

移动应用隐私权声明内容合规性检验方法^{*}

王寅¹, 范铭¹, 陶俊杰¹, 雷靖蕙¹, 晋武侠¹, 韩德强¹, 刘炅¹



¹(西安交通大学 电子与信息学部, 西安 陕西 710049)

通讯作者: 范铭, E-mail: mingfan@mail.xjtu.edu.cn

摘要: 移动应用的隐私权声明作为用户与应用的协议条款, 是用户信息采集前必须向用户披露的关键文档。近年来, 国家出台多部政策法规明确要求移动应用需要配备清晰和规范的隐私权声明。然而, 如今隐私权声明存在诸多问题, 如缺失核心条目的披露, 省略信息采集的目的和使用模糊的表述等。另一方面, 随着法律条款数量增多, 条款间要求各不相同, 隐私权声明合规检验工作愈加繁重。本文提出一种移动应用隐私权声明的多标签分类方法, 这一方法通过比较四部核心法律法规对隐私权声明的要求, 总结梳理得到 31 类核心条目标签及特征。在该标签体系下, 本文设计实现了一个隐私权声明语句的分类模型, 该模型可以实现 94% 的条目分类准确率。基于该模型, 本文结合句法结构解析和实体识别方法, 在安卓应用和小程序场景中进行合规性检验, 发现 79%, 63% 和 94% 的隐私权声明分别存在条目缺失、目的省略和表述模糊问题。

关键词: 移动应用; 隐私权声明; 合规性

中图法分类号: TP309; TP311

中文引用格式: 王寅, 范铭, 陶俊杰, 雷靖蕙, 晋武侠, 韩德强, 刘炅. 移动应用隐私权声明内容合规性检验方法. 软件学报. <http://www.jos.org.cn/1000-9825/7121.htm>

英文引用格式: Wang Y, Fan M, Tao JJ, Lei JY, Jin WX, Han DQ, Liu J. Compliance Detection Method for Mobile Application Privacy Policy Content. Ruan Jian Xue Bao/Journal of Software. <http://www.jos.org.cn/1000-9825/7121.htm>

Compliance Detection Method for Mobile Application Privacy Policy Content

WANG Yin¹, FAN Ming¹, TAO Jun-Jie¹, LEI Jing-Yi¹, JIN Wu-Xia¹, HAN De-Qiang¹, LIU Ting¹

¹(Faculty of Electronic and Information Engineering, Xi'an Jiao Tong University, Xi'an 710049, China)

Abstract: The privacy policy statement of a mobile application serves as a crucial document that must be disclosed to users before collecting their information. However, current privacy policy statements face various issues, such as missing key disclosure items, omitting information collection purposes, and using vague descriptions. With an increasing number of legal provisions, the requirements for privacy policy statements vary, making compliance verification more burdensome. This article proposes a multi-label classification method for mobile application privacy policy statements. This method compares the requirements of four core laws and regulations regarding privacy policy statements, summarizes and organizes 31 categories of core item labels and features. Under this label system, the article designs and implements a classification model for privacy policy statement sentences, which achieves a 94% accuracy rate in item classification. Using this model, compliance verification was conducted in Android applications and mini-program scenarios, revealing issues such as missing items (79%), omitted purposes (63%), and vague descriptions (94%) in privacy policy statements.

Key words: mobile application, privacy policy, compliance detection

移动应用是当前使用最广泛、涉及用户隐私数据最多的程序, 是隐私保护的关键领域。近年来, 移动应用隐私泄露事件频发, 对国家安全、政治稳定和人民生命财产造成严重威胁。检测移动应用的数据收集和使

^{*} 基金项目: 国家重点研发计划资助项目(2022YFB2703501); 国家自然科学基金(62272377, 62232014, 72241433, 61721002, 62032010, 62002280); 中央高校基本科研业务费专项基金; CCF-蚂蚁隐私计算专项科研基金; 陕西省科学技术协会青年人才托举计划项目。

收稿时间: 2023-09-11; 修改时间: 2023-10-30; 采用时间: 2023-12-15; jos 在线出版时间: 2024-01-05

用是否符合国家标准与法律规定,是隐私保护中的关键问题。

隐私权声明是用于申明应用程序使用的数据类型和具体过程的一类文本,是用户了解程序行为的最直接途径。一篇清晰、完整的隐私权声明包括如何收集、使用、保存和保护用户数据等内容。然而,在实际环境中,隐私权声明可能会面临多种违规问题。以“我们会收集您的姓名、手机号码等个人信息”为例,该句隐私权声明提及了用户个人信息的范围,但是没有提及个人信息的使用目的,除此之外,也没有对个人信息的详细范围进行列举,而是使用了“等”进行模糊化表述。国家工信部自 2020 年起开展多次关于 APP 侵害用户权益行为的通报,通报多个主流应用商店均有上万个移动应用存在违规情况,如隐私权声明文档缺失、难以访问以及缺少信息采集处理的描述等。

针对移动应用中的隐私权违规问题,我国采取以移动应用运营商自律检查为主,政府监督监管为辅的模式,出台多部法律法规加以约束,如《中华人民共和国个人信息保护法》^[1](后文简称《个保法》)、《信息安全技术-个人信息安全规范》^[2](后文简称《安全规范》)、《APP 违法违规收集使用个人信息自评估指南》^[3](后文简称《评估指南》)以及《APP 违法违规收集使用个人信息行为认定方法》^[4](后文简称《认定方法》)等。然而,这些法律条例对隐私权声明的要求较为分散,如《个保法》中要求“公开个人信息处理规则,明示处理的目的、方式和范围”,但并没有解释信息处理规则的具体要求;《评估指南》则详细指出“隐私政策中应当将收集个人信息的业务逐项列举,不应使用‘等’、‘例如’字样”,但未提及《认定方法》中提到的不应该“模糊不清”。除此之外,在上述提到的法律条例中没有对表述模糊问题的定义、特征和判别规则进行探讨。这种法律条例间的在相同要求上的详略差异和不同要求的侧重差异增加了隐私权声明的违规判别的难度。

