

大雅相似度分析

论文标题：基于机器学习的安卓移动用户情绪分析系统的设计与实现

检测日期：2018年06月15日

作者：王旭东

正文字符数：21600

正文字数：18734

检测范围：大雅全文库

一、总体结论

文献相似度	重复字符数	最密集相似处	密集相似处	非密集相似处	前部相似度	中部相似度	尾部相似度
6.07%	1312	1	3	14	8	5	5

二、相似片段分布



三、典型相似文献

相似图书

作者	题名	出处	相似度
刘爽	支持向量机在自动文本分类中的应用	大连：大连海事大学出版社，2014.07	0.77%
袁宏永	突发事件及其链式效应理论研究与应用	北京：科学出版社，2016.05	0.77%
张小超;吴静珠;徐云	近红外光谱分析技术及其在现代农业中的应用	北京：电子工业出版社，2012.04	0.77%
罗泽举	数据挖掘理论、方法与应用	北京：电子工业出版社，2014.12	0.77%
王仁武	Python与数据科学	上海：华东师范大学出版社，2016.03	0.76%
王立国;赵春晖	高光谱图像处理技术	北京：国防工业出版社，2013.05	0.69%
徐珉久	R语言与数据分析实战	北京：人民邮电出版社，2017.01	0.67%
托尔戈	数据挖掘与R语言	北京：机械工业出版社，2013.04	0.66%
吴定海	柴油机振动信号分析与故障诊断研究	北京：国防工业出版社，2012.12	0.66%
王志良;孟秀艳	人脸工程学	北京：机械工业出版社，2008.07	0.64%
西蒙 蒙策尔特 (Simon Munzert)	基于R语言的自动数据收集 网络抓取和文本挖掘实用指南	北京：机械工业出版社，2016.03	0.64%
郑志明;缪绍日;荆丽丽	银行业信息化丛书 金融数据挖掘与分析	北京：机械工业出版社，2016.01	0.64%
高木干雄;下田阳久;孙卫东	图像处理技术手册	北京：科学出版社，2007.08	0.63%
林强	行为识别与智能计算	西安：西安电子科技大学出版社，2016.11	0.63%
王营冠	无线传感器网络	北京：电子工业出版社，2012.06	0.63%
李颖宏;张永忠;王力	道路交通信息检测技术及应用	北京：机械工业出版社，2014.01	0.63%

梁循	面向社会化媒体大数据的社会计算	北京：清华大学出版社，2014.12	0.63%
刘龙	基于振动测试的结构损伤识别若干方法研究	上海：上海交通大学出版社，2014.04	0.63%
刘富强	数字视频图像处理与通信	北京：机械工业出版社，2010.06	0.63%
RobFarber	高性能CUDA应用设计与开发 方法与最佳实践	北京：机械工业出版社，2013.01	0.63%
夏虹;刘永阔;谢春丽	设备故障诊断技术	哈尔滨：哈尔滨工业大学出版社，2010.03	0.58%
阿培丁	机器学习导论 2版	北京：机械工业出版社，2014.04	0.56%
刘尚旺	物联网视频图像感知新技术	北京：科学出版社，2016.05	0.56%
侯媛彬;杜京义;汪梅	神经网络	西安：西安电子科技大学出版社，2007.08	0.56%
李蕾;王小捷	机器智能	北京：清华大学出版社，2016.06	0.56%
谭磊;范磊	Hadoop应用实战	北京：清华大学出版社，2017.01	0.56%
刘知远	大数据智能 互联网时代的机器学习和自然语言处理技术	北京：电子工业出版社，2016.01	0.56%
韩敏	基于神经网络的监督和半监督学习方法与遥感图像智能解译	北京：中国水利水电出版社，2015.12	0.56%
陆旭	文本挖掘中若干关键问题研究	合肥：中国科学技术大学出版社，2008.12	0.56%
吴今培;孙德山	现代数据分析	北京：机械工业出版社，2006.02	0.56%
丁世飞	人工智能 第2版	北京：清华大学出版社，2011.01	0.56%
戚文静	网络安全与管理 第2版	北京：中国水利水电出版社，2008.06	0.56%
邓书斌	ENVI遥感图像处理方法	北京：科学出版社，2010.06	0.56%
刘君华	智能传感器系统	西安：西安电子科技大学出版社，2010.05	0.56%
中国石油管道公司	油气管道安全预警与泄漏检测技术	北京：石油工业出版社，2010.07	0.55%
徐丽敏	鲁棒性说话人识别技术 在移动商务中的应用研究	南京：南京大学出版社，2011.12	0.55%
张颖伟;S.Joe Qin	复杂工业过程的故障诊断	沈阳：东北大学出版社，2007.12	0.55%
中国人工智能学会	中国人工智能进展 2001 中国人工智能学会第9届全国学术年会暨中国人工智能学会成立二十周年庆祝大会论文集	北京：北京邮电大学出版社，2001.11	0.55%
白鹏	支持向量机理论及工程应用实例	西安：西安电子科技大学出版社，2008.08	0.55%
杨东方;陈豫	数学模型在生态学的应用及研究 23	北京：海洋出版社，2013.05	0.55%
刘晓明	中国电子学会电子系统工程分会第十五届信息化理论学术研讨会论文集	北京：中国市场出版社，2008.09	0.55%
叶枫;周根贵;吕旭东	基于规则与案例推理的临床决策支持	北京：科学出版社，2014.01	0.55%
王小川;史峰;郁磊	MATLAB神经网络43个案例分析	北京：北京航空航天大学出版社，2013.08	0.4%
马希青	机械制图	北京：机械工业出版社，2015.09	0.14%
粟芳;魏陆	瑞典社会保障制度	上海：上海人民出版社，2010.01	0.12%

王秀娥;夏冬	市场调查与预测	北京：清华大学出版社，2012.01	0.12%
杨清德;冉洪俊	学电工电路就这么简单	北京：科学出版社，2015.07	0.11%
朱祥庭;张晶	数控机床电气装调	北京：清华大学出版社，2017.01	0.11%
付钰	信息对抗理论与方法	武汉：武汉大学出版社，2016.07	0.11%
赵志儒	色铅笔的宠物萌绘	北京：机械工业出版社，2014.02	0.11%

相似报纸

作者	题名	出处	相似度
	李嘉诚斥资千万美元投资“3D打印肉”	呼和浩特晚报，2014.06.26	0.34%
	李嘉诚斥资千万美元投资“3D打印肉”	天天商报，2014.06.26	0.34%
	这些年李嘉诚投资过的另类项目	南国早报，2014.06.26	0.34%
	李嘉诚斥资千万美元投资“3D打印肉”	新京报，2014.06.25	0.34%
	李嘉诚千万美元投资“3D打印肉”	贵州商报，2014.06.26	0.29%
	李嘉诚牌“人造肉”你会买吗？	东方今报，2014.06.26	0.29%
薛飞	2015青岛 中国财富论坛昨举行郭树清出席李群致辞张新起主持	青岛晚报，2015.06.21	0.06%

相似期刊

作者	题名	出处	相似度
马喜波;阎爱侠	基于支持向量机的有机化合物水溶解度的分类和预测的研究	计算机与应用化学，2008，第12期	0.88%
孟娇茹	航空发动机孔探损伤识别方法	黑龙江科技学院学报，2009，第1期	0.77%
齐保林;李凌均	基于SVR的设备状态趋势预测方法	矿山机械，2007，第2期	0.72%
马超	基于Web信息使用改进的无监督关系抽取方法构建交通本体	计算机系统应用，2015，第12期	0.7%
费胜巍;孙宇	用SVRM预测变压器油中溶解气体量	高电压技术，2007，第8期	0.69%
郭俊文;李开成;何顺帆;张明	基于改进不完全S变换与决策树的实时电能质量扰动分类	电力系统保护与控制，2013，第22期	0.67%
刘路;李弼程;张先飞	基于正反例训练的SVM命名实体关系抽取	计算机应用，2008，第6期	0.66%
谷玉霞	基于支持向量机的企业财务风险评价研究	会计师，2011，第2期	0.66%
翟嘉;胡毅庆;徐尔	用于多分类问题的最小二乘支持向量分类-回归机	计算机应用，2013，第7期	0.66%
赵向军;路梅	垃圾邮件过滤算法研究	徐州师范大学学报(自然科学版)，2006，第4期	0.66%
段丹;郭绍忠;李志博;刘沙	基于邮件分类的敏感社团挖掘技术	计算机应用，2007，第12期	0.64%
周杰;李弼程;林琛;韩永峰	基于网络资源的实体知识库系统设计研究——以政府相关实体知识库为实例	情报科学，2016，第1期	0.64%
巩军;刘鲁	一种k-NN文本分类器的改进方法	情报学报，2007，第1期	0.64%
梁宏斌;严正俊	LS-SVM在垃圾邮件过滤中的应用	现代电子技术，2007，第17期	0.64%
罗显科;柴毅;李华锋;梁奕欢	半监督增量式SVM在故障诊断中的应用研究	世界科技研究与发展，2013，第4期	0.64%
咎红英;张军琿;朱学锋;俞士汶	副词“就”的用法及其自动识别研究	中文信息学报，2010，第5期	0.64%
杨立东;王晶;谢湘;匡镜明	基于Tucker分解的音频分类研究	信号处理，2015，第2期	0.63%
陈晗;戴在平	家电控制系统的语音关键词识别算法研究	电声技术，2008，第4期	0.63%
黄丽芬;黄大荣;赵玲	基于多元集对分析联系数AHP方法的网络安全评估	指挥控制与仿真，2015，第4期	0.63%

