

УДК 519.226: 534.44

АЛГОРИТМ КЛАССИФИКАЦИИ ТОН/ШУМ ПО ЧАСТНЫМ АВТОКОРРЕЛЯЦИЯМ

А.В.Кельманов

Пусть речевой сигнал $s(t)$ преквантован с частотой $F_s = 1/T$ так, что $s_n = s(nT)$, где T - интервал квантования. Разобьем сигнал на L сегментов $\{s_n\}^L_{n=1}$, $n=1, N$, $1=1, L$ длительностью $T_a = (n-1)/F_s$ со сдвигом от сегмента к сегменту на интервале ΔT . Задача - определить, к какому типу звуков (голосовых или неголосовых) относится каждый из L сегментов.

Для решения этой задачи уже построены [1] два параметрических алгоритма, использующие вероятностные распределения. В настоящей работе предлагается еще один алгоритм подобного типа и проводится его сравнение с вышеуказанными.

I. Алгоритм решения

В основе предлагаемого алгоритма (назовем его алгоритм I), как и в алгоритмах (II и III) из [1], лежит статистический критерий диагностической проверки на адекватность описания речевого сигнала моделью авторегрессии. Воспользоваться этим критерием для автоматической классификации тон/шум позволяет тот факт, что голосовые звуки в отличие от неголосовых описываются моделью авторегрессии неадекватно [1].

Предположим, что каждый из сегментов речевого сигнала аппроксимирован моделью авторегрессии K -го порядка:

$$s_n = \sum_{i=1}^K a_i s_{n-i} + \alpha_n, \quad n = \overline{1, N}, \quad (I)$$

где $\{\alpha_n\}$ - остаточный ряд ошибок, распределенных нормально с нулевым средним и дисперсией σ_α^2 ; и найдены оценки Шла-Уакера [2] параметров авторегрессии и частных автокорреляций:

$$\hat{a}_{j+1,i} = \hat{a}_{j,i} - \hat{a}_{j+1,j+1} \times \hat{a}_{j,j-i+1}, \quad i=1,2,\dots,j;$$

$$\hat{a}_{j+1,j+1} = \frac{\hat{r}_{j+1} - \sum_{i=1}^j \hat{a}_{ji} \hat{r}_{j-i+1}}{1 - \sum_{i=1}^j \hat{a}_{ji} \hat{r}_i}, \quad j=0,1,\dots,K-1.$$

Если порядок процесса авторегрессии равен p , то для всех $j > p$ выборочные частные автокорреляции \hat{a}_{jj} : 1) распределены приблизительно нормально с нулевым средним даже при небольших объемах выборки N (см. [2]); 2) имеют дисперсию $\sigma^2[\hat{a}_{jj}] = N^{-1}$ (см. [2]) и 3) статистически независимы [2]. Отсюда нетрудно показать, что статистика

$$\Theta = N \sum_{i=p+1}^K \hat{a}_{ii}^2 \quad (2)$$

распределена как $\chi^2(K-p)$. Равенство (2) можно использовать для проверки модели авторегрессии на адекватность: если модель соответствует данным, то Θ распределена как $\chi^2(K-p)$, если нет, то среднее значение Θ увеличивается. Иными словами, необходимо проверить нулевую гипотезу $H_0: \forall i (i > p \Rightarrow a_{ii} = 0)$ против альтернативной $H_1: \exists i (i > p \Rightarrow a_{ii} \neq 0)$. Для заданного уровня значимости γ можно найти порог χ^2_γ такой, что $P(\Theta > \chi^2_\gamma | H_0) = \gamma$. Поэтому если $\Theta > \chi^2_\gamma$, то построенная модель неадекватна, а соответствующий сегмент голосовой.

Напомним, что в алгоритме II решение принимается также по критерию χ^2 , но статистика построена по автокорреляциям ряда остаточных ошибок $\{\alpha_n\}$, а в алгоритме I решение принимается по одному значению выборочной частной автокорреляции при помощи двухстороннего критерия для нормально распределенных наблюдений.

Для сокращения объема вычислений, как и в алгоритмах II и III, на первом шаге в алгоритме I предусмотрена никакочастотная фильтрация, осуществляемая при помощи чебышевского фильтра 3-го порядка с частотой среза $F_c = 0,8$ кГц, понижение исходной частоты квантования сигнала F_s до 2 кГц, а также извлечение профильтрованного сигнала окном Хэмминга перед авторегрессионной аппроксимацией (I). Блок принятия решения и коррекции ошибок классификации идентичен блоку, включенному в алгоритмы II и III.

Выбор "оптимального" [5] порядка модели p производится из соображений минимума ошибки аппроксимации и минимума объема вычислений с учетом формальной структуры речевого сигнала.