现有的隐私权声明检测分析研究主要可以分为完整性分析、模糊性分析和矛盾性分析。完整性分析指通过文本处理技术提取隐私权声明特征,进而检测隐私权声明中是否包含了用户关心的隐私权声明内容。模糊性分析指通过定义模糊性特征,来衡量隐私权声明表述的清晰性。矛盾性分析又可以分为文本自身矛盾和应用行为矛盾两类。文本自身矛盾指隐私权声明存在上下文矛盾的情况。应用行为矛盾指隐私权声明的描述内容与实际的应用程序行为不符。然而,这些研究并没有过多关注法律条例相关的要求,大多是简单地将隐私权声明表述划分为几种类型。例如 Costante 等人^[5]将隐私权声明划分为 6 种核心类别和 11 种附加类别,但是仍不够精细,没有对信息采集者是第一方还是第三方进行区分;Bhatia 等人^[6]将信息收集相关类语句进行拆解,区分了信息采集时的主体、客体、条件、目的、来源和目标,但是没有对其他类型的隐私权声明条目进行分析。这些研究的出发点是隐私权声明的质量评估,并没有紧密结合法律法规的要求,难以与当前开展的合规性检测工作有效结合。

隐私权声明合规性检测是移动应用隐私保护的重要环节,其难点在于:(1)合规性分析的检测方法应当对标于法律条例规定,然而不同的法律条例在内容和要求上都存在差异,且隐私权声明的涵盖范围广泛,条目种类复杂,难以通过少数标签完成系统性的隐私权声明的条目分类;(2)根据规范对文本要求的层级粒度不同,可以分为整体条目要求(如条目完整性)和某一具体条目的要求(如目的完整性),难以使用单一方法完成差异化粒度的合规性检测;(3)已有的隐私权声明文本分析工作都是围绕文本开展的,而文本和移动应用存在多种附加方式,不同的附加方式下文本的获取难度和途径存在差异。

针对以上问题,本文通过梳理归纳四部核心法律中对隐私权声明的要求,总结得到不同法律法规对隐私权声明的相同要求,以 100 篇规范应用程序的隐私权声明为参考,提出一种含 31 类标签的隐私权声明分类标准。通过人工标注隐私权声明文本,结合预训练语言模型,本文实现一种多标签的隐私权声明文本分类模型。针对倍受关注的条目缺失、目的省略和表述模糊三类隐私权声明违规问题,基于前述分类模型,以及对信息采集类条目语句的句法分析和实体标注,本文提出一种自上而下的隐私权声明合规检测方法,并在 APP 和小程序场景下的隐私权声明下进行了合规性检测与统计分析。

本文的主要工作和创新性贡献包括:

(1) 本文首次提出一种系统的隐私权声明分类划分标签体系,该体系严密贴合我国四部核心移动应用隐私保护法律法规,可以完整、准确地将隐私权声明结构进行分解。

(2) 根据分类标签体系, 本文构建了含 100 篇、共 86 万字的中文隐私权声明分类语料库, 并基于预训练模型实现了多标签分类模型, 该模型的平均准确率为 94%。

(3) 根据多标签分类模型, 本文进一步结合句法分析和实体分析方法, 提出条目完整性、目的完整性和表述清晰性的评估方法, 并在移动应用和小程序场景下开展合规检测和统计工作, 发现 79%, 63% 和 94% 的隐私权声明分别存在条目缺失、目的省略和表述模糊问题。该方法协助支持了监管人员进行隐私权声明评估, 部分违规检出结果通过提交给小程序平台得到了确认。

本文第 1 节介绍隐私权声明合规性分析的相关方法和研究现状。第 2 节介绍本文构建的基于多标签分类模型体系的隐私权声明内容合规性检验方法。第 3 节通过实验验证了分类方法的有效性, 并介绍在 APP 和小程序场景下的分析结果。最后总结全文。

1 隐私权声明分析相关工作

过去的隐私权声明合规性分析研究根据研究目标的不同, 主要可以分为完整性分析^[5-8]、模糊性分析^[4,9,10]和矛盾性分析^[11-19]三类:

(1) 完整性分析通常指通过文本处理技术提取隐私权声明特征, 进而检测隐私权声明中是否包含了用户关心的隐私权声明内容。Costante 等人^[5]提出一种结合文本分类技术和机器学习方法的隐私规则提取技术, 为用户提供处理后的结构化内容表示, 以此来评估隐私权声明的完整性。Bhatia 等人^[6]则从语义结构出发, 对语句中主体、客体、数据来源、条件和目的 5 种不同的语义角色以及它们之间的关系进行统计分析。Liu 等人^[7]总结了 GDPR 中要求的 10 类语句特征, 采用不同的分类模型检测隐私权声明是否包含这 10 类语句, 以此作为完整性结果。

(2) 模糊性分析通常使用几种模糊性特征, 来衡量隐私权声明表述是否清晰。[6]在滤除掉有具体指向的表述后, 将“其他”, “如果必要”, “任何人”等作为角色或者条件的模糊标志。[9]从自然语言的模糊范畴出发, 模糊性总结为不明确的条件、行为和类型的抽象与一般化、使用概率含义的情态和使用模糊的数量词四类。