郑继明;魏国华;吴渝	有效的基于内容的音频特征提取方法	计算机工程与应用, 2009, 第12期	0.63%
陈靓;张成;陈尚姿	基于支持向量机的电费信用评估模型	电力信息化, 2008, 第10期	0.63%
蔡文学;邱珠成;黄晓宇;陈康	基于改进个性诊断和热门路段的路况估计	系统工程, 2015, 第7期	0.63%
崔萌;张春雷	LIBSVM, LIBLINEAR, SVMmuticlass比较研究	电子技术, 2015, 第6期	0.63%
蒋先刚;张盼盼;盛梅波;胡玉林	基于级联与组合属性形态学滤波的模糊边界目标识别	计算机工程, 2016, 第3期	0.63%
王守觉;孙华;柳培忠;廖英豪;丁兴号;郭东辉	基于仿生形象思维方法的图像检索算法	电子学报, 2010, 第5期	0.63%
文守逊;王仁杰	基于PCA-SVM的商业银行高级管理人员能力评价	科技管理研究, 2008, 第9期	0.63%
姚慧;孙颖;张雪英	情感语音的非线性动力学特征	西安电子科技大学学报, 2016, 第5期	0.63%
黄浩;哈力旦	大间隔高斯混合模型的快速参数更新算法	计算机工程, 2010, 第3期	0.63%
梁春燕;杨琳;汪俊杰;张建平;颜永红	音子配列学语种识别系统中特征选择方法的研究	声学学报, 2013, 第2期	0.63%
井小沛;汪厚祥;聂凯;罗志伟	面向入侵检测的基于IMGA和MKSVM的特征选择算法	计算机科学, 2012, 第7期	0.63%
施洁斌	基于支持向量机的文本自动分类试验研究	现代图书情报技术, 2004, 第7期	0.63%
徐冰;郭绍忠;黄永忠	基于朴素贝叶斯分类算法的活跃网络结构挖掘	计算机应用, 2007, 第6期	0.63%
周强;张晓俊;顾济华;赵鹤鸣;朱俊杰;陶智	嗓音多频带非线性分析的声带病变识别	声学学报, 2014, 第1期	0.63%
张旭梅;石瀚凌	基于分类挖掘方法的商业银行个人理财业务客户流失分析	工业工程, 2011, 第6期	0.63%
伏晓丽	基于内容的视频检索中相关反馈应用综述	知识管理论坛(网络版), 2009, 第12期	0.63%
姜贤林;郭秀清	基于支持向量机的质量控制软测量建模	计算机应用, 2008, 第9期	0.63%
彭芳青;厉力华;徐伟栋;刘伟;张娟;邵国良	基于Multi-Agent的乳腺钼靶图像肿块分类方法	传感技术学报, 2010, 第2期	0.63%
	动态	现代图书情报技术, 2004, 第7期	0.63%
徐冉冉;琚昊霖;李朝锋	非平衡二叉树主动学习支持向量机	微电子学与计算机, 2013, 第5期	0.63%
郝文峰;顾建祖;汤灿	应用小波支持向量机的预应力混凝土碳化深度研究	工业建筑, 2009, 第A1期	0.63%
严良达;陶剑文	稀疏表示亲近支持向量机	计算机工程与应用, 2014, 第19期	0.61%
杨红敏;何丕廉	基于增量式拉普拉斯嵌入和SVM的图像识别	计算机仿真, 2007, 第11期	0.6%
邢笛;葛洪伟;李志伟	模糊支持张量机图像分类算法及其应用	计算机应用, 2012, 第8期	0.6%
纪四维;李熊达;唐静远;师奕兵	模拟电路故障诊断的子空间集成方法	计算机工程, 2011, 第17期	0.59%
刘康;钱旭;王自强	基于流形主动学习的遥感图像分类算法	计算机应用, 2013, 第2期	0.58%
孙进;张征;周宏甫	基于脑机接口技术的康复机器人综述	机电工程技术, 2010, 第4期	0.58%
刘忠宝;赵文娟	面向大规模信息的用户分类方法研究	微电子学与计算机, 2013, 第6期	0.58%
王安娜;李明;李华;栾峰	基于支持向量机的容差电路故障诊断	华北电力大学学报, 2005, 第A1期	0.56%
赵立权;谢妮娜	基于小波变换和改进的RVM的电能质量扰动分类	电工电能新技术, 2013, 第4期	0.56%
范莹;计华;张化祥	一种新的基于模糊聚类的组合分类器算法	计算机应用, 2008, 第5期	0.56%

其他网络文档

作者	题名	相似度
吕克洪	基于时间应力分析的BIT降虚警与故障预测技术研究	0.87%
	支持向量机算法用于癌症数据建模	0.78%

姜林	化学信息学方法研究及其在环境、生物学中的应用	0.77%
饶蓝	基于支持向量机的网络攻击检测研究	0.77%
朱少华	电网智能报警研究	0.77%
王晓瑜	基于支持向量机的入侵检测技术	0.77%
谭海龙	基于聚类分析与SVM的电力短期负荷预测研究	0.77%
魏志静	基于人工神经网络的分类方法研究及其在个人信用评估中的应用	0.77%
吕冰	基于核技术的人脸识别应用研究	0.77%
	农业机械学报	0.77%
张茂雨	支持向量机方法在结构损伤识别中的应用	0.77%
徐玉兵	基于支持向量机的概率密度估计及其在分布估计算法中的应用	0.77%
张世荣	支持向量机文本分类算法研究	0.77%
荣光	中文文本分类方法研究	0.76%
吴烜;沙明;李智毅	支持向量机算法诊断测厚仪CS值电压自动漂移故障分析	0.72%
侯铁民	基于支持向量机的多属性大规模数据分类算法的研究	0.71%
庄振华	基于Laplace谱的基因表达谱数据分类研究	0.71%
白宇	网络文本信息的自动分类方法研究	0.69%
杨进军	基于支持向量机的高分辨率雷达目标识别技术研究	0.69%
王占强	基于SVM的模型选择和参数优化方案研究与实现	0.68%
杨家勇	城市道路路面使用性能预测及养护决策研究	0.67%
刘佳	基于离散增量结合二次判别法预测蛋白质相互作用及DNA甲基化位点	0.67%
周永	一种复合的双引擎智能垃圾邮件过滤方法	0.66%
杨培洁	基于支持向量机的商业银行信用评估方法研究	0.66%
	新闻搜索引擎的设计	0.66%
	开题20091222-yangx	0.66%
陈伟	数据挖掘在养老基金决策支持系统中的应用	0.66%
王坤华	基于PSO和SVM的上市公司财务危机预警研究	0.66%
林青	垃圾邮件过滤技术研究	0.66%
高妮	支持向量机及其在乳腺癌辅助诊断系统中的应用研究	0.66%
尹志喜	基于内容的垃圾邮件过滤技术研究	0.66%
王志强	WEB文本信息抽取和分类研究	0.66%
段艳华	基于基因表达谱的肿瘤分类特征基因选择研究	0.66%
任金龙	基于支持向量机的铁路泥石流危险性评价方法研究	0.63%
常青	基于机器学习算法的Web文本挖掘应用研究	0.63%
张菁	探地雷达地雷图像处理与目标识别方法	0.63%
徐燕子	核空间聚类算法及其在大规模支持向量机应用中的研究	0.63%
骆健	基于茧丝纤度序列的多总体分类判别及其比较	0.63%
杨柯	基于关联规则的中文文本自动分类算法研究	0.63%
刘林	基于支持向量机的醋酸乙烯聚合率软测量研究	0.63%
童婧	蛋白质序列中RNA结合位点的预测	0.63%
	基于SVDD的烟叶异物检测技术研究	0.63%
彭敏晶	区域经济发展决策支持建模与仿真技术研究	0.63%
	一种 RBF 核SVM 的参数选择方法#	0.63%
侯凤勤	基于改进的组合模型的科技人才数量预测	0.63%
陈禹伶	基于部件融合的“我”体字库的建立	0.63%
王耿	反垃圾邮件算法的设计和实现	0.63%
米少华	基于机器学习的microRNA基因预测	0.63%



司德睿	基于文本内容的网页过滤技术研究	0.63%
	一种rbf核svm的参数选择方法	0.63%

四、全文相似情况

摘要

情绪时时刻刻伴随着人们的生活，并且每分每秒都在影响着人们的工作和生活的状态。而随着我国近些年经济快速发展，人们也逐渐意识到情绪对于生活的重要影响，越来越多的人开始关注自身情绪的状态和变化。

由于人的情绪与人体各项数据密切相关，安卓设备作为天然的数据收集器条件优越，本文以安卓手机为数据收集的载体，通过收集用户的运动数据、环境数据、设备使用情况等数据，基于机器学习分类器原理构造情绪分类模型，从而实现对用户情绪的识别。文章首先介绍了情绪分析的研究进展，并指出了不同情绪研究的优缺点。然后，对实验用到的平台和主要技术进行了简要的介绍，接着又对机器学习分类器的原理进行了分析。然后，结合分析和讨论详细介绍了本实验的数据收集、数据预处理、特征工程、模型构建过程，其中模型构建过程中基于K近邻算法、SVM算法、朴素贝叶斯算法，通过参数调整构建了三类七种不同的模型。

最后，本文还针对构建出来的模型，进行了准确性测试和性能测试，测试的结果表明，构建出来的情绪分析模型对于情绪的识别较为准确，同时具有良好的性能，达到了实验的预期目的。

关键字：机器学习 情绪分析 分类算法 安卓

ABSTRACT

Emotions are accompanied by people's lives all the time, and affects the people's work and life every minute. With the rapid economic development in China in recent years, people are gradually realizing the important influence of emotion on life. More and more people begin to pay attention to the state and changes of their own emotions.

