2] K.H. Davis, R.Biddulph, and S.Balashek, "Automatic recognition of spoken digits", J.Acoust. Soc. Amer., Vol. 24, pp.637-642, 1952.
- [3] P.B. Denes and M. V. Mathews, "Spoken digit recognition using time frequency pattern matching", J.acoust. Soc. Amer., vol. 32, pp. 1450-1455, 1960.
- [4] H. Dudley and S.Balashek, "Automatic recognition of phonetic patterns in speech", J. Acoust. Soc. Amer., vol. 30, pp. 721-7439, 1958.
- [5] D. H. Klatt, "Review of the ARPA Speech Understanding Project", J. Acoust. Soc. Amer., vol.62, no.6, pp. 1345-1366, Dec. 1977.
- [6] Г. Фант. Акустическая теория речеобразования / Пер. под ред. В.С. Григорьева. М.: Наука, 1964.
- [7] V. W. Zue and R. A. Cole, "Experiments on spectrogram reading" in Proc. ICASSP-79, pp. 116-119, 1979.
- [8] Т.К. Винцюк. Распознавание слов устной речи методами динамического программирования. Кибернетика. - 1968. - №1, с. 81-88.
- [9] Р. Беллман, Динамическое программирование. М.: ИЛ, 1960.
- [10] F. Itakura, "Minimum prediction residual principle applied to speech recognition", IEEE Trans. Acoust., Speech, Signal Processing, vol. ASSP-23, pp. 67-72, 1975.
- [II] J.K.Baker, "Stochastic modeling for automatic speech understanding", in Speech Recognition, D.R. Reddy, Ed. New York: Academic Press, 1975, pp. 521-542.
- [12] Ф. Джелинек [Елинек]. Распознавание непрерывной речи с помощью статистических методов. ТИИЭР, 1976, т. 64, №4, с. 131-160.

- [13] L.E. Baum and T. Petrie, "Statistical inference for probabilistic functions of finite state Markov chains", *Ann. Math. Stat.*, vol. 37, pp. 1554-1563, 1966.
- [14] А.А. Марков. Пример статистического исследования над текстом "Евгения Онегина", иллюстрирующий связь испытаний в цепь. *Известия Академии наук, СПб.*, VI, т.7, 1913, №3, с. 153-162.
- [15] S. E. Levinson, L.R. Rabiner, and M.M. Sondhi, "An introduction to the application of the theory of probabilistic functions of a Markov process to automatic speech recognition". *Bell Syst. Tech. J.*, vol 62, no.4, pp. 1035-1074, Apr. 1983.
- [16] B.H. Juang, "On the hidden Markov model and dynamic time warping for speech recognition - A unified view", *AT&T Tech. J.*, vol. 63, no.7, pp.1213-1243, Sept. 1984.
- [17] L.R.Rabiner and B.H. Juang, "An introduction to the hidden Markov models", *IEEE ASSP Mag.*, vol.3, no.1, pp.4-16, 1986.
- [18] Л.Р. Рабинер. "Скрытые марковские модели и их применение в избранных приложениях при распознавании речи: обзор". *ТИИЭР*, т.77, №2, февраль 1989 г.
- [19] L.E. Baum and J.A..Egon, "An inequality with applications to statistical estimation for probabilistic functions of a Markov process and to a model for ecology", *Bull. Amer. Meteorol. Soc*, vol.73, pp.360-363, 1967.
- [20] A. P. Dempster, N.M. Laird, and D.B. Rubin, "Maximum likelihood from incomplete data via the EM algorithm", *J.Roy. Stat. Soc*, vol. 39, no.1, pp.1-38, 1977.
- [21] A.J. Viterbi, "Error bounds for convolutional codes and an asymptotically optimal decodin algorithm", *IEEE Trans. Informat. Theory*, vol. IT-13, pp. 260-269, Apr. 1967.
- [22] Дж. Д. Форни-мл. "Алгоритм Витерби". *ТИИЭР*, 1973, т.61, №3, с.12-25.