(3) 矛盾性分析包括内部矛盾和外部矛盾两类, 内部矛盾通常指隐私权声明对数据的声明存在上下文矛盾的情况, 外部矛盾则是指隐私权声明对数据的声明与应用程序的实际行为存在不一致。Andow 等人^[11]使用句子级别的自然语言处理模型, 来捕获在数据收集和共享中的积极和消极陈述, 通过定义了 5 种矛盾和 4 种窄小定义的模式, 检测隐私权声明中数据共享和收集之间存在的矛盾。Slavin 等人^[12]提出了一个从 API 方法到隐私权声明短语的映射集合的半自动框架, 来检测隐私权声明是否正确地涵盖了程序行为。

然而, 这些方法存在如下问题:

(1) 对于完整性分析而言, 忽略了隐私权声明合规分析的层级关系。隐私权声明的合规性可以根据粒度不同, 分为条目间的合规性分析和条目内的合规性分析^[20-21]。已有的研究都是基于一种层级的合规性问题进行研究, 没有将二者结合进行完整性分析。例如, Bhatia 等人^[6]将完整性定义为单个语句是否包含了 5 种语义角色, 然而这种完整性评估仅对信息采集处理行为类的语句可行。要区分这类语句, Bhatia 等人采用了人工标注的方式, 这难以在大规模的隐私权声明文本分析需求中进行应用。

(2) 对于模糊性分析而言, 中文场景下的隐私权声明模糊性评估标准还存在空白。英文场景下的模糊性包括了概率含义的情态动词, 而中文没有情态动词, 而是以能愿动词表示可能性。除此之外, 隐私权声明中一些抽象实体在以往的工作中也被忽略了, 如隐私权声明经常使用“改善服务”、“完善体验”的表达方式, 这模糊了具体的服务形式。

2 隐私权声明内容合规性检验方法

2.1 多标签分类体系构建

根据法律法规的重要性不同,本文以2021年11月1日起施行的《中华人民共和国个人信息保护法》为

表1 隐私权声明标签设置规则及对应法律条文

首级标签	次级标签	标签设置说明	法律条文*
第一方收集/使用	个人信息类型	是否明确了运营者收集使用个人信息的规则,包括个人信息的类型、范围及相关功能业务	G7, A5, Z5, R2.1
	目的	是否逐一列出运营者收集使用个人信息目的、方式、范围等	G17, A5, Z6, R2.1
	用户选择	在个人信息的收集和使用过程中是否征得用户的授权同意,是否说明根据相关要求,运营者可能会无需征得授权同意的情况,是否说明若用户拒绝授权后相关信息的处理及功能的使用情况	G14, A5, Z23, R3.4
第三方收集/使用	个人信息类型	是否明确了运营者收集使用个人信息的规则,包括个人信息的类型、范围及相关功能业务	G21, A5, Z5, R2.1
	目的	是否逐一列出运营者收集使用个人信息目的、方式、范围等	G21, A5, Z5, R2.1
	权责/约束	是否明确了第三方介入个人信息收集的情形下,对第三方的权责说明或对第三方权限功能的约束细则	G21, A9, Z7, R5.3
第三方共享/转让/公开	SDK 信息	是否明确了第三方 SDK 的信息及功能,包括第三方 SDK 的主体、功能、收集个人信息的类型、方式、目的、范围等	G23, A9, Z22, R5.1
	公开披露	是否明确了运营者对外公开披露的方式、接收方身份信息、是否征求同意及特殊情况,以及涉及的个人信息的目的、场景、类型和对应的服务内容	G23, A9, Z14, R5.1
	共享	是否明确了运营者对外共享的方式、接收方身份信息、是否征求同意及特殊情况,以及涉及的个人信息的目的、场景、类型和对应的服务内容	G23, A9, Z14, R5.1
	转让	是否明确了运营者对外转让的方式、接收方身份信息、是否征求同意及特殊情况,以及涉及的个人信息的目的、场景、类型和对应的服务内容	G23, A9, Z14, R5.1
用户操作	用户可以进行的操作	是否说明用户是否具备对个人信息的访问、编辑、删除、查询、注销等操作的权利	G17/44, A8, Z15, R6.1
	操作的途径	是否说明用户在行使对个人信息的访问、编辑、删除、查询、注销等操作的权利时的具体操作方法或操作途径	G44, A8, Z15, R6.1
	运营者动作	是否说明当用户在行使对个人信息的访问、编辑、删除、查询、注销等操作的权利时或产品服务发生停止运营的情况时运营者的行为动作和对应所提供的服务	G22, A8, Z15, R3.8
数据安全	安全措施	是否包含保护个人信息的安全措施与技术,包括数据加密、进行安全评估/审核、对个人敏感数据的处理等	G9, A10, Z13, R5.2
	发生安全事故的处理方式	是否明确了个人信息发生安全事故时的处理方式及补救方案,包括为应对个人信息可能出现的风险而制定的制度及数据分类标准等	G57, A10, Z13
	跨境传输	是否明确了个人信息的存储地点及是否有跨境传输的问题,包含个人信息跨境传输的方式、目的和接收方等	G39, A8, Z12
	数据存储期限	是否明确了个人信息存储的时间期限及其依据,包括相关法规要求的特殊情况	G19, A5, Z10
Cookie 信息相关	Cookies 及同类技术	是否明确了运营者如何使用 Cookie 及其同类技术的方式方法等	G7, A5, Z21
	定义	是否包含了对某一名称或术语的概念或规定	/
	反馈渠道	是否明确了个人信息的安全投诉、举报渠道,包括运营者的反馈方式(办公地址、联系方式等)	G17, A5, Z16, R6.5
其他通用信息	免责声明	是否明确了运营者在满足某些情况下的责任说明	/
	目录	是否为隐私权声明文档的目录	/
	隐私协议适用范围	是否明确了适用该隐私权声明的产品和功能服务	G7, A5, Z20
	隐私权声明时效	是否明确了该隐私权声明的生效、失效或更新等时间	A5, Z17
	用户对条款的选择	是否明确了用户对该隐私权声明条款的选择方式	G14, A5, Z23, R3.4
面向特定人群用户条款	运营者信息	是否明确了运营者的名称、注册地址等信息	G17, A5, Z9, R6.5
	未成年人	是否明确了处理和保护未成年人个人信息的机制,包含收集未成年人个人信息的情况以及监护人的授权同意情况	G31, A5, Z8, R2.1
	告知方式	是否明确了隐私权声明修订和更新后告知用户的方式	G17, A5, Z18, R2.2
条款更改	更改原因	是否明确了隐私权声明修订和更新的原因	G17, A5, Z18, R2.2
	用户选择	是否明确了隐私权声明修订和更新后用户选择的方式	G14, A5, Z18, R2.2
无关语义句子	其他	是否为除以上细则内容外的描述文本	/