People's emotions are closely related to people's data, and Android devices have advantage in collecting data information. For these reasons, we choose Android devices as carriers to collect data information, including user's motion data, environmental data, and device usage data. The emotion classification model is constructed to recognize user emotions based on the principles of machine learning classifier. The article first introduced the research progress of emotion analysis and pointed out the advantages and disadvantages of different emotion studies. Secondly, the platform and main technology used in the experiment are briefly introduced. Then the principle of the machine learning classifier is analyzed. Next, the data collection, data preprocessing, feature engineering, and model construction process of this experiment are described in detail in connection with the analysis and discussion. The K-nearest neighbor algorithm, SVM algorithm, and Naive Bayes algorithm are used to build the models, including three types of seven different models.

Finally, this article also carried out the accuracy test and the performance test on the constructed model. The test results show that the constructed emotion analysis model has good accuracy in the recognition of emotions, and has good performance at the same time, achieving the intended purpose of the experiment.

Keywords: Machine learning Emotion analysis Classifier Android

目录

第一章绪论.....	1
1.1研究背景.....	1
1.2研究意义.....	2
1.3研究现状.....	3
1.3.1非生理信号情绪分析.....	3
1.3.2基于生理信号情绪分析.....	4
1.4本文研究内容.....	5



第二章开发平台及原理分析.....	7
2.1开发工具简介.....	7
2.2Numpy和Pandas.....	7
2.3scikit-learn.....	8
2.4机器学习分类算法.....	9
2.4.1K最近邻算法.....	9
2.4.2支持向量机算法.....	10
2.4.3朴素贝叶斯算法.....	11
2.4.4决策树算法.....	13
2.4.5Adaboost算法.....	14
第三章数据处理和特征工程.....	17
3.1用户数据收集.....	17
3.2用户数据预处理.....	19
3.2.1数据集成.....	20
3.2.2缺失值处理.....	21
3.3数据属性观察.....	22
3.4属性的相关性分析.....	22
3.5加速度合成.....	23
3.6数据特征提取.....	23
3.7主成分分析.....	24
3.8数据标准化.....	25
3.9本章小结.....	26
第四章构建情绪分析模型.....	27
4.1K近邻分类模型.....	27
4.2SVC分类模型.....	30
4.3朴素贝叶斯分类模型.....	32



第五章运行与测试.....	35
5.1准确性测试.....	35
5.2性能测试.....	37
5.3本章小结.....	37
第六章总结与展望.....	39
致谢.....	41
参考文献.....	43

第一章 绪论

从2001年世界上第一部智能手机诞生开始,智能手机的发展速度可谓是日新月异。从以诺基亚为代表的塞班一家独大,到Android, ios, WindowsOS三足鼎立,到如今Android和ios几乎已经占领了智能手机操作系统95%以上的份额,而其中又以Android份额最多。伴随着智能手机数量的爆发式增长,手机的硬件也得到了极大的升级,功能也变得更加丰富。手机已经不仅仅是用来打电话,发信息,而是集通信,娱乐,办公,社交等诸多功能于一体的智能终端设备。随着手机内置的传感器越来越多,手机所能收集到的信息也越来越多,而且因为多数人都有随身携带的习惯,让手机俨然成为一个天然的“人体数据收集站”。而且,近些年机器学习热度越来越高,“情绪分析”作为其中的一个重要方向也吸引了大量研究人员的兴趣。鉴于智能手机的优势和特性,通过分析智能手机的数据来进行情绪分析,正逐渐成为机器学习的一个研究热点。

1.1 研究背景

情绪是人在探索和发现周围环境时,对外界环境刺激做出的心理或生理反应,是一种综合了人的思想,感觉,和行为的状况[[endnoteRef:2]]。常见的情绪有高兴,平静,悲伤,生气等。积极的情绪可以让人精神焕发,干劲十足,消极的情绪可以让人萎靡不振,日渐消沉,人的情绪极大影响了人们的日常生活和工作的效率。因此在这种情况下,对情绪的了解显得极为重要。 [2: [] 孟昭兰. 人类情绪. 上海人民出版社, 1989]

人的情绪变化时通常伴随着人体的一系列生理表现,或者行为表现,比如心跳加快,血压升高,脑电波活跃,激素分泌,又比如欢呼雀跃,手舞足蹈,唉声叹气,愁眉苦脸等等。正是因为情绪变化时,人体会产生这些变化,如果可以通过愈来愈发达的技术手段准确的获取人体的这些信息,再加上科学正确的分析,就可以判断出用户此刻的心情。情绪分析和识别是一项多领域交叉的科研问题,目前,多个领域都在对于情绪进行相关的研究,涉及的领域有医学,心理学,神经科学,人工智能等。

智能手机在发展过程中,为了丰富功能,增强用户体验,内置了越来越丰富的传感器,比如光线传感器,温度传感器,磁力传感器,重力传感器,陀螺仪,及速度传感器等等,这为研究人员进行实验提供给了非常有利的条件。

之所以选用安卓手机来收集这些信息,是因为,首先,越来越多的人拥有智能手机,其中绝大部分是安卓手机(下文“手机”,“智能手机”都指代“安卓手机”),安卓手机的大规模使用保障了数据的充足,在之前是不存在这样优越的研究条件的。其次,随着手机越来越智能,功能越来越多,浏览网页,聊天阅读,拍照娱乐,移动支付等等,使用手机的场景越来越多,人们越来越离不开手机,这使得很多人有了随身携带手机的习惯。并且,手机的携带对被测试者的正常生活影响较小,这使得持续获取数据成为了可能。

机器学习经过多年的发展,相关理论已经非常成熟。随着科技水平的飞速发展,计算机的处理能力大幅提升,并且伴随着数据量的暴增,人们进入“大数据”时代,在这样的时代背景下,机器学习具有传统统计难以匹及的优势,因此本文选用了机器学习的方式来进行情绪的识别。

1.2 研究意义

情绪与人们的生活密切相关,人生活中的每时每刻都对应着一种情绪,或开心,或难过,或平静。心理学研究表明,情绪很大程度上影响着人们的工作和生活,积极的情绪可以提高人体的机能,能够形成一种促进人体活动的动力,激励人去努力,提高工作的效率。消极情绪会让人觉得难受,会降低人的活力,比如活动起来动作迟缓,工作效率低下,消极的情绪还会对人的精力和体力造成影响,浑身无力,疲乏易累。因此,时刻了解自己的情绪状态和变化是很有必要的。

而且,随着近年来我国经济飞速发展,人们对于自己身体健康的关注也渐渐地从只关注生理健康到生理健康与心理健康并重,越来越多的人开始关注自己的情绪状态和情绪变化。但现代社会生活节奏较快,多数人常常无暇顾及自己身体信号、行为状态的微小变化,比如行走的速度、活动的幅度、室内室外活动时间等等,但相关研究表明人的这些行为举止的微小变化常常对应着人的情绪变化。传统的情绪识别一般分为基于生理信号和基于非生理信号两种。基于生理信号的情绪识别一般需要依赖于精密的仪器,对人体的心率、皮肤阻抗等生理信号进行采集并进而向相关分析,这种方式的情绪分析往往存在着一定的实施困难,比如仪器对人正常生活影响较大不便于采集真实的数据,又比如,这些仪器一般很容易受到干扰,收集来的信息往往掺杂了大量噪声,提高了分析工作的难度等。而基

于非生理信号的情绪分析一般集中在语音语调方面,并且往往容易被人的刻意伪装而欺骗。为了帮助人们随时了解自己的情绪状态,使人能够及时意识到这些变化,并进行有效的自我调节,同时考虑到智能手机已经非常普及并且大多数人已经养成随身携带智能手机的习惯,本文作者通过安卓设备对用户的运动数据、环境数据、手机使用情况等信息进行采集,将数据进行特征提取后,利用已经较为成熟的机器学习分类算法构建模型,实现对用户情绪的识别。通过这种方式可以减小数据收集设备对用户的干扰,并且人的行为活动不容易伪装,更容易获取到真实的数据,而且通过收集多方面的信息更能全面地反映一个人目前的情绪状态。

情绪识别这项技术的应用前景十分广泛,比如听音乐时,如果能够识别出用户此刻的心情,进而推送给用户符合当前心情状态的歌曲,那么用户的使用体验就会得到极大的提升。再比如,在医疗护理过程中,如果可以准确分辨出患者此刻的心情,尤其是具有表达沟通障碍的患者,进而对其提供针对性的护理,无疑对于治疗过程是有极大帮助的。再比如,如果可以准确获知用户使用一款产品时的情绪状态,那么开发人员就可以针对用户使用时的情绪进行更有针对性的优化处理,从而持续提供更好的服务。

1.3 研究现状

目前有关人的情绪状态分析识别的研究主要有两类方向,一类是基于生理信号进行情绪识别,一类是基于非生理信号进行识别[[endnoteRef:3]]。[3: [] 张迪等. 基于生理信号的情绪识别研究进展. 生物医学工程学报. 2015,2,13(1)]