- [23] J.W. Carlyle. Reduced forms for stochastic sequential machines. - J. Math. Analysis and Applic, 1963, 7, p. 167-175.
- [24] Р.Г. Бухараев. Некоторые эквивалентности в теории вероятностных автоматов. - Уч. записки Казан, университета, 1964. 124, №2, с. 45-65.
- [25] Р.Н. Starke. Theorie Stochastischen Automaten. I, II. - Elektron Informationsverarb. und Kybern., 1965, 1, №2. [26] Р.Г. Бухараев. Основы теории вероятностных автоматов. М.: "Наука", 1985.
- [27] R. Bakis, "Continuous speech word recognition via senti-second acoustic states" in Proc. ASA Meeting (Washington, DC), Apr. 1976. [28] D. Kanevsky, M. Monkowsky, J. Sedivy. Large Vocabulary Speaker-Independent Continuous Speech Recognition in Russian Language. Proc. SPECOM'96, St.-Petersburg, October 28-31, 1996. [29] А. Н. Колмогоров, С. В. Фомин. Элементы теории функций и функционального анализа. М.: Наука, 1989.
- [30] Дж.Бендат, А.Пирсол. Измерение и анализ случайных процессов. М.: "Мир", 1974 г., стр. 372, 342, 368.
- [31] Ипатов Я.В., Коваль С.Л., Лапшина А.В., Сапожкова И.Ф. Компилятивный синтезатор речи по правилам. Параметрическое представление. Лингвистическое обеспечение.: АРСО-15, стр. 156-157.
- [32] Л.Р.Рабинер, Р.В.Шафер, Цифровая обработка речевых сигналов. Пер. с английского. М., "Радио и связь", 1981 г. [33] Т.К. Винтцук. Анализ, распознавание и интерпретация речевых сигналов. Киев, Наукова думка, 1987.
- [34] В.Б. Кудрявцев, СВ. Алешин, А.С. Подколзин. Введение в теорию автоматов. М.: Наука, 1985. [35] Академическая программа Intel. <http://www.intel.ru/education/Grants/grants.htm>

- [36] R. Singh, Bh. Raj, R. M. Stern. Structured definition of sound units by merging and splitting for improved speech recognition. 1999.
- [37] T. Yoshimura, T. Masuko, K. Tokuda, T. Kobayashi, T. Kitamura. Speaker interpolation in HMM-based speech synthesis system. 1999.
- [38] B.T. Tan, Y. Gu, T. Thomas. Word confusability measures for vocabulary selection in speech recognition. 1999.
- [39] M. C Nechyba, Y. Xu. Stochastic similarity for validating human control strategy models. 1999.
- [40] S. Siruguri. Modelling a navigation task. 1999.
- [41] Н. В. Зиновьева, Л. М. Захаров, О. Ф. Кривнова, А. Ю. Фролов, И. Г. Фролова. Автоматический транскриптор. АРСО-92.
- [42] Л.М. Захаров. Транскрипция текстов при синтезе и анализе русской речи. АРСО-96.

## **Публикации автора по теме диссертации**

[a1]. Мазуренко И.Л. Одна модель распознавания речи. В сб.:

Компьютерные аспекты в научных исследованиях и учебном процессе.

-Издательство Московского университета, Москва, 1996 г., стр.107— 112.

[a2]. Мазуренко И.Л. О сокращении перебора в словаре речевых команд в составе системы распознавания речи. В сб.: Интеллектуальные системы, т.2, вып. 1-4, Москва, 1997. Стр. 135-148.

[a3]. Мазуренко И.Л. Компьютерные системы распознавания речи. В сб.: Интеллектуальные системы, т.3. вып. 1-2 - Москва, 1998 г. Стр. 117-134.

[a4]. Мазуренко И.Л. Многоканальная система распознавания речи. В сб.: VI всероссийская конференция "Нейрокомпьютеры и их применение" Сборник докладов. Москва 16-18 февраля 2000 г. Стр. 222—225.

[a5]. Бабин Д.Н., Дудецкий В.Н., Мазуренко И.Л., Уранцев А.В. и др.  
Устройство для синтеза и анализа речевых сигналов. Патент N  
94045004/09 (045160) решение от 20.12.96.

[aб]. Бабин Д.Н., Мазуренко И.Л., Уранцев А.В., Холоденко А.Б.Способ  
идентификации факта речевой активности оператора. Патент N  
99103468/ 28(003418) решение от 29.12.99. МІЖ 6 В 60 К 28/06, G 10 L  
5/06

В патентах [5], [6] основная заслуга автора заключается в разработке и программной реализации алгоритмов распознавания речевых сигналов.