

ПОДСТРОЙКА ПОД ДИКТОРА ПРИ РАСПОЗНАВАНИИ
ОГРАНИЧЕННОГО НАБОРА УСТНЫХ КОМАНД

Н.Г.Загоруйко, В.С.Лозовский

Для всех известных в настоящее время алгоритмов распознавания ограниченного набора устных команд характерно более или менее значительное ухудшение надежности распознавания при росте числа дикторов, участвующих в эксперименте. Еще худшие результаты получаются при распознавании слов, произнесенных дикторами, не участвовавшими в формировании обучающей последовательности. С другой стороны, на практике встречаются случаи, когда допустима подстройка под диктора перед распознаванием произносимых им команд, если этот процесс будет простым и быстрым. Если бы в этом направлении были получены хорошие результаты, то для ряда применений распознающих автоматов можно было бы отказаться от универсальности, подстраивая параметры автомата при непродолжительном предварительном знакомстве с диктором. Рассмотрению такой возможности и посвящена эта работа.

Пусть $X = \{X_j\}$, где $j = 1, \dots, n$, n - мерное пространство параметров распознавания. Произносимые команды отображаются в точки пространства X . Будем называть их реализациями. Каждую из реализаций учитель относит к одному из K классов S_1, \dots, S_K . В настоящей работе исследовалась подстройка под диктора для алгоритма распознавания по миниму-

му евклидова расстояния до эталонов.

Под надежностью распознавания контрольной последовательности реализаций будем понимать отношение правильно распознанных к их общему числу.

Пусть τ_j ($j=1, \dots, \ell$) - обучающая последовательность для класса S_i : $\tau_j \in S_i$. Если экспериментатору известны априорные вероятности каждой реализации P_j и нормированные веса W_j , учитывающие типичность каждой реализации, то эталонный вектор определяется из соотношения:

$$\mathcal{E}_i = \frac{1}{\ell} \sum_{j=1}^{\ell} W_j P_j \tau_j$$

Здесь мы не будем учитывать W_j и P_j , так что

$$\mathcal{E}_i = \frac{1}{\ell} \sum_{j=1}^{\ell} \tau_j. \quad (I)$$

Будем считать, что нами получены $\mathcal{E}_i^{(1)} \in X$ ($i=1, \dots, k$) по обучающей последовательности одного диктора и $\mathcal{E}_i^{(2)} \in X$ ($i=1, \dots, k$) - по обучающей последовательности второго. На рис. I этот случай изображен для $n=2$, $k=3$.

Расстояния между одноименными эталонами этих дикторов выражаются n -мерными векторами $\mathcal{E}_i^{(2)} \mathcal{E}_i^{(1)} = \mathcal{E}_i^{(2)} - \mathcal{E}_i^{(1)}$.

Далее, мы вправе написать:

$$\mathcal{E}_i^{(2)} \mathcal{E}_i^{(1)} = C_i \mathcal{E}_i^{(1)} + \mathcal{E}_i^{(2)} C_i, \text{ где, } C_1 \mathcal{E}_1^{(1)} = C_2 \mathcal{E}_2^{(1)} = \dots = C\mathcal{E},$$

где

$$\mathcal{E}_i^{(2)} \mathcal{E}_i^{(1)} = C\mathcal{E} + \mathcal{E}_i^{(2)} C_i.$$

Здесь постоянный для всех классов вектор смещения $C\mathcal{E}$ можно трактовать как поправку, учитывающую глобальную разницу между дикторами: высоту голоса, особенности произношения, темпа речи и т.п. Векторы же $\mathcal{E}_i^{(2)} C_i$ характеризуют имеющиеся локальные отклонения, обусловленные особенностями произношения вторым диктором конкретных команд из данного набора.

Итак, высказывается следующая гипотеза об относительной стабильности положения эталонов, в которой предполагается, что

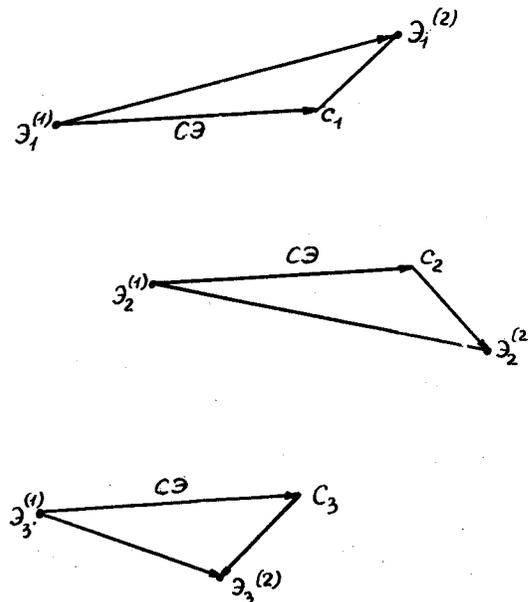


Рис. I

для набора эталонов одних и тех же команд, произнесенных двумя дикторами, в пространстве X существует такой вектор

смещения $CЭ$, что $|\mathcal{E}_i^{(2)} C_i| < |CЭ|$ для $i=1, \dots, k$. Другими словами, предполагается, что глобальные различия между дикторами могут быть велики, но взаимные расположения эталонов для каждого диктора в пространстве X будут различаться не столь сильно.