Длина ряда частных автокорреляций K выбирается с учетом того, что при уменьшении K возрастает объем вычислений, однако при этом уменьшается разность $K - p$, равная числу частных автокорреляций, по которым принимается решение, и поэтому увеличивается вероятность ошибочной классификации. Наоборот, при увеличении K вероятность ошибок уменьшается, но возрастает объем вычислений.

2. Трудоемкость алгоритма

Трудоемкость алгоритма I (без учета низкочастотной фильтрации) состоит из трудоемкости оценивания K частных автокорреляций, равной $(K+1) + K^2$ операций (сложения и умножения), трудоемкости вычисления статистики θ , равной $K - p$ операций, и трудоемкости принятия решения и коррекции ошибок, равной 5 операциям.

Т а б л и ц а I

В табл. I приведены трудоемкости алгоритмов I-II.

После низкочастотной фильтрации, которая требует примерно 5 и операций, объем выборки N для каждого сегмента сокращается до величины:

$$N_1 = \frac{2(N-1)}{F_s} + 1 \approx \frac{2N}{F_s}.$$

При этом, вследствие уменьшения частотного диапазона, уменьшается порядок модели p , а стало быть, уменьшается и K (предполагается, что разность $K - p$ остается такой же, как и при исходной частоте квантования). Для подсчета трудоемкости алгоритмов с учетом фильтрации в каждой строке табл. I следует заменить N на $2N/F_s$ и прибавить 5 и .

3. Эксперименты и их результаты

Речевой материал. При проведении эксперимента использовались 2 словаря. Первый состоял из 100 русских слов [3], второй — из 63 слов [4]. По первому словарю одним диктором-мужчиной (диктор I) были сформированы две акустические последовательности

тельности, а по второму – пятью дикторами-мужчинами (дикторы I, 2, 3, 4, 5) было сформировано по одной последовательности. Общее число сегментов, на которых принималось решение тон/шум, составило величину, равную примерно 24,5 тысячи.

Условия эксперимента. Эксперимент проводился на ЭВМ "Минск-32". Ввод слов осуществлялся в машинном зале (уровень шума ~ 60 дБ) через микрофон типа МД-59 и семиразрядный аналогово-цифровой преобразователь с частотой квантования $F_s = 20$ кГц.

Условия анализа. Акустические реализации были обработаны при помощи алгоритма I, описанного в данной статье, а также при помощи алгоритмов II и III. Разбиение слов на сегменты длительностью $T_a = 25,6$ мсек производилось со сдвигом (с перекрытием) от сегмента к сегменту на время $\Delta T = 16$ мсек так, что $N = 512$, а $N_1 = 52$. Значения параметров K и r равнялись 30 и 4 соответственно. Поэтому статистика (2) распределена как $\chi^2(26)$, и при уровне значимости $\gamma = 0,001$ значение порога $\chi^2_{\gamma} = 54,1$.

Результаты классификации. Все акустические реализации слов были выведены на печать в виде картинок видимой речи с целью проведения визуальной классификации. После этого те же реализации были обработаны алгоритмически. Для примера на рис. I приведен график изменения статистики χ^2 во времени и результаты классификации слова "мифический" по алгоритму I. В табл. 2 приведены данные по классификации сегментов для отдельных видов звуков при помощи трех алгоритмов, а в табл. 3 – результаты классификации тональных и шумовых сегментов.

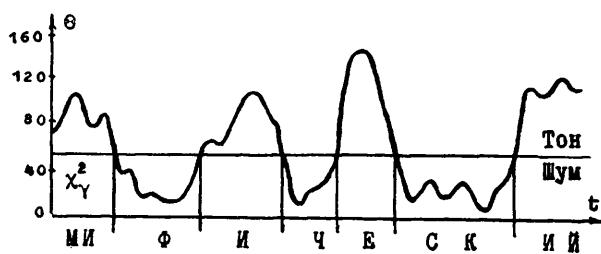


Рис. I

Таблица 2

Ошибки классификации сегментов для отдельных видов звуков

Фонемы	Число реализаций	Алгоритм I		Алгоритм II		Алгоритм III	
		число ошибок	% ошибок	число ошибок	% ошибок	число ошибок	% ошибок
а, о, у, э, и, и	10069	31	0,60	185	1,84	640	6,36
м, н, л, р, в	4695	64	2,68	148	3,15	191	4,07
б, д, г	2145	164	7,64	161	7,50	133	6,20
з, ж	1511	58	3,84	84	5,56	25	1,65
п, т, х	2465	5	0,20	25	1,01	144	5,84
ц, ч	717	0	0	6	0,84	37	5,16
ф, с, ш, х	2882	2	0,07	6	0,21	235	8,15
Всего:	24485	324	1,32	615	2,51	1405	5,74