*G 表示《中华人民共和国个人信息保护法》,A 表示《信息安全技术-个人信息安全规范》(GB/T 35273—2020),Z 表示《App 违法违规收集使用个人信息自评估指南》,R 表示《App 违法违规收集使用个人信息行为认定方法》,/表示无对应条文。

主要参考依据,同时将《信息安全技术-个人信息安全规范》(GB/T 35273—2020)、《App 违法违规收集使用个人信息自评估指南》、《App 违法违规收集使用个人信息行为认定方法》等文件纳入参考的评价体系中。

首先,本文参考 Polisis^[22]中的分类方法,确定分类体系包括“第一方收集”、“第三方收集”、“访问与编

辑”等 10 类基本条目标签及基本特征。然后，安排三人阅读移动应用榜单排行榜“豌豆荚”前 100 名 APP 的隐私权声明，根据每类基本条目标签下不同语句的表达形式和表达内容，对照四部法律中的不同规则要求，在第一轮标注前确定主要的二级标签类别，以此标准完成第一轮标注。最后，对于第一轮标注中存在争议较多的标签类别，经过协商对标签类别进行微调，确定最终的二级标签并完成第二轮标注。例如，Polisis 定义了“第一方收集”的表达特征为包含“信息采集方式”、“信息类型范围”和“信息采集的目的”，经过本文标注者两轮标注与协商后，确定“第一方收集”的三个二级条目标签，即“个人信息类型”、“目的”和“用户选择”，其中“个人信息类型”和“目的”与 Polisis 原有特征对应，而“用户选择”则与四部法律法规中要求的“是否征得用户同意”相关。本文构建的多标签分类体系包括 10 个一级标签类别和 31 个二级标签类别，其标签细则和对应的法律法规条目如表 1 所示。

2.2 基于预训练模型的多标签分类方法

在模型数据集构建阶段，本文将 2.1 中的 100 篇隐私权声明文档作为多标签分类模型的数据集。首先，本文使用 Jieba 工具进行中文分词，使用常用中文停用词表中的符号部分，去除语句中的无意义符号，得到未标注的初始段落语料。然后，本文以表 1 为依据搭建了一个文本标注平台^[23]，如图 1 所示，标注人员通过导入未标记的初始语料，可以通过点选进行一级条目和二级条目选项进行语料标注。标注过程由三人经过两个轮次完成，在第一轮次，两人在无交流的情况下独立对语料进行标注，在第二轮次，两人对标签存在差异的语料条目进行协商，对于难以协商达成一致的语料文本，文本标签由第三人从争议标签中选择确定。最后，本文共标记得到 21204 条语句，本文构建的语料数据集的不同标签分布如图 2 所示。



图 1 隐私权声明文本标注系统

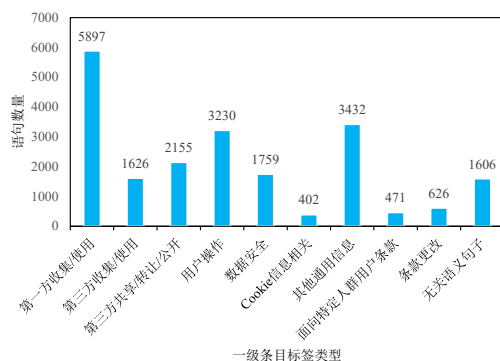


图 2 数据集中不同标签的分布

在模型训练阶段，本文首先使用预训练词表将语句转换为“[CLS]+文本”的形式，其中“[CLS]”为分类任务的特殊 token，如果文本内存在分割关系，用[SEP]代替分隔符。然后，在向量转换阶段，文本标注标签被转换为 1×31 维的特征向量，转换规则为对于向量中的第 i 维，如果语句标注结果中有第 i 个标签，则该位为 1，否则为 0。最后，预处理后的语句数据被输入给预先训练好的 BERT 模型，完成特征向量的分类预测任务。本文使用的预训练模型为 BERT-wwm-ext，预训练模型的结果通过 Relu-Dropout-Sigmoid 的网络层，获得最终分类结果。为解决训练数据中语句标签不均衡的问题，损失函数使用 Focal Loss 来衡量模型的分类结果。

在模型预测阶段，本文首先按照制定的正则表达规则，根据句号、问号、换行符等对文本 F 进行分段和分句处理，得到整篇文本分句后的子句集合 S。然后，对 S 中的每个子句，通过文本长度归一化 scaleLength 处理后输入到训练好的模型 M 中，得到该子句每一类标签的概率预测结果。最后，隐私权声明文本分类的最终结果 P 由每个子句的预测结果组成。

