

УДК 621.391:519.2.

МЕТОД АВТОМАТИЧЕСКОГО ВЫДЕЛЕНИЯ
УДАРЕНИЯ В ПОТОКЕ РЕЧИ

А.Г. Хайретдинова

При решении таких задач, как автоматическая сегментация речевого потока, нормализация речи по темпу и другие возникает необходимость автоматического выделения ударения в потоке речи.

Целью данной работы является исследование и разработка метода выделения ударных гласных в речевом потоке с помощью ЭВМ.

Многочисленные эксперименты, проведенные нами и другими исследователями речи, позволили подробно изучить фонетические особенности ударных и безударных гласных в акустическом сигнале и указать достаточно большое число их отличительных признаков.

Выбор признаков

При составлении машинного алгоритма из всей совокупности известных признаков были выбраны те, которые, на наш взгляд, являются наиболее устойчивыми и удобными при выделении их на ЭВМ. Рассмотрим эти признаки.

Признак χ_1 . Ударная гласная в отличие от остальных гласных внутри одного слова характеризуется наибольшей длительностью χ_1 . Это свойство справедливо как для слов, произнесенных изолированно, так и для слов, произнесенных во фразе. Лингвистическая частота появления этого признака для нормального темпа произношения приближается к 0.9. Однако определение границ гласных с помощью ЭВМ осуществляется с некоторыми ошибками, которые в конечном счете снижают эту цифру до 0.8-0.85.

Признак χ_2 . Ударная гласная имеет наибольшую интенсив-

ность χ_2 по сравнению с другими гласными фонемами внутри рассматриваемого слова. (здесь под интенсивностью понимается среднее значение амплитуд огибающей речевого сигнала). Это правило соблюдается почти всегда, если ударной является открытая гласная (а, о, э); если же ударной является закрытая гласная (и, у, ѿ), а среди безударных в слове имеются открытые гласные, то, как правило, такие ударные гласные по интенсивности уступают последним [1]. Поскольку в русской речи частота появления открытых ударных гласных значительно больше частоты появления закрытых ударных [2], то признак χ_2 является характерным признаком ударности.

Признак χ_3 . В повествовательном предложении внутри одного слова среднее значение частоты основного тона (ОТ) на участке ударной гласной имеет большую величину, чем на остальных участках. Сказанное можно проиллюстрировать таблицей I.*). Здесь через χ_3 обозначена разность между средними значениями ОТ двух соседних участков. В таблице приведены частоты появления ударных и безударных гласных по признаку χ_3 . Исключением являются слова, стоящие в конце фразы и произносимые обычно с понижением голоса. При этом замечено: если ударение падает на первый слог, то происходит плавный спад средней частоты ОТ к концу слова; если же ударным является другой из слогов (второй, третий), то перед ударной гласной спад средней частоты ОТ происходит скачкообразно (не менее 20 Гц).

Таблица I

гласные	$\chi_3 < -15$ Гц	-15 Гц $< \chi_3 < -5$ Гц	-5 Гц $< \chi_3 < 5$ Гц	5 Гц $< \chi_3 < 15$ Гц	$\chi_3 > 15$ Гц
ударные	0,1	0,11	0,1	0,19	0,52
безударные	0,34	0,26	0,3	0,064	0,04

Признак χ_4 . В повествовательных предложениях наблюдается повышение частоты ОТ в течение времени звучания ударной гласной, в то время как на безударных гласных чаще всего происходит спад частоты ОТ. Это свойство можно проиллюстрировать таблицей 2. Здесь χ_4 — величина, равная разности частот ОТ в

*). Таблицы I и 2 составлены на основе анализа гласных, встретившихся в 50 повествовательных фразах, содержащих 175 ударных и 274 безударных гласных (в произнесении одного диктора).