1.3.1 非生理信号情绪分析

非生理信号的情绪分析研究主要包括基于面部表情的情绪分析,基于语音语调的情绪分析两种。

基于面部表情的情绪分析:不同的表情常常会反映不同的情绪,所以研究人员可以通过分析人的面部表情以及伴随的面部肌肉动作来判断人的情绪[[endnoteRef:4]]。比如嘴角上扬时,并且眼角出现皱纹,通常可以判断此时的情绪状态为高兴;眼睛瞪大,眉头紧皱,通常可以判断此时情绪状态为愤怒。基于面部表情的情绪识别既可以是基于局部特征进行的情绪分析,也可以是基于整体特征的情绪识别。前者是考虑到人在不同情绪下,人脸五官的形状,大小,以及相对位置会有差异,并以此作为情绪分析的数据特征,进行情绪分析。后者则是从整体出发,考虑不同情绪下整体面部特征的区别,提取特征的范围是整个人脸。目前基于面部表情的情绪识别已经取得一定进展,比如美国的情绪识别公司Affectiva, Affectiva通过将大量不同种族、年龄、性别的面部表情图片输入计算机,通过算法观察到所有脸部的纹理与皱纹,以及形状的变化等面部重要特征点变化,从而识别出当前人的情绪,该算法可以准确识别出高兴、难过、惊喜、生气、蔑视、厌恶、恐惧七种情绪。据称Affectiva的情感AI还将被集成到Voxpopme平台,使研究人员更容易获知观众的喜好,深入了解观众行为,减少视频工作的传统挑战。[4: [] 傅栩栩,叶健东,王鹏等. 人脸面部表情识别. 计算机与网络,2015,(10):70-71]

基于语音语调的情绪分析:不同的语言表达方式通常也意味着不同的情绪状态,比如,人的心情愉悦时,语音会轻快,语调会上扬,心情失落时,语调会变得低沉,语调会放缓。基于语音的情绪识别研究一般多关注于基音的频率、声音携带的能量、语速快慢、表达是否流利等特征。目前,基于语音的情绪识别技术已经实现商用,如2012年成立的Beyond Verbal就是一家语音情绪识别领域领先的公司,该公司实现了通过算法识别音域变化,从而对用户的情绪进行判断,可以识别出400个复杂情绪,甚至能够识别其中的微小差别。

基于非生理信号的情绪分析研究简单易行,操作难度小,但是存在识别不准确和不可靠的问题,毕竟人是种复杂的动物,常常会为了隐藏自己的情绪而伪装自己的声音表情。另一方面,对于一些病患人群可能无法提供非生理信号情绪识别所需的特征,所以对于这一人群也不适用。

1.3.2 基于生理信号情绪分析

另一大类是基于生理信号的情绪分析。这类研究关注的重点通常是人体的心率,血压,脑电波等生理信号[[endnoteRef:5]]。人处在不同的情绪状态时,表现出来的生理信号也是不同的。研究人员通过收集人在喜、怒、哀、乐等不同情绪下的各项生理指标,比如通过测量人体心脏跳动的频率、呼吸的频率、脑内电波的活跃程度、皮肤外表的阻抗等数据,然后与对应的情绪状态高兴还是悲伤建立对应联系,然后分析数据间的内在关系,从而建立数据模型实现对情绪的判断。生理信息相比非生理信号而言,不易被伪装,并且通过这种方式,在获得准确信息的前提下,分析结果准确率高。但是该方法也存在着收集困难,信号准确率无法保障,易被干扰的缺点。[5: [] 赵国朕,宋金晶,葛燕等. 基于生理大数据的情绪识别研究进展. 计算机研究与发展,2016,(1):80-92]

1.4 本文研究内容

下文主要研究基于机器学习的安卓移动用户情绪分析系统的设计与实现,即将从安卓设备采集来的数据,通过经典的机器学习算法进行分析,从而判断出使用者现在的心情状态。

下面对后边章节将要讨论的内容进行简要介绍。

第二章主要介绍了用到的技术和工具软件,对情绪分类模型的构建原理进行了简要介绍,主要讲了五种机器学习的分类算法,包括K近邻分类算法、SVM分类算法、朴素贝叶斯分类算法、决策树分类算法、Adaboost分类算法。

第三章详细介绍了实验过程中用户数据收集以及对收集来的原始数据的预处理操作,以及在构建情绪分类模型过程前进行的特征工程,包括了特征提取、相关性分析、特征选择、主成分分析等。

第四章介绍了情绪分析模型的构建，情绪分析模型主要包括K近邻分类模型、SVC分类模型、朴素贝叶斯分类模型三种。其中，K近邻分类模型中通过交叉验证选择出了准确率最高的K值，SVC分类模型中通过设置不同的核函数构建了三个不同的分类模型，朴素贝叶斯分类模型中按照对数据分布的假设不同设置了三个模型。

第五章对实验结果进行了准确性测试，输出了每个模型的平均准确率，并对测试集样本的最后十个数据的预测结果进行输出，然后与真实标签进行比对。本章同时对不同分类模型进行了性能测试，比较了处理相同数据集时不同的模型的耗时情况。

第六章对本文工作进行了总结，分析了本文工作的缺点与不足，并对后续工作的一些待完善与继续深入研究的方向进行了展望。

第二章 开发平台及原理分析

2.1 开发工具简介

本文中实现情绪分析系统使用的开发工具是Anaconda和PyCharm，接下来将一一介绍。

Anaconda是一个用于科学计算的Python发行版，它包括250多种流行的数据科学软件包，比如大名鼎鼎的numpy和pandas，以及适用于Windows，Linux和MacOS的conda软件包和虚拟环境管理器。Conda使安装，运行和升级复杂的数据科学和机器学习环境（如Scikit-learn，TensorFlow和SciPy）变得简单快捷。

PyCharm是一款开发者常用的Python集成开发工具，由Jetbrains公司开发完成。PyCharm除了提供了调试，语法高亮，代码跳转等一些基本功能，还提供了智能代码完成，代码检查，即时错误突出显示和快速修复，以及自动代码重构和丰富的导航功能，而且PyCharm内置了多种集成的调试器和测试运行器。PyCharm还与IPython Notebook集成，具有交互式Python控制台，并支持Anaconda以及多个科学软件包，包括Matplotlib和NumPy等。实际项目的开发过程中，通过使用PyCharm IDE可以极大地提升开发者的开发效率。

2.2 Numpy 和 Pandas

NumPy是用Python进行科学计算的基础软件包，它代表“Numeric Python”，是一个由多维数组对象和用于处理数组的例程集合组成的库。Numpy前身是Numeric，最初是由Jim Hugunin开发的，他同时开发了另一个包Numarray，这个包提供了一些额外的功能。2005年Travis Oliphant通过将Numarray的功能集成到Numeric包中来创建了NumPy包。

NumPy除了提供了矩阵数据类型、矢量处理等，还包含了精密的运算库等许多数值编程工具，因此十分擅长严格的数据处理任务[[endnoteRef:6]]。Numpy内部解除了Python的PIL(全局解释器锁)，运算效率极好，是大量机器学习框架的基础库。除了明显的科学用途外，NumPy还可以用作通用数据的高效多维容器，可以定义任意数据类型，在Python科学计算中应用十分广泛 [6: [] 姚建盛,李淑梅. Python在科学计算中的应用. 数字技术与应用,2016(11):76]

Pandas (Python Data Analysis Library) 是Python的一个数据分析包，最初由AQR Capital Management于2008年4月开发，并于2009年底开源出来，目前由专注于Python数据包开发的PyData开发团队继续开发和维护，属于PyData项目的一部分。Pandas是一种基于NumPy的工具，该工具诞生之初就是为了解决数据分析任务。Pandas中提供了高效地操作大型数据集所需的工具，其中包含了大量库和一些标准的数据模型[[endnoteRef:7]]。Pandas提供了大量能使开发者快速便捷地进行数据处理的函数和方法。Pandas里面还定义了Series和DataFrame两种数据类型，这使得数据操作变得更方便，并且Pandas可以轻松处理数据中的缺失数据，插入和删除数据，数据对齐，标签切片等等。此外，由于Pandas提供快速、灵活和富有表现力的数据结构，这使得它还可以处理很多不同类型的数据，比如具有异构类型列的表格数据，有序和无序的时间序列数据，以及任何其他形式的观测统计数据。 [7: [] 齐伟. 跟老齐学Python——从入门到精通. 北京: 电子工业出版社,2016]

2.3 scikit-learn

scikit-learn简称sklearn，于2007年问世，是Python最重要的机器学习库之一，常被用于机器学习和数据挖掘等应用中[[endnoteRef:8]]。 [8: [] Gavin Hackeling. Mastering Machine Learning with scikit-learn. Birmingham: Packt Publishing,2017.7]

sklearn依赖于matplotlib、NumPy和SciPy，内置了丰富的机器学习算法，有效提高了机器学习的效率。此外，sklearn内置的大量的标准数据集也为开发者节省了不少获取数据和处理数据的时间。而且sklearn文档完善，API丰富，上手难度小，颇受开发人员的喜爱。

sklearn库主要包含了分类（Classification），回归（Regression），聚类（Clustering），降维（Dimensionality reduction），模型选择（Model selection），预处理（Preprocessing）六大功能。分类，即识别对象属于哪个类别，包含的算法有K最近邻，支持向量机分类，决策树，朴素贝叶斯，随机森林等，常用于垃圾邮件检测，图像识别等方向。回归，最主要是预测与对象相关联的连续属性，多应用于药物反应，股价预测等方面，主要包括线性回归，多项式回归，支持向量回归等算法。聚类就是将对象根据数据间的特征划归为不同分类，常用方法有K均值、mean-shift等。跟分类不同的是，分类是用带标签的数据训练出模型，然后判断新数据哪种类别，聚类的数据是不带有标签的，是完全根据算法分析数据中之间的相似性来对数据进行自动归类。降维，就是降低样本特征的维度，常用来提高计算效率，或者进行可视化。模型选择即比较、验证、选择参数和模型，目标是通过参数调整提高模型精度。预处理则是对数据进行一些操作，如提取数据特征、归一化、标准化、白化、去均值化、二值化等操作来满足计算需求。