2.3 合规性检验方法

2.3.1 条目完整性

条目完整性指隐私权声明中应当包含法律法规要求列出的条目。在本文构建的多标签体系中，法律条文中明确要求指出的标签有 27 个，我们称之为“必要标签”。因此，当一篇隐私权声明中缺少这些“必要标签”的内容时，可以认为这篇隐私权声明是不完整的。本文首先使用多标签分类体系模型，将输入的隐私权声明执行分句和预测，得到所有语句的全部标签集合，然后计算该集合中缺失的标签数量，以及这一数量在全部必要标签中的比例。一篇隐私权声明的内容完整性评价值的取值范围为[0,1]，越靠近 1 表示该篇隐私权声明包含的必要标签数量越多，也就越完整。举例而言，如果一篇隐私权声明的所有分类标签结果缺少“其他通用信息”的首级标签和“反馈渠道”的次级标签，则该隐私权声明一共缺少 6 个“必要标签”，因此它的内容完整性评价值为 77.8%。

2.3.2 目的完整性

目的完整性指数据收集使用类别的语句在描述数据使用方式时，应当列出使用目的的情况。这里的数据收集使用类别具体指“第一方收集或使用”和“第三方收集或使用”。一句完整的数据收集使用语句至少包括数据类型、数据操作和数据目的三个要素，其中数据类型指收集使用的个人信息类型，如“手机号码”，数据操作指收集使用的方式，如“填写”、“传输”，数据目的指信息收集使用的目的，如“完成注册”。数据目的与数据类型和数据操作紧密相关，因此，在确定数据目的时需要同时确定数据类型和数据操作，本文通过构建词表和使用句法结构匹配分析来确定不同要素的内容。

根据要素划分，本文构建三类词表：（1）在数据类型词表构建环节，本文参考《网络安全标准实践指南-网络数据分类分级指引》，确定常用个人信息的基本范围，然后使用 Word2Vec 对语料库进行相似词语分析，最后经过人工整理得到数据类型词表。如表 2 所示，类型词表包含 80 类，共 1547 个不同的信息相关词语，如“真实姓名”、“性别”、“病史”等。（2）在数据操作词表构建环节，本文参考隐私权声明矛盾性分析相关工作，将数据操作归纳为收集、使用和传输三类，然后提取语料库文本中的全部谓词，最后经过人工整理得到数据操作词表。如表 3 所示，操作词表包含 3 类，82 个不同的操作相关词语，如“提供”、“读取”、“分享”等。（3）在数据目的提示词表构建环节，通过分析语料库中的介词、使役动词等用于提示目的的特定词素，本文构建了目的提示词词表，如表 4 所示，包含 143 个不同的目的提示词语，如“为了”、“以实现”、“出于”等。

表 2 数据类型词表示例

一级类别	二级类别	数量	词语示例	数量
基本信息	姓名,生日,地区,...	11	真实姓名,出生年月,住址,街道,...	352
身份认证	身份证,护照,驾驶证, ...	8	实名认证,有效身份证件,行驶证,驾照,...	72
生物信息	识别特征	1	指纹,虹膜,面部识别,掌纹,...	20
网络信息	账号,密码,照片,...	14	用户名,登录名,口令,密钥,相簿, ...	271
健康信息	身高,体重,病例,...	6	BMI,住院记录,就医信息,病史,...	62
工作教育	职业,薪资,单位,...	5	工作经历,公司,毕业院校,基本工资,薪酬,...	108
财产信息	银行账户,交易记录,信贷,...	6	订单,收入,支出,付款,贷款,...	183
社交信息	通讯记录,好友列表,...	5	聊天记录,通话记录,群列表,群聊,...	41
网络记录	日志,浏览记录,点击记录,...	5	崩溃,搜索历史,收藏列表,...	140
设备信息	手机型号,设置,应用信息	3	iOS, android, 安卓, 系统信息, ...	78
位置信息	位置信息,住宿	2	GPS,地理位置,定位,导航,...	99
硬件信息	摄像头,录音机,网络,...	11	WiFi,NFC,RAM, ...	68
其他信息	婚姻,快递,旅游	3	目的地,出发日期,收货人,收件物品,...	53
总计		80		1547

表 3 数据操作词表示例

操作类别	操作词语	数量
收集行为	提供,获取,读取,收集,输入,获得,记录,接收,扫描,...	33
使用行为	处理,展示,存储,支付,分析,调用,保存,搜索,拨打,...	42
发送行为	分享,传输,传递,上传,提交,传播,...	9

表 4 数据目的的定位词表示例

目的定位词语类别	目的定位词语	词语数量
为	为了, 为注册, 为登录, 为支持, 为改进, 为你, 为您, ...	46
以	以实现, 以创建, 以保证, 以完成, 以助于, 以保护, 以保存, ...	48
来	来使用, 来验证, 来优化, 来满足, 来提高, 来巩固, 来开展, ...	24
于	出于, 用于, 由于, 基于, 鉴于, 有助于	6
使	使我们, 使得, 使其能, 使您的, 使您能, 使您可以	6
特殊定位词	才能实现, 才能得以, 从而, 之目的, 的目的, 所必要, 所必需, 所需, ...	13

在句法结构匹配环节，本文首先使用 Hanlp^[24]对隐私权声明语句进行分词和词性标注，得到每一个分词的词性标注。然后对词性列表展开依存句法分析，得到语法结构树。最后，从语法结构树上抽取主谓宾关系对和介词词语，与前述构建的三要素词表进行匹配，得到要素关系对。当要素关系对中的目的成分存在缺失时，认为该句隐私权声明目的不完整。

