

## 基于高斯混合模型的说话人确认系统

杨澄宇, 赵文, 杨 鉴

(云南大学 信息与电子科学系, 云南 昆明 650091)

**摘 要:** 由于在人的语音频谱中, 低频和较高频段含有较多说话人的个性信息, 本文提出一种 LPC 倒谱的改进算法用于与文本无关的说话人识别。该改进算法通过语音频谱的各频段进行加权, 突出说话人的个性信息, 从而使说话人更易于区分。

**关键词:** 说话人识别; LPC 谱加权倒谱; 混合高斯模型

**中图分类号:** TP391.42 **文献标识码:** A

## SPEAKER IDENTIFICATION SYSTEM BASED ON GAUSSIAN MIXTURE MODELING

YANG Cheng-yu, ZHAO Wen, YANG Jian

(Department of Information and Electronics, Yunnan University, Kunming Yunnan 650091, China)

**Abstract:** We proceed a text-independent speaker identification system using Gaussian Mixture modeling. Since the lower and higher portion of the spectrum contain more reliable idiosyncrasy information on speaker than the middle portion, this paper presents a novel LPC cepstral algorithm for speaker identification task. By weighting the different portion in frequency domain, it can makes idiosyncrasy information be prominent and speaker recognition be easier.

**Key words:** speaker identification; LPC spectrum weighting cepstral; GMM

### 1 引言

说话人识别研究是传统语音识别的一个分支, 可用于安全系统、多媒体数据库检索等领域。它又可分为说话人辨识 (verification) 和说话人确认 (identification) 两种不同的处理方法。给定一个说话人集合  $S$  和未知说话人  $X$  的一段语音, 说话人辨识是指先判断  $X$  是否是  $S$  中的元素, 如果是的话, 再判断  $X$  是哪个元素。如果预先已知  $X$  属于  $S$ , 只需判断  $X$  是  $S$  中的哪个人, 则称为说话人确认。另外, 有的说话人识别系统通过规定训练和测试文本的内容来提高识别率, 称为以文本有关的说话人识别。如果对文本内容不做要求, 则称为以文本无关的说话人识别系统<sup>[1]</sup>。本文讨论的是与文本无关的说话人识别。

说话人识别系统大致包括两个主要部分, 即语音特征参数提取部分和识别算法部分。本文用改进的 LPC 倒谱作特征参数, 用概率模型—混合高斯模型 (GMM) 作识别算法。

混合高斯模型是一个状态的 CHMM, 该模型用多个高斯分布的概率密度函数的组合来描述矢量在概率空间的分布情况。在说话人识别系统中, 用混合高斯模型的参数描述说话人语音的特征矢量的概率分布。不同说话人将对应不同的混合高斯模型参数。

常用于说话人识别的语音特征参数主要有 LPC 倒谱和 Mel 倒谱, 两者都属倒谱系数。在所有声道的音段参数中, 基

于 LPC 的倒谱系数最能有效反映说话人特征<sup>[2]</sup>, 本文将使用 LPC 倒谱。由于在人的语音频谱中, 低频段和高频段含有较多的说话人个性信息<sup>[3]</sup>, 所以本文考虑在求 LPC 倒谱的过程中, 对频谱各段加权, 获得“LPC 谱加权的倒谱”参数, 突出说话人的个性特征。

### 2 LPC 谱加权的倒谱

#### 2.1 LPC 倒谱

(1) LPC 是线性预测分析的简写。按照线性预测分析, 语音信号在时刻  $n$  的样点值  $S(n)$ , 可以由过去的  $P$  个样点值的线性组合来近似。即

$$S(n) \approx a_1 S(n-1) + a_2 S(n-2) + \cdots + a_p S(n-p)$$

上式可看成由过去的样本值的线性组合来预测当前样本的值。这里  $a_1, a_2, a_3, \dots, a_p$  是常数, 称为 LPC 系数。按照语音信号产生模型, 声道、声门激励及辐射的全部谱效应简化为一个时变的数字滤波器。其稳态系统函数为:

$$H(z) = \frac{S(z)}{U(z)} = \frac{G}{1 - \sum_{i=1}^p a_i z^{-i}} \quad (1)$$

上式中,  $S(z)$  为系统输出,  $U(z)$  为系统的输入。清音的输入是随机噪声, 浊音的输入是一定频率的脉冲, 即声门波。  $G$  为一常数,  $a_i$  ( $i=1, 2, \dots, p$ ) 为 LPC 系数。上式又称为 LPC 全极点模型。对于具有短时平稳性的语音段, 该模型能对声道谱的