По определению вектора смещения $CЭ$, последний должен определяться при минимизации некоторого функционала F от локальных отклонений $\mathcal{E}_i^{(2)} C_i$. Учитывая принятую нами процедуру распознавания по минимуму евклидова расстояния до эталонов, имеем:

$$F_1 = \sum_{i=1}^k |\mathcal{E}_i^{(2)} - (\mathcal{E}_i^{(1)} + CЭ)|^2.$$

В данной работе с целью экономии времени вычислений процедура минимизации F была заменена процедурой непосредственного счета для более простого функционала:

$$F_2 = \sum_{i=1}^k [\mathcal{E}_i^{(2)} - (\mathcal{E}_i^{(1)} + CЭ)]^2. \quad (2)$$

Полагая (2) равным нулю и осуществляя элементарные преобразования, имеем:

$$CЭ = \frac{1}{2} \sum_{i=1}^k \mathcal{E}_i^{(2)} - \frac{1}{2} \sum_{i=1}^k \mathcal{E}_i^{(1)} = \mathcal{E}^{(2)} - \mathcal{E}^{(1)}, \quad (3)$$

где $\mathcal{E}^{(1)}$ и $\mathcal{E}^{(2)}$ - центры распределений эталонов первого и второго дикторов.

Гипотеза об относительной стабильности положения эталонов получит подтверждение, если надежность распознавания контрольной последовательности одного диктора по эталонам другого с учетом вектора смещения $CЭ$ окажется выше, чем без учета последнего.

Эксперимент заключался в распознавании десяти классов - произносимых цифр от "0" до "9" в пространстве огибающих на

выходе пяти октавных полосных фильтров. Аппаратурно были выполнены: пятиканальный полосный анализатор, коммутатор аналогового сигнала и цифратор (преобразователь напряжение-код). Остальные операции выполнялись программно на ЭЦВМ М-20. Опрос пяти каналов коммутатором происходил примерно за 2 мсек, отсчеты снимались каждые 20 мсек, точность цифратора - 6 дв. разрядов. Программным методом находилось начало и конец каждого слова, затем производилась линейная нормализация по динамическому диапазону (максимальный отсчет за время звучания слова вытягивался до фиксированного уровня) и линейная нормализация по длительности (до 50 отсчетов). Распознавание велось в пространстве 250 измерений. Дикторов двое: мужчина с предварительной тренировкой в произнесении слов и женщина без особой подготовки. Дикторам было предложено произнести каждое из 10 слов по 40 раз в нормальном темпе и полным стилем.

Положение эталонов $\mathcal{E}_i^{(1)}$ и $\mathcal{E}_i^{(2)}$ определялось в соответствии с (1). Поскольку цель эксперимента - лишь исследование возможности подстройки под диктора, в качестве контрольной последовательности использовалась обучающая. Вектор смещения эталонов $CЭ$ определялся по формуле (3). В табл. I приведены результаты распознавания слов, произнесенных первым диктором, а в табл. 2 - вторым. На пересечении каждой строки и столбца таблиц помещены четыре числа (нули опущены):

1. надежность распознавания в процентах реализаций данного диктора по эталонам первого диктора;
2. то же, но по эталонам второго диктора;
3. то же, но по эталонам первого диктора плюс вектор смещения;
4. то же, по эталонам второго диктора минус вектор смещения.

Пример: Слово "ноль", произнесенное первым диктором, распознано, как "ноль" по эталонам первого диктора с надежностью 100%, по эталонам второго - 26%. Смещение эталона I-го диктора снизило надежность распознавания его реализаций до 87% и улучшило распознавание реализаций второго диктора до 66%. Отдельно на каждой таблице указаны четыре числа, характеризующие надежность распознавания для каждого из четырех случаев, усредненную по десяти словам.

Из приведенных данных следует, что учет векторов смещения повысил надежность распознавания реализаций первого дикто-

Таблица I.

Результаты распознаваемых команд, произносимых первым диктором.

