

用音乐指纹匹配作为研究路线,该路线优点是特征精简,对原音乐匹配精准,缺点是单当音频片段发生失真或非线性变化,音乐指纹也将演变成另一个截然不同的版本,从而失配。

违规音乐识别方案需要具备对音色、音调、节奏变化的鲁棒性,所以在违规音乐识别的场景中,音频特征的选择采用了 CENS 特征与节奏特征相结合的方式,以 CENS 特征模糊化音高信息,能够代表音乐在不同音阶上的能量分布特性,使用节奏特征代表人耳听力的感知效果,和 CENS 特征结合产生相互纠正的效应。

3 违规音乐识别方案

解析不同的违规音乐文件构建违规音乐数据库,以等长时间切片的音频存储到数据库,下文中称之为匹配样本 S 。数据库中每个匹配样本包括 CENS 特征、节奏特征、音乐时长和名称等基本信息。其中,匹配样本的 CENS 特征表示为 SC ,节奏特征表示为 ST 。待识别的音频片段称为目标样本 Q ,目标样本的 CENS 特征表示为 QC ,节奏特征表示为 QT ,目标样本和匹配样本默认时间等长。

3.1 特征提取

CENS 特征来自于色度特征的变换,色度特征体现为 12 维向量序列,每一个维度代表等律音阶的一个音高。将量化的色度特征向量序列与长为 w 的汉明窗进行分量卷积,然后基于参数 d 进行降采样,再次产生了一系列的向量,最后根据欧几里德范数进行归一化得到。

除此之外,仅依靠色度特征代表的声波能量分布情况并不能完全代表一首歌曲的全部信息,人耳对于不同乐曲的一个较高的区分度来自于音乐的节奏,需要结合音频的节奏特征进行违规性判定。

3.2 匹配策略

(1) CENS 特征匹配

CENS 特征具有时间分辨率属性,如图 1 所示,纵向排布的 12 维度对应 12 个八度等效的音高,而横向对应着每个音高在时间维度上的能量变化情况。每一个变化的单位是一个小方格,也是计算的基本单位,小方格的宽度即为时间分辨率。图中的分辨率是 500ms,即 500ms 内的能量分布变化会被忽略,从而消除了发音条件带来的无法预期的局部时间偏差。

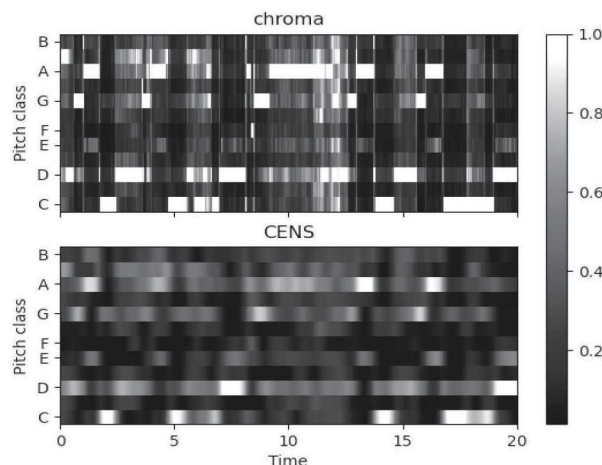


图1 色度特征及CENS特征示意图

CENS 特征纵向的一列代表音频在此单位时间内能量在不同音阶上的分布方式,也是检索过程中匹配的对象,在匹配样本中此列向量表示为 sc_i ,匹配样本和目标样本时间定长,在时间分辨率不变的情况下,包含的向量数目均为常数 N ,表示为 $SC=(sc_1, sc_2, \dots, sc_N)$ 。 sc 由 12 个维度构成,即 $sc=(v_1, v_2, \dots, v_{12})$ 。由于 CENS 特征经过归一化处理,所以 sc 向量是单位向量,几何上可以表示为在单位球体表面的一个点。

在时间维度上,对于 CENS 特征向量的对比具体化为计算 sc 和 qc 向量之间的差异, sc 、 qc 序列间的差异计算采用 DTW (Dynamic Time Warping) 方案,以兼容局部节奏的差异性。在音阶维度上,对于向量 sc 和 qc 的对比具体化为计算其中元素 $v(j)$ 之间的差值。对比的结果称为样本距离。对于原歌曲相同位置的片段,样本距离趋近于 0,其他同乐谱不同演绎方式的歌曲,歌曲风格越相似、音色特点越相似,样本距离越小。

(2) 节奏特征匹配

节奏特征 $T=(t_1, t_2, \dots, t_M)$,体现为一维数组, ST 和 QT 的匹配采用 DTW 算法。对于原歌曲相同位置的片段,样本距离接近于 0,其他同乐谱不同演绎方式的歌曲,节奏特征越相似,样本距离越小。