Таблица 3

Ошибки классификации тоновых и шумовых сегментов

Тип сегмента	Число реализаций	Алгоритм I		Алгоритм II		Алгоритм III	
		число ошибок	% ошибок	число ошибок	% ошибок	число ошибок	% ошибок
Тон	18421	317	1,72	578	3,14	989	5,37
Шум	6064	7	0,12	37	0,61	416	6,86
Всего:	24485	324	1,32	615	2,51	1405	5,74

4. Обсуждение результатов

Сравнение алгоритмов по трудоемкости. Обращаясь к табл.1, прежде всего отметим, что наименее трудоемкий - алгоритм III. Далее, учитывая, что K должно быть больше r , замечаем, что, начиная с $K = r+1$ до $K = \sqrt{2N(r+1)} + r^2 + r$, трудоемкость алгоритма I меньше трудоемкости алгоритма II, при дальнейшем увеличении K алгоритм I становится более трудоемким. Не трудно убедиться в том, что при $N = 52$, $K = 30$ и $r = 4$ алго-

рите I требует больше времени на принятие решения, чем алгоритм II. Для примера приведем трудоемкости алгоритмов с учетом низкочастотной фильтрации при исходной частоте квантования $F_s = 20$ кГц и при $K = 30$ и $r = 4$. При этих значениях параметров трудоемкости алгоритмов I, II и III приближенно равны 10И_h, 9И_h и 5,5И_h соответственно.

Сравнение алгоритмов по надежности. Сравнивая алгоритмы по табл.2, нетрудно заметить, что для неголосовых групп звуков: глухих взрывных (и, т, к), аффрикат (ч, ч), глухих целях (ф, с, х, ш) – наименьшую вероятность ошибочной классификации дает алгоритм I, затем – алгоритм II и, наконец, – алгоритм III. Аналогичная картина наблюдается для голосовых групп звуков: гласных (а, о, у, э, и, я) и сонорных (м, н, л, р, в). На звонких взрывных (б, д, г) наилучшие результаты дает алгоритм III, а наихудшие – алгоритм I, однако отличия в вероятностях ошибки незначительны. Звонкие целях (з, ж) наилучшим образом классифицируются алгоритмом III, а наихудшим – алгоритмом II. Из табл.3 следует, что как по отдельности, так и в целом на голосовых и неголосовых звуках наилучшие результаты показывает алгоритм I. Для алгоритма II вероятности ошибочной классификации (обеих родов) возрастают. Наконец, алгоритм III дает наибольшие вероятности ошибочной классификации.

Таким образом, наиболее надежным является наиболее трудоемкий алгоритм I, а наименее надежным – наименее трудоемкий алгоритм III.

5. Выводы и рекомендации

Результаты экспериментального тестирования позволяют сделать следующие выводы и рекомендации.

1. Предложенный алгоритм автоматической классификации звуков речи на звонкие и глухие спрообразован на пяти дикторах и показал высокую надежность. Ошибки классификации тоновых и пульсовых сегментов равны 1,72% и 0,12% соответственно.

2. Алгоритм аналогичного типа, с которым проводилось сравнение, менее надежный, однако позволяет производить классификацию за меньшее время.

3. Окончательный выбор того или иного алгоритма классификации тон/пульс рекомендуется проводить, исходя из конкретных условий работы речевой системы, в которой задействован классификатор по-

добного типа, и учитывая претворенные требования уменьшения вероятности ошибочной классификации и повышения трудоемкости.

В заключение отметим, что в настоящее время описанный алгоритм успешно используется в системе распознавания изолированных слов.

Л и т е р а т у р а

1. КЕЛЬМАНОВ А.В. Алгоритм классификации тона/музыки, основанный на критерии адекватности модели авторегрессии. - В кн.: Методы обработки информации. (Вычислительные системы, вып.74.) Новосибирск, 1978, с.129-148.
2. БОЖС Дж.: ДЕННИС Г. Анализ временных рядов, прогноз и управление. - М.: Мир, т.1, 1974. - 206 с.
3. КЕЛЬМАНОВ А.В. Система распознавания изолированных слов по частной автокорреляционной функции. - В кн.: Эмпирическое предсказание и распознавание образов. (Вычислительные системы, вып.76.) Новосибирск, 1978, с.132-143.
4. ДЕБЕДЕВ В.Г. Автоматическое распознавание речи по маск-признакам. - В кн.: Эмпирическое предсказание и распознавание образов. (Вычислительные системы, вып.67.) Новосибирск, 1976, с.136-140.
5. MARKEL J.D. Application of digital inverse filter for automatic formant and F_0 analysis.- IEEE Trans. Audio Electroacoust, 1973, v.AU-21, N 3, p.154-160.

Поступила в ред.-изд. отд.
20 октября 1979 года