以隐私权声明语句“我们收集手机号码是为了帮你完成注册”为例，经过词性标注和句法结构分析，可以得到如下关系对：（1）主谓宾关系对“我们”，“收集”，“手机号码”。其中“收集”为谓语，对应于数据操作中的收集行为，“手机号码”为宾语，对应于数据类型中的“电话号码”；（2）介词词语关系对“为了”，“完成注册”。其中“为了”为介词，其后的“完成注册”对应于数据要素中的数据目的。由此分析得知，该隐私权声明语句存在一个完整的数据三要素，即（电话号码，收集，完成注册）。

2.3.3 表述清晰性

表述清晰性指一句隐私权声明语句不包含表达不清或含糊其辞的表述，阅读者可以完整清晰地把握语句表达的全部信息^[25]。例如，“为了改善服务和提升用户体验”未提及具体的服务类型，以及“我们可能会收集您的身份证号码”中未明确说明是否会收集身份证号码，因此它们都是表述不明确的。为量化一篇隐私权声明的模糊性，本文参考了[9]中提到的模糊性定义，并针对隐私权声明这一特定语料环境，将常见的模糊性表述总结归纳为条件约束、泛化性约束、可能性约束、数量涵盖关系和抽象实体五类，如表 5 所示。

表 5 隐私权声明表述模糊类型与举例

表述模糊类型	关键字词	实例
条件约束	根据需要、视情况而定、如适用、必要时、有时	我们将 视情况 收取一定成本费用。
泛化性约束	通常地、普遍地、一般情况下、大多数情况下、多数情况中	一般情况下 ，我们不会与第三方共享儿童的个人信息。
可能性表达	可能、可能会、将以	您的个人信息 可能会 被转移到境外。
数量涵盖关系	最、一些、某些、大多数、多数的、部分、大部分、小部分、一部分、等、相关	向我们提供的相关个人信息，例如电话号码 等 ；
抽象实体	改善服务、提供功能、完善体验、保障业务、维护权益	目标是让你更便利地 体验 小金平台的 产品与服务 。

如图 3 所示，根据五类不同的模糊类型特征，本文具体分析方法如下：

（1）对于因条件约束和泛化性约束词语引起的模糊表述，首先使用正则匹配得到语句中对应类型的模糊词，然后通过否定词提取，结合否定语义判断是否存在模糊表述。图 5 中示例由于存在“不需要”这一否定词，因而不存在泛化性约束引起的模糊表述问题。

（2）对于因可能性表达引起的模糊表述，首先通过词性标注获得语句中不同分词的词性，然后对其中的能愿动词进行提取分析，确定语句中是否存在可能性表达。

（3）对于因数量涵盖关系引起的模糊表述，计算数量涵盖关系词语与数据类型在语句中的分词个数距离，如果距离小于阈值 θ ，则认为数量涵盖关系词语修饰了数据类型，也即存在数量涵盖关系模糊。 θ 的大小设置与信息类型词长度有关，当 θ 较小时，无法判断较长的信息类型词（如“电话号码”）使用数量涵盖关系词语的情况，因此 θ 应大于常见数据类型词语的分词长度；而当 θ 过大时，可能会判断错误（如“使用设备标识符完成广告的推送、推广和推荐等”中“设备标识符”与“等”不构成修饰关系）。综合常见表述形式，本文中 θ 取 5。

（4）对于因抽象实体引起的模糊表述，计算目的谓语与宾语间的分词个数距离，如果距离小于阈值 γ ，则认为宾语前没有足够的实体解释修饰词，即存在抽象实体模糊。 γ 的大小设置与修饰词长度有关，应大于常见无实体含义的修饰词（如“我的”、“我们的”、“正常的”等），同时也不能过大，避免滤除过多有实体含

义的修饰词（如“基于位置提供的”、“电话客服”等）。综合常见表述形式，本文中 γ 取 3。

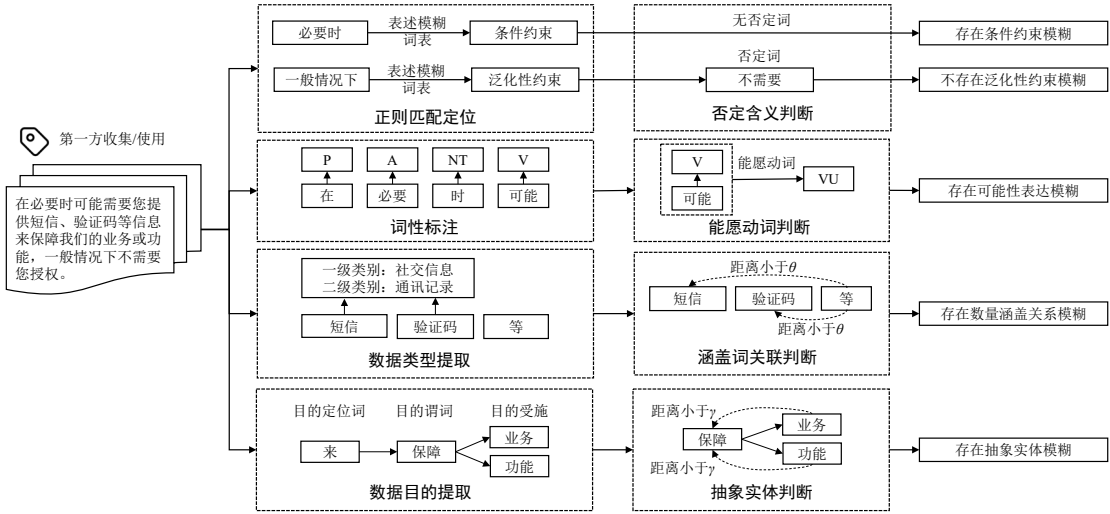


图3 隐私权声明语句的模糊性检验

3 实验分析

3.1 数据集构建

隐私权声明根据位置来源不同可以分为应用商店附加和应用内部携带两类。其中，应用内部附加又可以分为静态附加和动态附加两类：（1）应用商店附加指移动应用在提交发布至应用商店时，必须要提交的附加文本。这一分布多见于移动应用软件，即 APP。如图 4(a)所示，该类分布的隐私权声明可以直接在应用商店中按固定路径访问获取。（2）应用内部携带指用户实际使用时可以访问到的协议文本，该文本通过静态文本或以动态链接的方式添加在代码中。这一分布多见于小程序中。如图 4(b)所示，对于静态代码附加的文本，我们直接提取文本内容，而对于图 4(c)所示的动态附加型隐私权声明，我们采用动态测试的方式进行获取。