包络提供很好的近似。

(2) 倒谱是对一帧短时语音进行同态解卷得到的一组时域值。因为一个线性时不变系统的激励信号和系统的冲击响应是按卷积方式结合起来的,语音信号作为线性时不变系统的输出,通过对其进行解卷积处理,使语音中包含的激励源和声道冲击响应分离开,从而更清晰表示出语音包含的各种特征。同态解卷是解卷中的一种<sup>[4]</sup>,其过程如图所示:

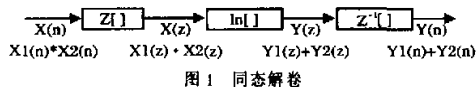


图1 同态解卷

由图1可见,对具有卷积关系的信号  $X_1(n)$  和  $X_2(n)$ ,经同态解卷过程后,对应的时域输出信号  $Y_1(n)$ ,  $Y_2(n)$  之间是相加的关系。 $Y(n) = Y_1(n) + Y_2(n)$  即称为  $X(n)$  的倒谱。

(3) LPC 倒谱是指由 LPC 系数推导出的倒谱系数。在图2所示倒谱的定义中,第一步骤是算出  $X(n)$  的 Z 变换  $X(z)$ 。而 LPC 系数与  $X(n)$  的 Z 变换  $X(z)$  有式(1)的对应关系。因此,  $X(z)$  可由 LPC 系数求出,然后再按上图的步骤则可求出倒谱系数。这样的倒谱系数就叫 LPC 倒谱系数。

## 2.2 LPC 谱加权的倒谱

由参考文献[3]知,在人的语音中,说话人声门的特征分布在中/低频带内。摩擦音的特征主要分布在高频段,另外,高频段还包含说话人声道全长的信息及声带的各种交叠方式(cross-modes of the vocal tract)的信息。所以对高频段和低频段进行增强,有助于说话人识别任务。另外,中频段内往往能量较高,而包含的主要是与话语内容有关的信息。在与文本无关的任务中,对中频段的能量适当减弱,将有利于识别性能的提高。图2给出了 LPC 谱加权倒谱的计算过程。在用 LPC 系数求出该段语音的频谱后,对各频段分别加权,最后再转化为频域加权的 LPC 倒谱。

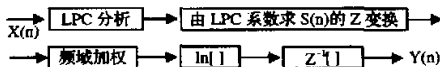


图2 LPC 谱加权的倒谱的计算过程

## 3 混合高斯模型(GMM)

GMM 目前常被用于“与文本无关的说话人识别”,具有良好的识别性能<sup>[5]</sup>。

定义:一个混合高斯模型的概率密度函数是由  $M$  个高斯概率密度函数加权求和而得到的,如下式所示:

$$P(x|\lambda) = \sum_{i=1}^M C_i b_i(x) \quad (2)$$

这里,  $P(x|\lambda)$  是某个 GMM 模型(说话人)的观测值  $x$  的概率密度函数。 $x$  是  $d$  维的随机矢量,  $C_i, i=1, 2, \dots, M$  是高斯分布混合时的权重值。 $b_i(x), i=1, 2, \dots, M$  是第  $i$  个单个高斯分布概率密度函数,如下式:

$$b_i(x) = \frac{1}{(2\pi)^{\frac{d}{2}} |\Sigma_i|^{\frac{1}{2}}} \exp \left\{ -\frac{1}{2} (x - \mu_i)^T \Sigma_i^{-1} (x - \mu_i) \right\} \quad (3)$$

上式中  $\mu_i$  是均值矢量,  $\Sigma_i$  是协方差矩阵。另外,混合权重值

$$\sum_{i=1}^M C_i = 1 \quad (4)$$

$\lambda$  为混合高斯模型的参数,由各个高斯模型的均值矢量、协方差矩阵、和混和权重值共同组成,如下式:

$$\lambda = \{C_i, \mu_i, \Sigma_i\} \quad i = 1, 2, \dots, M$$

在本实验中,每个说话人用一个混合高斯模型来参数化,即求出一个模型参数  $\lambda$  去对应一个说话人。

判别准则:给定  $M$  个说话人  $Y_1, Y_2, \dots, Y_M$  的高斯混和模型的参数集  $\{\lambda_i\} (i=1, 2, \dots, M)$  及测试语音的特征矢量集  $X = x_1, x_2, \dots, x_T$  时,可算出  $X$  属于第  $j$  个模型  $\lambda_j$  的概率,即后验概率,见下式:

$$P(X|\lambda_j) = P(X_1|\lambda_j)P(X_2|\lambda_j)\dots P(X_T|\lambda_j) \quad (5)$$