2.4 机器学习分类算法

2.4.1 K最近邻算法

K最近邻算法[[endnoteRef:9]] (K-Nearest Neighbors, KNN)工作原理：一组每个数据都带有标签的数据集，被称为样本集。样本集中的数据和其对应分类是已知的。输入样本集后，通过对样本集数据特征进行分析，训练出模型。等再输入不带标签的新数据时，提取出新数据的数据特征，与训练好的数据模型进行比对，从训练集中提取出K个与新数据最相似的样本的标签，选出这K个数据里面比例最高的标签作为新数据的标签，从而实现分类的目的。通常情况下，K的取值不大于20。 [9: [] Peter Harrington. Machine Learning in Action. New York: Manning Publications,

1988.8]

K最近邻的主要过程：

- 1、计算测试对象到训练集中每个对象的距离。
- 2、按照距离远近排序，这个距离可以是欧式距离，马氏距离，曼哈顿距离。
- 3、选取与当前测试对象最近的K个训练对象，作为该测试对象的邻居。
- 4、统计这K个邻居的类别的出现频率。
- 5、K个邻居里出现频率最高的类别，即为测试对象的类别。

K值选取较小时，也就意味着使用待分类点周围较少的邻居点进行预测，比如极端值K=1,那么也就是由待分类点最近的一个点来对该点进行分类，如果一旦该邻居点是噪声点，那么就会对预测结果造成较大误差。K值选取较大时，可以有效降低噪声的影响，但是容易导致分类界限不明显。K的具体取值一般与数据情况有关，不合理的K值会在一定程度上产生过拟合或者欠拟合的问题。具体的K值选择可以通过经验判断或者交叉验证来确定。

K近邻算法优点在于算法成熟，容易理解，准确率高，对异常值不敏感，既可以处理离散型数据，又可以处理连续型数据。缺点在于，样本分布不均匀时容易误分，数据量比较大时计算复杂度和空间复杂度都比较高，因此一般只应用于样本量比较小且分布比较均匀的情况。

2.4.2 支持向量机算法

支持向量机 (Support Vector Machine, SVM),本身是一种二分类模型，处理多分类问题时可以转换为一对多（选择某一类，样本中除这一类的其他所有类为一类，两者之间构建支持向量机）和一对一（任意两类之间构建支持向量机）来处理。支持向量机算法可以处理的数据可以分为三类：线性可分，近似线性可分，线性不可分。支持向量机的分类原理是寻找一个最优的超平面将数据进行分类，使边界最大[[endnoteRef:10]]。根据数据可分的三种情况，支持向量机有以下几种分类方式： [10: [] Nello Cristianini, John Shawe Taylor. 支持向量机导论. 李国正,王猛,曾华军译. 北京: 电子工业出版社, 2004.3]

1、如果数据线性可分，那么选择硬间隔方法使边界最大化，通过学习线性分类器来完成这一分类过程。硬间隔可以将所有样本正确分类，也正因为如此，受噪声样本影响很大。

2、如果数据近似线性可分，那么选择软间隔方法使边界最大化，通过学习线性分类器完成这一分类过程。软间隔对应于近似线性可分或线性不可分的数据集，允许一些超平面附近的样本被错误分类，从而提升了泛化性能。因为不是所有情况都需要把点全部分对，有时候样本点中存在一些本来就是错误的数据，也就是噪声，学习过程中如果学习了这些噪声，就会出现过拟合的情况，降低模型预测的准确性。这个过程需要加入惩罚因子C，使得点被错分的情况更合理。

3、如果数据线性不可分，那么就通过核函数以及软间隔最大化的方式，学习非线性分类器实现这一分类过程。当数据线性不可分的时候，SVM通过将低维数据向高维空间转化实现线性可分。

对于线性不可分的数据集的任意两个实例：，。选取特定映射f之后，使得与在高维空间中线性可分，运用上述（近似）线性可分问题的求解方法，可以发现目标函数和分类决策函数只涉及内积。由于高维空间中的内积计算非常复杂，通过引入核函数，内积问题变成了求函数值问题，由高维运算变成了低维运算，有效降低了计算复杂度。常用的核函数有：

多项式核函数：多项式核函数应用于将低维数据向高维空间转化，缺点是当多项式阶数比较高时，核函数计算复杂度会非常高，甚至无法计算。

式 (2-1)

线性核函数：主要用于线性可分的情况，其参数较少，运算复杂度低，通常选择线性核函数作为首先尝试的核函数来观察效果。

式 (2-2)

高斯径向基核函数：也是将低维数据转换到高维空间，但本身参数相比多项式核函数要少，所以在大样本和小样本时都有很好的性能。使用最广泛的一种核函数，在不确定用哪种核函数时，推荐使用高斯径向基核函数。

式 (2-3)

sigmoid核函数：

式 (2-4)

支持向量机算法的优点在于可以以较低的计算量解决高维问题，泛化错误率低，计算开销不大，结果容易解释，并且错误率较低。缺点在于，对参数调节和函数的选择敏感，对于核函数的高维映射解释力不强，尤其是径向基函数。

2.4.3 朴素贝叶斯算法

朴素贝叶斯算法一种简单而且效果比较不错的弱分类器，其理论基础是概率论中的贝叶斯理论。朴素贝叶斯算法虽然构建简单，分类效果却很优秀，甚至比许多复杂算法还要高效，尤其是在大型数据集，表现更佳。之所以说其朴素，是因为朴素贝叶斯算法是基于各个样本特征相互独立的假设的[[endnoteRef:11]]。举个例子，比如一个男生具有长得高，皮肤白，性格好的特点，可以得出结论该男生受女生喜欢，虽然可能这些特征之间具有一定的关联，或者相互依赖，但在朴素贝叶斯算法看来，这些特征在判断男生是不是受女孩喜欢的问题上，特征之间是相互独立的，并且对事件的影响是相同的，即权重相同。 [11: [] Peter Harrington. 机器学习实战. 李锐译. 北京: 人民邮电出版社, 2013]

先验概率：在事件发生前，基于历史事件统计，或者背景常识，或者人的经验判断得出的事件可能发生的概率。比如，天空中阴云密布，历史上天空中出现阴云密布的情况会下雨的可能性是70%，那么70%就是先验概率。

后验概率：在事件发生后，依据事件发生的结果反推该事件是由某因素引起的概率，即执果寻因。举个例子，中午吃了苹果，下午肚子疼，想算一下肚子疼是由吃苹果导致的概率，这就是后验概率。

贝叶斯理论：

事件X在事件Y发生的条件下的概率，与事件Y在事件X发生的条件下的概率是不一样的，贝叶斯公式就是用来描述这种关系的。

式 (2-5)

其中 $P(X|Y)$ 是在Y发生的情况下X发生的概率。X，Y都是事件，并且 $P(Y)$ 不为0。 $P(X)$ 是事件X发生的先验概率， $P(X|Y)$ 是X的后验概率。

当式 (2-5) 表示为式 (2-6) ，，y表示类变量，X是特征向量：

式 (2-6)

由朴素贝叶斯的朴素假设，也就是每个特征变量之间相互独立，即有：

式 (2-7)

所以，式 (2-7) 又可以表示为：

式 (2-8)

由于输入数据与分母是常量相关，于是进一步推导：

式 (2-9)



所以，目的是选择出类变量y的所有可能值的中使概率最大的那个可能值，可以用公式表示为：

式 (2-10)

最后，通过与的计算得出结果。

朴素贝叶斯分类器具有算法逻辑简单，实现难度小的优点，而且因为朴素贝叶斯分类器建立在特征相互独立的假设下，只涉及二维存储，所以分类过程中时间复杂度低。朴素贝叶斯分类器在样本特征比较少，并且特征之间相关性小时，具有较为良好的分类效果，而且理论上，相比其他分类方法，朴素贝叶斯算法误差率最低。但实际应用中，情况往往比较复杂，样本特征比较多时往往难以保证特征间相互独立的假设，分类效果也会受到影响。

2.4.4 决策树算法

决策树应用十分广泛，既可以用于回归也可以用于分类。分类问题中，决策树采用树形结构对样本的属性进行分类，既可以处理离散（if-then）的特征空间，也可以处理连续的特征空间（只需将连续空间通过阈值化变为if-then形式即可）。决策树由边和结点构成，内结点代表属性和特征，外结点代表类别[[endnoteRef:12]]。边代表判别的规则，即if-then规则。决策树主要分为特征选择、生成、剪枝三步。根据特征选择度量方式不同，分为ID3、C4.5、CART三种决策树算法，它们对应的度量方式分别是**信息增益（Information gain）**、**增益比率（gain ratio）**、**基尼指数（Gini index）**。[12: [] Willi Richert, Luis PedroCoelho. 机器学习系统设计. 刘峰译. 北京: 人民邮电出版社,2003.1]

信息增益：

首先是熵的定义，

式 (2-11)