		РЕЗУЛЬТАТЫ РАСПОЗНАВАНИЯ									
		ноль	один	два	три	четыре	пять	шесть	семь	восемь	девять
ПРОИЗНЕСЕНО ПЕРВЫМ ДИКТОРОМ	ноль	100 26		66			5	3			
	один	87 66		13 32			2				
	два		100 50		48	2					
	три	8 3		92			76	21			
	четыре	3 16		91 5		3	65	14			3
	пять				100 37		5	27	6	25	
	шесть					100 40		8	32	7	
	семь						100 65			35	
	восемь							100 73		27	
	девять								100 35		48
ПРОИЗНЕСЕНО ВТОРЫМ ДИКТОРОМ	ноль										
	один										
	два										
	три										
	четыре										
	пять										
	шесть										
	семь										
	восемь										
	девять										

98,9 41
97,8 50,4

Таблица 2.

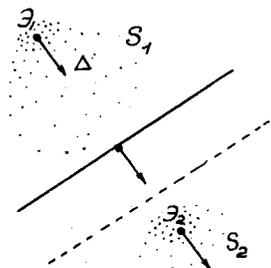
Результаты распознаваемых команд, произносимых вторым диктором.

		РЕЗУЛЬТАТ РАСПОЗНАВАНИЯ									
		ноль	один	два	три	четыре	пять	шесть	семь	восемь	девять
ПРОИЗНЕСЕНО ПЕРВЫМ ДИКТОРОМ	ноль	100 70		25			3			2	
	один	90 89		5 10		5				2	
	два		79 97				11 3			10	
	три		62 97				15 3			18 5	
	четыре			95 8		5 77		2 7		6	
	пять			67 30		80 65			8	5	5
	шесть				89 5 11		95				
	семь					81 8 6	8 92		5		
	восемь						80		20 100		
	девять						18		80 100		2
ПРОИЗНЕСЕНО ВТОРЫМ ДИКТОРОМ	ноль										
	один										
	два										
	три										
	четыре										
	пять										
	шесть										
	семь										
	восемь										
	девять										

45,6 89,1
53 89,9

ра по эталонам второго в среднем на 9,4%, а надежность распознавания реализаций второго диктора по эталонам первого — на 7,4%. Можно заметить, что учет вектора смещения "портит" свои эталоны незначительно: надежность распознавания "своих" реализаций для первого диктора ухудшилась лишь на 1,1%, а для второго даже улучшилась на 0,8%.

Объяснить замеченное улучшение можно несоответствием вида решающих функций характеру распределений. Представим себе, что плотность распределения реализаций для некоторых классов S_1 и S_2 неравномерна и имеет вид, изображенный на рис.2. Вычисление эталонов \mathcal{E}_1 и \mathcal{E}_2 по формуле (I) даст указанные на рисунке смещения их в сторону области с большими плотностями реализаций. Учет для каждого эталона некоторой векторной поправки Δ , как видно из рисунка, приведет к смещению разделяющей границы и повысит надежность распознавания



реализаций обоих классов. Поэтому для указанного случая может оказаться целесообразным помещать эталоны \mathcal{E}_i не в центры тяжести реализаций каждого класса, а в центры занимаемых ими объемов.

Заканчивая обсуждение результатов, следует отметить, что гипотеза относительной стабильности эталонов для проведенного эксперимента подтвердилась, хотя, может быть, и не в такой степени, как того бы хотелось. Есть надежда, что еще более эффективные результаты могут быть получены для других параметрических пространств и для более совершенных алгоритмов классификации.

Несколько слов о практическом использовании результатов. Пусть в распознающем автомате хранятся координаты эталонов для наиболее типичного произношения, и пусть в его памяти хранятся векторы смещения СЭ для каждого из дикторов, реализации которых будут входить в состав контрольной последовательности. В каждом конкретном случае, в зависимости от того с каким диктором имеет дело автомат, учитывается нужный вектор смещения. Если диктор неизвестен, то распознавание ведется с учетом каждого из имеющихся векторов смещения и состоит в вы-

боре того вектора смещения и того класса, при которых расстояние до поступившей реализации окажется минимальным. Так может выглядеть один из вариантов одновременного распознавания смысла сообщения и диктора.

Далее. Пусть алфавит распознавания велик. Представляется возможным определить вектор смещения СЭ не по всем классам, а лишь по тем, которых достаточно для нахождения СЭ с нужной степенью точности. Для этого необходимо провести дополнительную работу. На множестве из нескольких дикторов для данного алфавита распознавания вычисляются \mathcal{E}^i центры распределений эталонов для каждого диктора. Затем для каждого диктора выбираются эталоны, наиболее близкие к \mathcal{E}^i . Если выбранными окажутся одноименные эталоны, — это будет означать, что найдены слова, которые достаточно полно отражают индивидуальные особенности каждого диктора. В дальнейшем вектор смещения эталонов можно искать, "переучивая" автомат не на всем словаре, а только на этой его небольшой части.

Поступила в редакцию
I/X-1967г.