3.3 检索策略

匹配样本采集的过程中,设置片段的定长区间为 D ,从违规音乐起始时刻开始,以固定步长 $Step$ 向后逐一划分,得到的音乐片段进行特征和其他属性提取,然后入库。

将目标样本和匹配样本的 CENS 特征和节奏特征进

行逐一对比,在比对过程中,如果 CENS 特征与节奏特征的样本距离均趋近于 0,则停止对比,当前的匹配样本所属歌曲即为目标样本的来源,其所处于歌曲的位置也就是目标样本在原歌曲中的位置。否则,如果所有匹配样本得到的样本距离都不趋近于 0,将所有样本的 CENS 特征和节奏特征产生的距离降序排列,分别得到列表 L_{cens} 和 L_{tempo} ,对于 L_{cens} 取 $topK$,其中命中的歌曲对应应在 L_{tempo} 中的节奏距离若不在 $topK$ 之列,则进行剔除,在余下的命中歌曲中取 L_{cens} 和 L_{tempo} 的公共部分,并选取距离值最小的,即为判定的命中违规音乐。

3.4 容错机制

在 CENS 特征的匹配过程中,利用 CENS 特征的局部模糊化和 DTW 的动态偏差,保证了匹配方案的局部稳定性,但是对于不同的歌曲演奏版本,全曲的节奏快慢差异常能达到 10% 甚至 20% 之多,在这种情况下,DTW 的对比方案会在匹配结束的情况下产生不属于区间 D 的冗余内容。

基于以上考虑,要消除较大幅度的全局节奏差异,需要对匹配样本进行区间 D 或者采样率调整。在长度调整的情况下,在原有匹配的流程中,添加两组匹配流程,其中一组将目标样本进行时域上的伸缩,以比例 α 减掉或增加匹配样本长度尾部内容,由此得到的 L_{cens} 和 L_{tempo} 会增加至三组,分别为 (L_{cens}, L_{tempo}) 、 (L_{cens1}, L_{tempo1}) 、 (L_{cens2}, L_{tempo2}) 。在检索的过程中,对三组结果的 $topK$ 取并集,并用 L_{tempo} 的并集进行剔除。同理,在采样率调整的情况下,仍然添加两组匹配流程,其中两组的 CENS 特征生成过程中将 (w, d) 参数分别向上和向下调整,达到容错的目的。

4 实验分析

实验分两类共计三组:第一类测试出自于原曲的片段在原曲的 Step 递增过程中 CENS 特征距离和节奏特征距离产生的变化,共计一组;第二类测试第一组通过反复选取不同条件的目标样本,在匹配样本音乐数据库中进行检索,记录违规性判定的准确率,第二组通过改变匹配样本区间 D 进行容错方案的验证,同样记录判定准确率,共计两组。

4.1 实验条件

实验基于 python3 开发程序,在 windows 环境下执行,音乐样本库选型采用 MySQL5.7。

(1) 数据预备和特征提取。

一类实验选取电子琴版本回家,歌曲长度共计 173 秒,目标样本在歌曲中位置处于 13 秒,片段定长区间 D 为 30 秒。

二类实验选取 30 首不同的歌曲构建违规音乐数据库,采集音乐的基本信息和整曲特征,所有音乐总长度共计 5639 秒,片段定长区间 D 为 30 秒,Step 为 5 秒,进行切片,得到切片数量总计 955 片。二类实验三选取的目标样本长度分别为:30s+15%、30s、30s-15%。

(2) 匹配和检索。在 SC 和 QC 向量进行匹配的时候,DTW 算法的最大容忍长度差异设置成 16,超过 16 则不被允许。在结果判定过程中,原曲判定 CENS 特征的阈值为 2,节奏特征的阈值为 2,若没有命中原曲,需要从 L_{cens} 和 L_{tempo} 中进行比较筛选,候选范围 $topK$ 取 $K=20$ 。

(3) 测试目标样本。用于测试的目标样本总数为 16 个,其中 3 个来自于原曲,其他的来自于不同演唱者、不同演绎方式(独奏、合唱)、不同乐器。目标样本的片段区间 D 仍为 30 秒。

4.2 实验结果

一类实验第一组结果如图 2 所示,在 13~15 秒的位置,CENS 特征距离和节奏特征距离同时出现了最低点,在方案中,即可认定为出自于原曲并终止检索。图中的 CENS 特征距离在 85~90 秒位置也出现了低点,但是对应的节奏距离不在原曲判定的阈值范围内,故可以被舍弃。

二类实验第一组结果如表 1 所示,CENS 特征取最小作为判定音频相似的依据,准确率 75%,使用节奏特征作为唯一依据,正确率 62.5%,而将两者结合能达到 81.25%,说明此二种特征均能表现音频的一部分特点,但不足以完全判定。