Таблица 2

гласные	$x_1 \leq -10 \text{ Гц}$	$-10 \text{ Гц} < x_1 < -2,5 \text{ Гц}$	$-2,5 \text{ Гц} < x_1 < -2,5 \text{ Гц}$	$2,5 \text{ Гц} \leq x_1 < 10 \text{ Гц}$	$x_1 > 10 \text{ Гц}$
ударные	0,06	0,07	0,14	0,21	0,49
без- ударные	0,28	0,28	0,25	0,08	0,09

конце и начале гласной.

В словах, стоящих в конце фразы, на ударной гласной имеет место спад частоты ОТ (например, из 90 конечных слов в 73 случаях x_1 , больше 15 Гц., в 17 случаях — меньше 15 Гц.).

В дальнейшем каждая гласная будет характеризоваться этими четырьмя признаками. Следует подчеркнуть, что признаки x_1 и x_2 характеризуют ударные гласные в словах, произнесенных как изолированно, так и во фразах. Признаки x_3 и x_4 , относятся лишь к словам, стоящим во фразе.

Описание алгоритма

В основе алгоритма для выделения ударных гласных в речевом сигнале лежат методы распознавания образов. В пространстве выбранных признаков каждая гласная фонема, описываемая совокупностью признаков, может быть интерпретирована как точка; совокупностям всех ударных и безударных гласных соответствуют две односвязанные области пространства. Предварительные эксперименты показали, что точки в каждой из областей расположены по закону, близакому к нормальному. Для разделения этих областей была использована гиперплоскость:

$$S = \sum_{i=1}^n w_i x_i + w_0 = 0,$$

где x_i — значение i -го признака,

w_i — вес i -го признака; w_0 — свободный член,

n — число признаков.

Значения w_i и w_0 определялись как коэффициенты уравнения гиперплоскости, перпендикулярной к линии, соединяющей математические ожидания обоих областей и делящей эту линию на части таким образом, что проекции её на координатные оси пропорциональны средним квадратичным отклонениям соответствующих признаков этих образов [4].

При определении коэффициентов гиперплоскости использова-

лись точки обучающего материала. При принятии решения знак S указывает, по какую сторону от гиперплоскости находится контрольная реализация.

Определение признаков на ЭВМ

При вычислении значений выбранных признаков на ЭВМ требуется в первую очередь определить границы всех гласных фонем. При этом знание вида гласной не обязательно.

С целью определения границ гласных речевой сигнал пропускался через полосовой фильтр с граничными частотами 250 Гц и 1000 Гц (был использован рекурсивный цифровой фильтр, аппроксимированный по Баттерворту [5]). После фильтра осуществлялось детектирование сигнала и нахождение его огибающей (с интервалом дискретности 15 мсек.), которая представляет собой чередование "горбов" и "нулевых зон" (т.е. участков, где огибающая равна нулю). "Нулевые зоны" соответствуют высокочастотным фонемам (шипящие, взрывные и т.д.), спектральные составляющие которых в указанной полосе частот имеют никакий уровень. "Горбы" в большинстве случаев появляются там, где находятся гласные. На выходе фильтра они представлены первой или двумя первыми формантами, т.е. основной частью своей энергии. Наряду с гласными в этом сигнале присутствуют сонанты (I форманта). Длительность "горба" на уровне 0,5 от его пиковой амплитуды принималась за длительность гласной.

На практике встречаются случаи, когда несколько соседних фонем оказываются слитыми в один "горб". Сюда относится в первую очередь дифтонги (например, "ие", "ые"), сочетания некоторых гласных с сонантами (например, слог "умн" в слове "умножить"). Следует сказать, что даже привлечение таких тонких характеристик фонем как спектральные, не всегда помогает членению этих сочетаний.

В данной работе границы между фонемами в упомянутых случаях определялись приближенно, исходя из следующих соображений: если среди "горбов" имеется такой, длительность $\bar{\tau}$ которого существенно отличается от средней длительности $\bar{\tau}_f$ ближайших соседних "горбов" (например, находящихся на участке длиной в одно слово), то в таком "горбе" обычно имеет место слияние нескольких фонем. В большинстве случаев отношение $\bar{\tau}/\bar{\tau}_f$ указывает число фонем, входящих в это сочетание. При этом следует

учесть, что по явойе длительности сонанты близки к безударным гласным [3].