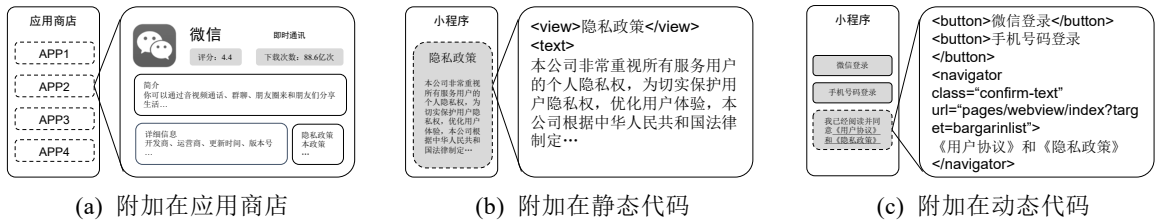


图4 隐私权声明在不同移动应用中的三种附加方式

针对小程序代码内部携带隐私权声明的爬取流程如图 5 所示。对于静态代码，首先，我们对小程序进行反编译，从源码中进行关键词匹配，当满足关键词匹配且单页面文本量超过阈值时，认为该页面文本含有隐私权声明文本。对于运行态中的动态链接，由于小程序采用 WebView 形式的展示结构，完全的遍历测试会产生某一页面测试无法跳出的情况，因此我们采用一种半固定的动态测试策略。首先访问小程序的不同导航标签，然后在不同页面内进行随机点击，最后使用关键词匹配方法检测当前页面场景，包括隐私权声明页面和进入隐私权声明的中继页面。对于访问的任意一个页面，当页面包含多个隐私权声明特征词如“隐私”、“政策”且文本量超过设定的阈值时，认为该页面文本含有隐私权声明文本。

我们在动态测试时使用动态测试工具 Airstest 进行测试，该工具结合了图像识别框架，可以识别页面内的组件文字信息。如图 6 所示，当我们进入小程序首页时，它会弹出“协议规则”窗口，其中包含了“相关法

律条款及隐私政策”。动态测试工具会返回该元素的文字内容和文本框的四角坐标。我们首先计算整个元素的实际页面长度和文字个数，然后计算得到“隐私政策”的相对位置坐标，最后点击这一坐标实现页面跳转。在样例小程序中，该页面文本数量没有超过设定文本字数阈值，即认为这一步骤没有真正访问到隐私权声明文本，因此我们会在新页面内尝试随机点击，直到访问到真正的隐私权声明页面。在隐私权声明页面中，我们采用模拟向下滑动的方式来捕获完整的隐私权声明文本，直到页面截图内容不再发生改变。