二类实验第二组结果如表 2 所示,其中以 30s+15% 和 30s-15% 的特性都会遗漏一定的相似样本,是因为将目标样本进行长度调整可能造成原本完全合适的相似音乐缺少或冗余内容,但是二者检索到节奏差异较大的匹配样本不同,故能在结合到原区间方案的时候产生正确率的提高。

4.3 分析总结

从一类实验第一组的结果可以看出,结合 CENS 和节奏特征进行原曲判定尤其是片段的位置判定,具有较优秀的准确率与可靠性,进一步证明了将特征结

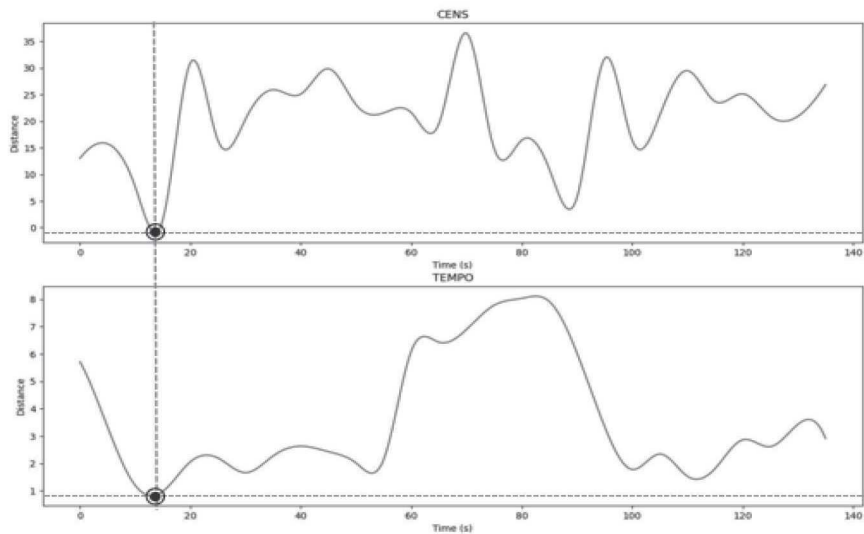


图2 CENS特征距离及节奏特征距离示意图

表1 判定特征对比实验准确率表

判定特征	CENS特征	节奏特征	二者结合特征
判定准确率	75.00%	62.50%	81.25%

表2 特征区间对比实验准确率表

特征区间	30s-15%	30s	30s+15%	三者结合
判定准确率	62.50%	81.25%	68.75%	87.50%

合是全方位解析音乐文件的一个重要方法。

从二类实验第一组的结果可以看出，基于 CENS 特征的违规音乐识别准确率稍显逊色，但是虚假识别的结果，与目标样本真正所属音乐的节奏特征和乐曲风格具备一定的相似度，这印证了 CENS 特征代表了不同音阶的能量分布特征的属性；基于节奏特征的音乐识别在准确率上有所提高，虚假识别结果的乐曲节奏感相似，但乐曲风格可以产生巨大的差异。而两种特征的结合在查准率上有着显著的提高。

从二类实验第二组的结果可以看出，将乐曲的全局节奏差异进行拉伸或者压缩，能够匹配到同一首歌曲不同演绎方式中的不同部分，将三种不同节奏的目标样本进行联合查询，能完善音乐检索的效果。

5 结束语

针对 5G 消息中的违规音乐提出了鲁棒的违规音乐识别方案，能够在单声部、多声部、多方式演绎场景中实现违规音乐的高效识别。在未来科研工作中，提高实时性和检索结果的广泛同类乐曲关联可以成为研究目标。

参考文献：

[1] Joó S, Jo S, Chang D Y. Yoo: Melody extraction from polyphonic audio signal MIREX 2009[J]. MIREX Audio Melody Extraction Contest Abstracts, 2009.

[2] Shih H H, Narayanan S S, Kuo C. Multidimensional humming transcription using a statistical approach for query by humming systems[C]//Acoustics, Speech, and Signal Processing, 2003. Proceedings. (ICASSP'03). 2003 IEEE International Conference on. IEEE, 2003.

[3] Wang A. An industrial-strength audio search algorithm[C]//ISMIR 2003, 4th International Conference on Music Information Retrieval, Baltimore, Maryland, USA, 2003. ★

作者简介

谢仪颀：硕士，现任中国移动通信集团设计院有限公司咨询设计师，主要研究方向为信息安全、计算机视觉和语音。

杜刚：高级工程师，博士，现任职于中国移动通信集团设计院有限公司，主要研究方向为信息安全、人工智能。

张晨：硕士，现任职于中国移动通信集团设计院有限公司，主要研究方向为信息安全、网络安全、内容安全、大数据分析。

杜雪涛：硕士，现任职于中国移动通信集团设计院有限公司，主要研究方向为信息安全、网络安全、内容安全。