Для рассмотренных нами двух дикторов оказалось: при $\kappa < 2$ имеет место одна фонема, при $2 < \kappa < 3$ имеют место две фонемы, при $\kappa > 3$ имеет место три фонемы. Как показали дальнейшие эксперименты, такой подход к определению границ гласных является пригодным.

Итак, из речевого сигнала были выделены участки (фонемы), каждый из которых представляет собой либо гласную, либо сонанту.

В дальнейшем для каждого из участков были вычислены признаки x_1, x_2, x_3, x_4 .

При определении x_1 предварительно подсчитывались текущие значения частоты ОТ. Для этого был использован сдвиговый метод определения ОТ [6] при интервале дискретности 15 мсек.

Нормирование признаков осуществлялось по следующему правилу:

- значения x_1 и x_2 относились к соответствующим максимальным значениям на участке фразы длиной в одно слово (для нормального темпа средняя длина слова около 500 мсек);
- для признака x_3 определялись средние значения соответственно по всем его положительным и отрицательным значениям;
- значения x_4 в зависимости от его знака относились к модулю соответствующих средних значений.

Признак x_1 нормировался аналогично x_3 .

Для фонем, расположенных в интервале 500 мсек. (длина слова) от конца фразы, признаки x_1 и x_2 были взяты с обратным знаком.

На ЭВМ "БЭСМ-6" по описанному выше алгоритму выделения ударения были обработаны 50 повествовательных фраз, составленных из наиболее употребительных слов. Среднее число слов во фразах составляло 5-6. Фразы произносились двумя дикторами (мужчиной и женщиной) в нормальном темпе. Общее количество гласных во фразах равнялось 898, из них 350 ударных и 548 безударных.

Обучение проводилось на ударных и безударных гласных, взятых из 20 фраз, после чего распознавались гласные, взятые из 30 новых фраз.

* Предварительные эксперименты показали, что по указанным признакам сонанты близки к безударным гласным.

Суммарная ошибка подсчитывалась по формуле:

$$\Delta = p_1 \alpha_1 + p_2 \alpha_2,$$

где α_1 - ошибка отнесения ударных гласных к безударным, α_2 - ошибка отнесения безударных гласных к ударным, p_1 и p_2 - априорные вероятности появления соответственно ударных и безударных гласных.

Так как в русской речи наиболее употребительными являются двух- и трехсложные слова [7], то было принято $p_1 = 0,35$, $p_2 = 0,65$. Суммарная ошибка получилась равной 0,08.

Результаты эксперимента показали достаточно высокую эффективность выбранной системы признаков при выделении ударной гласной из потока речи.

Предложенный алгоритм принятия решения позволяет достаточно просто и надежно (с вероятностью 0,92) решать рассматриваемую задачу.

Л и т е р а т у р а

1. Хайретдинова А.Г. Исследование свойств ударения. Отчет. Институт математики СО АН СССР, 1969.
2. Ильина В.Н., Юдина Л.С. Статистика открытых слогов русской речи. Сб. "Вычислительные системы", Новосибирск, вып. 14, 1964.
3. Дифференциальные признаки слогов. Отчет, ДГУ, 1967.
4. Загоруйко Н.Г. Структура проблемы распознавания слуховых образов и методы её решения. В кн.: Автоматическое распознавание слуховых образов. Новосибирск, "Наука", 1966.
5. Лозовский В.С. Программа синтеза рекурсивных цифровых фильтров (Р.Ц.Ф.). Тр. ИМ СО АН СССР. "Вычислительные системы", Новосибирск, вып. 36, 1969.
6. Соболев В.Н., Баронин С.Л. Исследование сдвигового метода выделения основного тона речи. "Электросвязь", № 12, 1968.
7. Сапожков И.А. Речевой сигнал в кибернетике и связи. М., Связьиздат, 1963..

Поступила в редакцию
7.1.1971 г.