变量的不确定性越大，熵也越大。

信息增益，即信息获取量（例如，通过a作为结点来分类获取了多少信息）：

式 (2-12)

依次比较各个**属性的信息增益**，**选择最大的那个属性**作为根结点，然后对于后面的结点，依次重复这个过程，直至最后给定**结点的所有样本属于同一类或者没有剩余属性**可以进一步划分样本，迭代停止。

增益比率：

式 (2-13)

式 (2-14)

其中代表的是a属性的x种情况中某一种在样本中的比例，即。增益比率，也就是某属性增加的信息熵与某属性自有信息上的比率。选择增益比率高的那个作为根结点。

基尼指数：

基尼指数（基尼不纯度）= 样本被选中的概率 * 样本被分错的概率

式 (2-15)

样本集分为K类，表示选中的样本在k类别中的比例，此处看作样本是k类别的概率，则这个样本被分错的概率是。

属性a的基尼系数，

式 (2-16)



选择基尼系数低的属性作为根结点。

决策树的优势在于它的数据形式非常容易理解，而且能够给出数据间的内在关系。除此之外，决策树计算复杂度不高，对中间值的缺失不敏感，可以处理不相关的数据特征。而且，决策树对数据形式要求简单，不必像其他分类方法一样统一数据属性，既可以是数值型，也可以是标称型。但缺点是容易拟合过度，处理连续变量效果不好，类别较多时，错误会增加的比较快等。

2.4.5 Adaboost算法

Adaboost算法，即自适应增强算法，是一种迭代算法[[endnoteRef:13]]。通过对同一数据集迭代训练不同的分类器，每次找到一个最优的分类器，然后下一次迭代时增大前一个分类器错误分类样本的权值，减小正确分类样本的权值。最后将得到的多个最优的分类器组合起来就得到一个强分类器。[13: [] Giuseppe Bonaccorso. Machine Learning Algorithms. Birmingham: Packt Publishing,2017.7]

Adaboost算法实现步骤：

- 1、初始化各样本数据的权值。假设有N个训练数据，第一次开始迭代时，各个样本被赋予相同的权值
- 2、对数据集迭代训练弱分类器。如果训练过程中，某个训练数据被准确分类，那么在下次迭代过程中，降低该训练数据的权值，同时提高被错误分类的训练数据的权值。一次迭代过程完成后，使用权重值更新后的训练数据集进行下一次迭代，构造新的弱分类器。如此迭代下去，完成整个训练过程。
- 3、集成各个弱分类器构建一个新的强分类器。为了让分类准确率高的弱分类器发挥更大的作用，按照分类过程中各个弱分类器的误差大小情况，为各个弱分类器分配权重。误差率越小的分类器，在构建强分类器的过程中，所占权重越高，否则，所占权重越小。这样，一个强分类器就构建完成了

Adaboost是一种简单有效的分类算法，很好地利用了不同弱分类器进行级联，并且在构建过程中充分考虑了不同分类器的权重问题，分类结果精度高。主要缺点有，分类精度可能会受数据不平的影响而下降，时间复杂度高，弱分类器的数目也就是迭代次数不易确定。

第三章 数据处理和特征工程

本实验主要研究内容是，以Android设备为信息收集载体，通过获取Android用户的运动信息、手机使用情况、用户情绪状态，并以此进行分析，探索发现用户的各项信息与情绪状态的内在联系，并建立模型，实现通过各项数据识别用户情绪。实验步骤主要包括：数据的获取、数据预处理、特征提取、PCA降维、训练模型与预测、参数调优、结果对比。

3.1 用户数据收集

本实验由10名使用Android设备的同学参与志愿活动，进行数据的收集。每日的数据收集分为早中晚三部分，在早上8：00，中午12：00，下午6：00，分别进行一次用户情绪录入，并与此同时进行离散采样，每隔一分钟进行一次采样，每次采样时间持续十分钟。志愿者的心情部分主要包含高兴、平静、难过、愤怒四种基本情绪类型，由志愿者凭感受主动录入。主要采样信息为Android设备各项传感器数据和设备基本情况信息。其中需要采样的传感器部分包括加速度传感器、方向传感器、陀螺仪传感器、磁场传感器、重力传感器、线性加速度传感器、GPS传感器、光线传感器、的数据。设备情况信息的采集部分包括手机的网速情况和情景模式。收集到的数据以txt格式存储，保存在用户手机内存卡根目录。整个采集过程持续两星期。

收集的各项信息数据格式如下：

（加速度x，加速度y，加速度z，方向x，方向y，方向z，陀螺仪x，陀螺仪y，陀螺仪z，磁场x，磁场y，磁场z，重力x，重力y，重力z，线性加速度x，线性加速度y，线性加速度z，GPS经度，GPS纬度，光照，网速，情景模式，情绪状态）

将收集到的原始数据转换成excel表格，情况如图3.2，图3.3，图3.4所示。

3.2 用户数据预处理

本实验中收集来的大量原始数据，是不能直接被使用的。收集来的原始数据有可能会存在包含大量噪声，数据不完整（比如有的属性可能缺失，或者不确定），数据不一致（比如在不同表中的同一属性名称不一致或者数据矛盾），度量单位不一致等问题[[endnoteRef:14]]。可能产生的原因有，比如收集过程中设备使用方法不正确，用户操作不当，或者设备出现故障和异常，信息收集过程受到干扰和中断，再比如数据收集后数据存储不当，工作人员误操作等等。这就需要通过科学的数据预处理技术对数据进行清洗，剔除坏值，非法值，异常值，填补空值等操作，来消除数据中的噪声，保证数据数据一致性及数据完整性，保证收集来的数据有效和可用。除此之外，还需要数据进一步进行集成，规约，变换等操作，使得数据满足研究分析的要求。[14: [] Jacqueline Kazil, Katharine Jarmul. Python数据处理. 张亮,吕家明译. 北京: 人民邮电出版社,2017.6]

数据的预处理可以显著提高数据的质量，同时可以有效地提高后续过程中数据分析的效率。Python中提供了强大的pandas和numpy库，本文作者使用这两个科学的数据分析库进行了数据的预处理操作。

3.2.1 数据集成

数据集成是把多组源数据融合成一组数据，这多组源数据可能来自多个不同的数据库或者不同的文件，所以集成的过程中要消除数据不一致和数据冗余。数据不一致主要表现为属性名称不一致或者数据矛盾。数据冗余一般是同一属性名称多次出现或者属性间存在线性关系。一般通过相关性分析来消除属性间线性相关。此处主要介绍相关的数据拼接和合并，有关属性间的相关性造成的数据冗余在后面单独会有介绍。

本实验由于是从多个用户收集信息，每个用户收集的信息保存在不同的文件中，所以需要多份数据合并在一起。本文作者首先将多份文件通过pandas读文件方法将txt文件读入，分别存储为不同的DataFrame（Pandas中的一种特殊数据结构，表现形式为二维数组），然后将不同的DataFrame拼接在一起。其中需要注意的是，防止数据不一致和数据冗余，不同文件的数据格式要保持一致，同一个属性在不同的文件中是否有不同属性名。因为数据分几次收集，前期收集到一部分数据后，又进行了数据收集工具的改进，数据格式存在一定的调整，在进行拼接时，本文作者使用pandas调整了属性顺序，然后进行拼接操作。

3.2.2 缺失值处理

通过numpy的isnan方法（numpy中的查询空值的方法，返回DataFrame的空值和非空值情况）发现，收集来的数据中明显存在一些缺失值，可能造成这一现象原因可能有用户忘记填写、应用被关闭或者后台清理、手机关机或没电等，也有可能是应用出现bug，未能正常收集数据。

如果有连续多个缺失值，应当采取的操作是将连续空缺的几组数据删除。如果是个别属性不连续的出现缺失，一般的处理方法有中位数替代法，平均值替代法，频率最高值替代法，默认值替代法，邻近值替代法或者根据需要直接删除属性缺失的行或者列。考虑到人的心情在一段时间内是相对稳定的，所以各项数据也应该是基本稳定的，在本实验中，本文作者采取同组的平均值进行缺失值的替换。

3.3 数据属性观察

通过对收集来的数据进行直观分析，初步判断用户数据中的GPS经纬度信息不具有明显变化，造成这一现象的原因可能有用户群体大多是在实验室进行毕业设计的大四同学，一天中待在实验室的时间较长，活动范围有限，基本上存在变化的时刻集中出现在中午就餐时刻，但由于手机内置的GPS传感器精度有限，无法明显捕捉到这样小范围的移动。然后，提取出每组数据中的GPS经纬度信息，使用MATLAB进行绘图，显示出的情况与人为判断基本一致，近乎集中于两三个点。鉴于这种情况，本文作者判断这两个属性对于后面模型构建以及预测分析不具备太大的参考价值，故选择直接删除这两项。

3.4 属性的相关性分析

本实验对获得的20个变量中，可能出现线性相关的变量之间，进行了相关性分析，以防止数据线性相关造成的数据冗余。

相关性的计算方法：

属性A和属性B的相关性计算，

式（3-1）

式（3-2）

式（3-3）

如果，则有A与B正相关。

如果，则有A与B负相关。

如果，则有A与B相互独立。

如果很大，那么就说明A，B相关性很强，可以删除其中一个。

3.5 加速度合成



本实验中使用到的加速度类传感器有加速度传感器和线性加速度传感器，每种传感器都有x，y，z三个轴，因为本实验中不需要区分具体的方向，所以把两类加速度传感器的x，y，z三轴的分量数据（即加速度x，加速度y，加速度z，线性加速度x，线性加速度y，线性加速度z六个属性），分别进行合成，合成后形成两个新的加速度属性——加速度和线性加速度，特征合成以后的操作都使用新特征进行操作，不再使用旧特征，并将旧属性从数据列表中剔除。将这几项特征合成后，使用合成后的值进行操作，可以有效减少训练模型时的时间复杂度，而且不损失精度。