为了避免运算溢出,上式等号两边取对数得:

$$\ln P(X|\lambda_j) = \ln P(X_1|\lambda_j) + \ln P(X_2|\lambda_j) + \dots + \ln P(X_T|\lambda_j) \quad (6)$$

根据 bayes 判别规则,从  $M$  个说话人集合中确认未知说话人得公式如下:

$$i^* = \arg(\max_{1 \leq i \leq M} (\ln P(X|\lambda_i))) \quad (7)$$

## 4 说话人确认系统的实现及实验结果

### 4.1 说话人确认系统实践

本实验共录制了 30 人的语音,18 个男声和 12 个女声。语音在普通实验室环境下录制,麦克风的音量在半刻度以上。在录音开始的 1~2 秒录的是室内噪音。录音人按照日常说话习惯说话,没有特殊要求,内容不限,本实验的语音中有方言还有英语。每个人分别录两段语音,一段长一分钟,用于训练说话人概率模型,另一段语音长 30 秒,用作测试语音。语音经 11025Hz 采样及 16Bit 量化后检测端点并分为 20ms 的短时帧。

### 4.2 系统框图

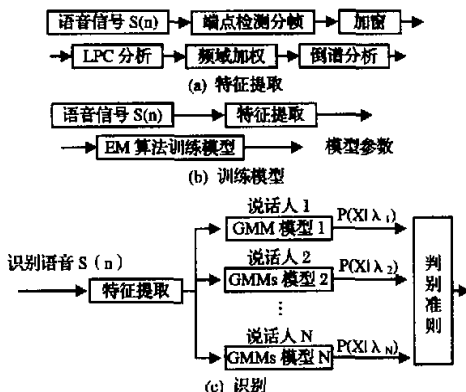


图3 系统框图

### 4.3 LPC 谱加权曲线图

选加权系数有很多方法,本次实验用的是一组经验值,原则上,应适当增强高频和低频部分的幅值,降低中频部分的幅值。如图4,图中曲线为  $0 \sim 2\pi$  的余弦曲线的近似。 $f_l, f_h$  为低、中、高频的分界。实验中,该组权值可使识别性能得到明显改善。

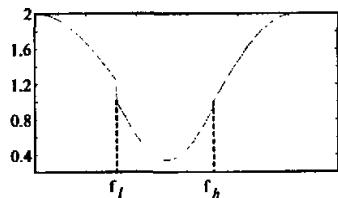


图4 频谱加权曲线

#### 4.4 实验结果

说话人的人数对识别结果的影响:

	10 人	20 人	30 人
使用 LPC 倒谱	100%	90%	80%
使用频域加权 LPC 倒谱	100%	100%	100%

当说话人集合的人数增多时,实验表明,改进后的方法取得较好的识别率。

测试语音长度对识别率的影响:

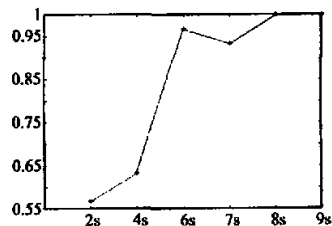


图5 正识率曲线

图5为正识率和测试语音音长之间的关系曲线。横轴表示测试语音音长而纵轴为正识率。这里,训练语音长20s,说话人集合为30人。由图可见,当语音长度超过6s时,正识率显著提高。当音长超过8s时,正识率可到100%。

## 5 结论

由以上实验可知,使用改进后的特征参数可以有效的提高正识率,而且所需测试语音的长度、训练语音的长度均不长,有利于实际应用。另外,实验表明,GMM可对随机变量的概率分布情况提供较好的近似。用于分类和模式识别时,该方法简单而有效。但是,当若干个模式的概率分布较接近时,将降低正确判别的比率。这时对各模式的概率分布进行预先调整,突出模式的特征可有效提高正识率。

## 参考文献

- [1] 牟晓隆,胡超秀,吴文虎.与文本无关的复合策略说话人辨识系统[J].清华大学学报(自然科学版),1997,37(3):16-19.
- [2] 李蘼华.倒谱参数基音信息有效结合进行说话人辨识[J].信号处理,2000,(1):16-19.
- [3] Qiguang Lin, En-Ee Jau, ChiWei che, et al.. Selective Use Of The Speech Spectrum And A VQGM Method For Speaker Identification [A]. ICSP 96[C]. CDROM, Volume 4, SuP1L2.
- [4] 拉宾纳, R. W. 谢弗语音信号数字处理[M]. 北京:科学出版社, 1984.
- [5] K. Markov, S. Nakagawa. Text - Independent speaker identification on TIMIT database[A]. Proceedings, Acous. Soc. Jap. [C], 1995, 83-84.