加速度合成公式：

式（3-4）

其中，a代表合成后的加速度，，，分别代表x，y，z三轴上的分量加速度。

3.6 数据特征提取

直接收集来的数据特征可能数据关系不够明显，考虑到人的心情是由一段时间的一个平均状态来表现的，所以对收集来的数据，采用分箱技术进行切片分组，以10为单位分组，对每一组数据求出平均值、最大值、最小值、方差作为新的特征。

最大值：一组数据中其他值都小于等于数据中的某个值，这个值就是最大值。

最小值：一组数据中其他值都大于等于数据中的某个值，这个值就是最小值。

均值：即一组数据中的平均数，

式（3-5）

方差：一组数据中每个样本与该组数据平均值的偏离程度，

式（3-6）

3.7 主成分分析

经过前面的数据预处理和特征工程，本文作者已将原数据整理成了的样本集，为了体现整个过程的科学性和严谨性，也为了进一步从中已有的样本中提取对结果影响比较大的有关变量，减小无关变量和噪声的影响，以及降低计算量，本文作者对获得的新样本集进行了主成分分析操作。

通过调用sklearn.decomposition模块的PCA方法，生成一个PCA实例。PCA降维可以将数据降到指定维数，但是考虑到将数据样本降低越多固然会降低更多的计算量，但是也可能导致预测准确度过低的问题，所以一般降维时多选择指定降维后的最小精度（即保证降维后，新数据集保留的原数据集信息在指定精度以上）或者设置为“mle”，由PCA函数自动确定降低的维数。因为本实验中不知道降低的维数对于精度的影响，所以不容易指定维数，本实验中选择设置为“mle”方式，PCA函数自动选择最优的降维处理。经过PCA降维处理后，得到了新的60维样本集。

3.8 数据标准化

收集来的数据因为单位不统一、量纲不统一，是无法直接用来分析的。因为分析时往往不清楚各个属性对于结果的影响，所以一般假设各个属性对于结果的影响是相同的，即权重相同。如果量纲不统一，就可能导致样本之间数量级不统一，有的数据很大，有的数据很小，直接进行分析的话就会导致数值大的属性对结果影响过高，会影响到模型的准确性。

数据标准化主要用来消除量纲的影响，比如把属性按比例缩小，把属性放到一个特定的区间[[endnoteRef:15]]。本文作者收集到的数据包括运动类信息、手机状态类信息、环境信息三大类，其中运动类信息包含了加速度、方向、陀螺仪、重力等传感器数据，手机状态类信息主要包含了网速，环境信息主要包含了光强、磁场强度等数据，经过特征选择和特征提取后，形成了新的64维特征向量，即使降维后仍然有60维，其中的数据量纲存在巨大差异。在本实验中，本文作者假设每种特征都是与人的情绪状态相关的，并且对情绪状态的影响系数是相同的，为了防止某一特征所占的权重过高，从而过多的影响模型的构建，需要对数据进行标准化操作。在本实验中，本文作者选择了Z-Sorce标准化将数据进行了规范化处理。常见的标准化操作还有最大最小规范化，以及对定量特征进行二值化。

[15: [] 沈祥社. Python数据分析入门——从数据获取到可视化. 北京: 电子工业出版社, 2018.3]

1. Z-score标准化：

根据属性的均值和方差来对属性进行规范化，一般在最大最小化规范化出现异常数据时使用。



式 (3-7)

其中的和分别为属性A的均值和方差。

2. 最小最大规范化：

已知属性区间，将属性的取值范围由[old_min,old_max]映射到[new_min,new_max]

式 (3-8)

该方法保留了原来数据中存在的关系，但如果将来遇到超过目前属性[old_min,old_max]取值范围的数值，将会引发错误。

3. 对定量特征二值化：

对定量特征进行二值化之前，需要预先设定一个阈值a，如果比阈值大就赋值为1，如果比阈值小就赋值为0。

式 (3-9)

3.9 本章小结

经过前面的数据集成、缺失值处理保证了数据的一致性和完整性，通过属性的相关性分析消除了线性相关的变量，通过加速度合成将重要性不高的属性进行了合并，有效降低了数据的冗余和计算复杂度，特征选择和特征提取减少了无用属性，从原有数据中提取出了最能代表数据特征的属性集合，有效降低了噪声带来的影响，保证了后面模型训练的精度。本实验直接收集来的原始数据涉及23个属性，经过这部分的处理，融合成了新的64维特征向量，然后经过主成分分析进行进一步降维，降低运算的复杂度，得到了60维特征向量。

第四章 构建情绪分析模型

通过前面数据进行数据预处理操作，本实验得到了一个干净的、有效的样本集，经过特征选择和特征提取操作，获得了最能表现样本特征的特征属性，接下来就需要通过这些特征构建情绪的分析模型。从本质上说，这属于机器学习中有监督学习部分的分类问题，所以我们的目的也就是选择一种合适的分类器，将情绪准确地分类。

机器学习中的分类算法主要有K近邻算法、支持向量机算法、朴素贝叶斯算法、决策树算法、随机森林算法、Adaboost算法[[endnoteRef:16]]。本实验主要选取了K近邻算法、支持向量机算法、朴素贝叶斯算法构建情绪分析模型并进行分类。 [16: [] Ethem Alpaydin. 机器学习导论. 范明等译. 北京: 机械工业出版社, 2014.4]

4.1 K近邻分类模型

K近邻分类模型是基于K近邻算法构建的。K近邻算法原理是通过比较待测点与带有标签的样本点的距离，选择K个与待测点距离最近的样本点，统计这K个点中标签情况，选择比例最高的标签作为待测点的标签。

K近邻算法具有理论成熟、简单好用、测试准确率较高、对异常值不敏感等优点。而且考虑到K近邻算法要求数据量不能太大，否则会导致计算量过大，也不能数据量太小，这样会容易导致误分，本实验数据量刚好满足这样的要求，基于上述考虑，本实验中选用K近邻分类模型来构建第一个情绪分析模型。

K近邻算法中最重要的部分是K值的选择，不同的K值会对K近邻模型的分类准确度产生较大影响。通常情况K值不应大于20，而且为避免在通过K个标签分类时产生相同比例的标签，K值一般选择奇数值。具体的K值选择可以通过经验判断和交叉验证选择。本实验中，本文作者通过交叉验证选择K值。该步骤通过sklearn.model_selection模块的cross_val_score方法实现，设置cv参数为整数，使用KFold或StratifiedKFold的方法进行数据集打乱。通过设置K值的范围，本实验中将该范围设置为了1~31，然后比较不同K值下分类准确率的变化，并通过matplotlib模块的pyplot方法将不同K值下的分类准确率用折线图的方式绘制出来，选择出准确率最高的K值。通过交叉验证得到不同K值下，K近邻分类模型准确率变化情况如图4.2所示，可以看出K值为16时，模型的准确率最高，同时泛化能力较强。

然后，通过sklearn.neighbors模块的KNeighborsClassifier方法生成一个K近邻分类器。首先，将K值设置为前一步通过交叉验证获得的K值。然后设置K近邻算法的实现方式。K近邻算法的实现方式有枚举实现、KD树实现。枚举实现，即暴力实现，通过挨个搜索待测点距离每个样本点的距离，然后选出K个最近邻，这种方式计算量较大，只适合小数据样本。KD树实现方式没有直接计算待测点距离每个样本点的距离，而是先把数据存储进一个KD树，此处的K是指K个特征，根据构建好的KD树模型再进行距离的计算，可以有效减少计算量，提高分类的效率。此处本文作者设置参数algorithm为auto，即由分类器自动选择效率最高的算法。默认情况下，K个近邻的权重是相同的，但实际情况中可能不是这样的，通常距离待测点越近的样本点的标签越具有参考性，本实验中，作者按照K个近邻距离反比为K个近邻赋予权重。

经过前面步骤，已经获得了一个K近邻分类器，然后传入训练集X_train和训练集标签y_train，即可获得K近邻分类模型。

4.2 SVC分类模型

支持向量机算法的参数只与支持向量有关，数量少，这使得该算法在解决小样本高维度数据集的机器学习问题时具有其他分类器难以媲美的优势，并且支持向量机解决非线性问题有较好效果，泛化能力强，所以本实验中选择支持向量机算法进行了情绪分类模型的构建。

SVC分类模型是基于支持向量机算法构建的。支持向量机算法理论上仅支持二分类，而情绪具有多种类型，将数据划分为不同情绪类型就是多分类问题。支持向量机解决多分类问题的策略是将多个二分类器组合起来实现一个多分类器，本实验中作者使用了基于支持向量机实现的SVC（C-Support Vector Classification）来解决情绪的多分类问题。

首先，将收集来的数据按照80%，20%比例划分为训练集和测试集。

然后通过sklearn.svm模块的SVC方法生成一个SVC分类器。支持向量机解决多分类问题有一对一（one-against-one）和一对多（one-against-rest）两种策略。一对一策略是在任意两类样本之间构建一个支持向量机，本实验中有4种情绪类型，4个类别的样本之间就要构造6个支持向量机。然后对未知样本进行情绪分类时，选择得票数高的情绪类型作为未知样本的情绪类型。由于这种方式需要构建的支持向量机较多，影响了模型构建的效率，所以本实验中，作者选择一对多（one-against-rest）的策略构建情绪分析模型，即依次将某个情绪类型的样本归为一类，然后将剩余的类别归为一类，这样4个情绪类型的样本就构造出了4个SVM，分类时选择具有最大分类函数值的那类情绪作为待测样本的情绪类型。SVC解决多分类问题流程图如图4.4所示。

实验中首先选择了线性核函数作为SVC分类器的核函数，因为线性核函数简单，计算量小，通常研究过程中常使用这个核函数进行第一次实验。

线性核函数：

式（4-1）

因为考虑到情绪相关的数据可能比较复杂，简单的线性核函数可能并不能得到最好的分类效果，实验中还选择了多项式核函数和高斯径向基核函数对情绪数据进行分类。

多项式核函数：

式（4-2）

高斯核函数：

式（4-3）

最后，将训练集数据及训练集对应的标签传入构造好的支持向量机分类器中，对分类器进行训练，就得到了基于支持向量机的情绪分类模型。

4.3 朴素贝叶斯分类模型

朴素贝叶斯分类模型基于朴素贝叶斯算法，这个模型比较简单。朴素贝叶斯算法基于各个属性对于结果都有影响，并且权重相同的假设。根据计算方式的不同，朴素贝叶斯有三种实现，高斯朴素贝叶斯、多项式分布朴素贝叶斯、伯努利分布朴素贝叶斯[[endnoteRef:17]]。本实验中选择使用这三种朴素贝叶斯算法进行模型构建。[17: [] 米歇尔. 机器学习. 曾华军等译. 北京: 机械工业出版社, 2008.3]

基于朴素贝叶斯分类模型的主要构建步骤如图4.5，

首先，本文作者将收集来的数据按照80%，20%比例划分为训练集和测试集。然后，通过sklearn.naive_bayes模块的GaussianNB、MultinomialNB、BernoulliNB生成三个朴素贝叶斯分类器，接下来，将已经得到的训练集数据和训练集标签数据分别传入三个分类器，就得到了三种朴素贝叶斯分类模型。

GaussianNB假设特征的先验概率为正态分布，即如下式：

式（4-4）

其中为Y的第k类类别。和需要从训练集估算得知，为在样本中，所有的平均值。为在样本中，所有的方差。



MultinomialNB假设特征的分布为多项式分布，即如下式：

式 (4-5)

其中是第k个类别的第j维特征的第m个取值条件概率，是训练集中输出为第k类的样本个数。为一个大于0的常数，常常取为1，即拉普拉斯平滑。也可以取其他值。

BernoulliNB假设特征的先验概率为二元伯努利分布，即如下式：

式 (4-6)

此时m只有两种取值，只能取0或者1。

第五章 运行与测试

5.1 准确性测试

准确性测试主要是针对不同情绪分析模型，输入相同的样本，比较输出结果与真实结果。将前面划分好的数据集，取出测试集最后十个数据，如图5.1所示。

将取出的10个数据分别输入训练好的七个不同的分类模型中，并输出七个分类模型的分类结果，输出结果见表5.1。

对于训练集最后的10个数据，K近邻、SVC(linear)、SVC(poly)、SVC(rbf)、高斯朴素贝叶斯、多项式朴素贝叶斯、伯努利朴素贝叶斯六个情绪分析模型，分别识别正确了6个、6个、8个、9个、5个、8个、3个，与预期效果基本一致。对于整个测试集，不同的模型识别准确率情况如图5.2所示。

由图5.2可以看出，SVC分类模型（kernel= ' rbf ' ）准确率最高，达到了74%，接下来依次是多项式分布朴素贝叶斯分类模型，SVC分类模型（kernel= ' poly ' ），K近邻分类模型和SVC分类模型（kernel= ' linear ' ），高斯朴素贝叶斯分类模型，准确率依次为71%，66%，65%，65%，57%，准确率最低的是伯努利分布朴素贝叶斯分类模型，为33%。

5.2 性能测试

性能测试部分主要测试的是，对于同一个样本集，采用不同方式构造情绪分类模型并进行识别整个过程的耗时情况。训练集是一个的数据集，测试集大小为。构建不同情绪分类模型的耗时情况如图5.3所示。

由图5.3可以看出，效率最高的是K近邻分类模型和高斯朴素贝叶斯分类模型，出现这种情况的原因是，样本数量比较小，所以K近邻分类模型分类效率最高。接下来，效率由高到低依次是多项式分布朴素贝叶斯分类模型和伯努利分布朴素贝叶斯分类模型，SVC分类模型（kernel= ' linear ' ），SVC分类模型（kernel= ' rbf ' ），效率最低的是SVC分类模型（kernel= ' poly ' ）。

5.3 本章小结

通过对构建出的模型进行了准确性测试和性能测试，结果显示SVC分类模型（kernel= ' rbf ' ）分类准确率最佳，K近邻分类模型和高斯朴素贝叶斯分类模型分类效率最高，但考虑到样本量比较小，计算时间上的差距不大，都在可接受范围。所以，综合来说，SVC分类模型（kernel= ' rbf ' ）分类效果最佳。

第六章 总结与展望

本文针对智能手机日渐普及以及人们越来越关注自己的情绪的情况，通过收集Android用户数据信息，采用机器学习常见分类器如K近邻分类器、SVC分类器、朴素贝叶斯分类器构建了情绪分析模型，实现了对Android用户的情绪识别，来帮助人们随时随地了解自己的情绪状态。

本文提出了一种基于机器学习的安卓移动用户情绪分析方法，对于数据的收集过程、处理过程和分析模型的构建都进行了详细的分析和研究，并进行了准确性测试和性能测试。在数据收集过程，针对传感器的选取、收集的数据种类等方面进行了一定的分析。数据收集完成之后，对数据进行了详细的预处理，包括对异常值、缺失值产生原因的分析以及相关的处理，并对不同的原始数据集进行了集成。将数据进行了基本的整理之后，对数据进行了详细的特征提取和特征选择。此过程中，对属性进行了直观观察，进行了相关性的分析，加速度的合成。考虑到一些数据可能表达的数据关系不明显，本文作者又将数据进行了特征提取操作，为使数据满足后面计算需要，又对数据进行了归一化操作。由于新特征集维度比较高，为避免维度灾难，作者又对数据进行了PCA降维。然后，利用机器学习算法中三种常见分类器，通过调整参数使之适合已获得的样本集，然后训练出了三种情绪分析模型。最后，对这三种情绪分析模型进行了准确性测试和性能测试，验证的结果表明模型对于情绪的识别效果良好，大多数模型的识别准确率都超过60%，尤其是采用“ rbf ”核方法的SVC分类模型准确率最高，达到了74%。性能方面，K近邻模型和高斯朴素贝叶斯模型处理相同数据集明显耗时更少。综合准确性和性能来看，采用“ rbf ”核方法的SVC分类模型分类效果最好。

取得了一定成果的同时，研究仍然存在着一些不足。首先，数据的收集方面可以进一步改进，本文的研究主要基于运动数据和环境数据以及少量的用户使用数据，后面的研究可以多增加一些用户使用手机情况的数据，如打开某些APP的次数，以及使用时长，点亮屏幕的次数，使用APP的时间等。其次是，数据的特征提取仍存在改进空间，比如可以提高一些数据的采集频率，模拟成连续数据，然后对连续数据进行特征提取，还可以针对不同的数据进行不同频率的采样，并针对性的进行特征提取。另外，模型构建时，本文只选用了三种机器学习的分类器进行模型的构建，后续的工作中还可以选用更多的模型进行预测分析，比如随机森林、Adaboost等。并且，本文只对部分参数进行了调优处理，实际上每种模型可以调节的参数还有很多，还可以在参数调节方面进行改进提升，比如，增加调优的参数个数、改进调优的方法。

本文主要的工作集中在基于机器学习的安卓移动用户情绪分析系统的数据处理与模型构建上，由于本人的水平有限，在某些方面未能进行深入的探讨与研究，难免存在一些不严谨和谬误，敬请各位评委老师指导。

致谢

时光如梭，转眼间即将离开生活了四年的大学校园，这几年的学习和生活当中，从周围的老师同学、好友亲朋身上得到了很多的帮助，从学习上到生活上，从专业知识到为人处世，在此向所有帮助过我的老师、同学、家人、朋友致以最衷心的感谢。

首先感谢我的指导老师董洛兵老师，感谢董老师在实验过程中为我答疑解惑，耐心细致的帮我解决遇到的困难，指明前进的方向和道路，感谢董老师提供给我的这个学习的平台，为我的实验过程提供了很多便利，接触到了很多前沿的理论和知识，开拓了眼见，增强了自我学习能力。

然后，要感谢所有帮助过我的师兄和师姐们，在百忙之中还对我的毕业设计进行了耐心的指导，并在我毕业设计完成过程中向我提出了不少中肯的意见，使我在完成毕业设计的过程中少走了不少弯路，再次向他们表示感谢。

接着，我要感谢我的父母，一路走来是他们一直在背后支持着我，任劳任怨，默默无闻，失落时的鼓励，沮丧时的宽慰，过往的一幕幕无一不温暖着我的心。我能顺利完成本科学业，我的父母功不可没，真心地感谢他们为我付出的一切，作为子女，我只能更加努力的生活和工作来回报他们，真心祝愿他们身体健康。

最后，感谢所有在我学习和成长道路上帮助过我的人。

参考文献