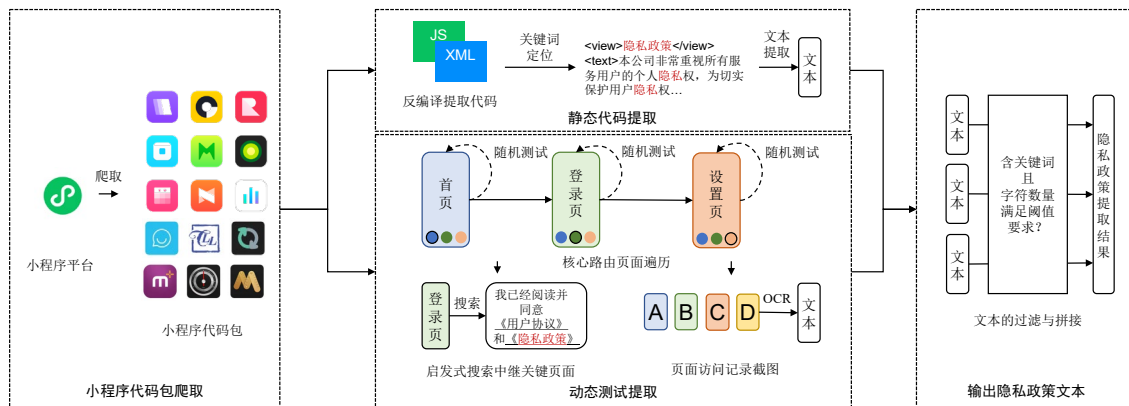


图 5 小程序中的隐私权声明爬取流程

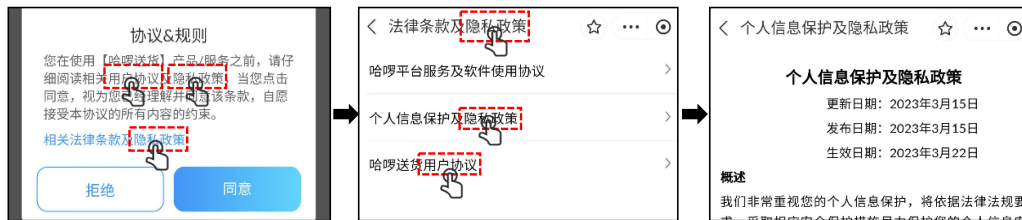


图 6 隐私权声明爬取示例

我们选取了一定数量的流行应用小程序，在小样本上检验我们的爬取方法。小程序缺乏可以提供访问量的应用市场，因此我们首先爬取了 App 排行榜前 400 的应用名称，然后在小程序平台检索对应的小程序，手动爬取隐私权声明文本作为基准数据集。在构建过程中，400 个 App 应用名有相应小程序的有 162 个，其中能手动访问到隐私权声明的有 88 个。经过动态测试爬取实验，有 58 个小程序能够成功爬取得到隐私权声明文本，平均耗时 190.6 秒。爬取失败的原因有：(i) Airstest 定位到的页面组件不全或坐标不准确，导致无法通过特定组件进入隐私权声明页面；(ii) 部分小程序访问隐私权声明路径过长，在进入过程中容易陷入循环而难以跳出；(iii) 少量小程序的隐私权声明嵌入方式为非网页形式而是文件形式，小程序在加载时渲染的页面文字无法被 OCR 工具正确识别。

通过以上方法，本文共爬取了不同类型应用共 22961 篇隐私权声明文本，涵盖“影音播放”、“系统工具”、“网上购物”等 16 个类别，如图 7(a)所示。在合规性检验章节，我们选择“金融理财”、“健康运动”和“旅游出行”三个类别的隐私权声明文本进行分析。相较于收集使用用户信息较少的应用，如“生活休闲”、“网络游戏”和“影音播放”等，这三类应用收集用户信息更加频繁，在隐私敏感问题上更具有代表性。其次，我们通过少量样本的预先人工分析过程中发现，不同应用间的违规情况和比例大致相同，例如内容缺失的标签分布情况大多呈现类似分布。限于篇幅原因，最终，我们在三类应用中每类随机选出 400 篇文件大小在 5 到 80KB 的隐私权声明文本，数据集的文字数量和分布如图 7(b)所示。

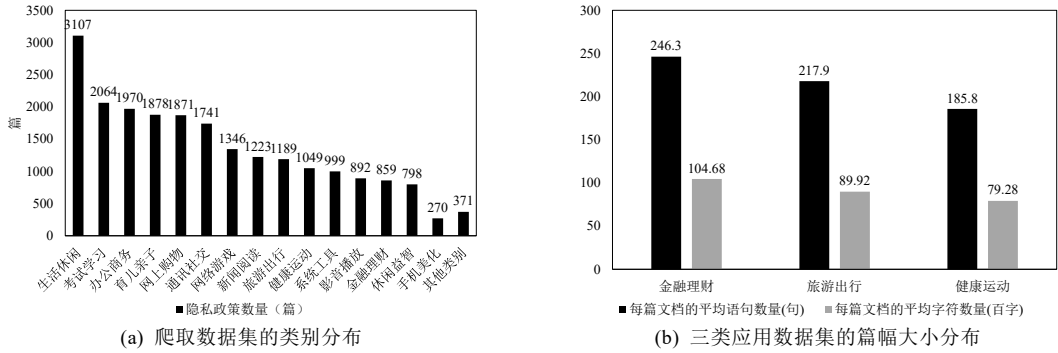


图 7 数据集的类型分布与篇幅大小分布

3.2 研究问题

本文针对多标签分类体系提出一种基于预训练模型的分类方法,该方法的效果决定了能否顺利完成合规性检验工作。因此,我们关心分类方法能否有效地完成预定的多标签分类工作。另外,本文针对内容完整性、目的缺失性和模糊性表述三类违规问题提出了检验方法,我们关心本文定义的问题在实际样本中的分布情况,以及产生这一分布现象的原因,以提出相应建议,帮助相关人员理解和解决这些问题。综上所述,本文围绕以下四个研究问题开展实验:

RQ1 本文提出的隐私权声明语句多标签分类方法的效果怎么样?

RQ2 在真实环境中,隐私权声明的内容完整性违规现象有多严重?

RQ3 在真实环境中,隐私权声明的目的缺失违规现象有多严重?

RQ4 在真实环境中,隐私权声明的模糊性表述违规现象有多严重?

3.3 实验方法与结果

3.3.1 多标签分类模型评估

在训练阶段,我们参照[26]中的参数设置进行初始值选择。为实现更快的收敛速度与更优的分类效果,基于本文多标签分类实验的特点,本文中分类器的初始学习率设置为 1×10^{-5} ,优化器使用 Adam 算法,损失函数使用 Focal Loss 函数,迭代次数设置为 30,Dropout 设置为 0.1。由于本文使用的预训练模型最大输入处理长度为 512,因此在处理长文本的过程中使用截断法将文本长度缩短至 512,并对过长的文本进行二次分句和分句分标签预测。训练后的模型通过十折交叉验证方法计算精确率和召回率。表 6 为首级标签的评价效果。根据实验结果,本文有如下结论:

(1) 基于预训练模型的多标签分类模型可以实现隐私权声明语句的有效分类。当分类目标为首级标签时,该分类模型的微平均精确率为 94.0%,宏平均精确率为 95.0%。

(2) 不同标签的语句分类效果与不同标签语句的文本特征紧密相关。(i)在所有的首级标签中,分类效果最好的为“Cookie 信息相关”标签,原因是该类别的文本表述普遍带有特征词“Cookie”,并且不同的隐私在这一部分的表述相似度较高,例如“我们不会将 Cookie 用于本政策所述目的之外的任何用途”。(ii)分类效果最差的为“无关语义句子”,原因是该类别包括了全部无法被划分到其他标签下的语句,这些语句难以被划分入其他特征明显的标签类型,自身内容也无专一特征。例如,表示引导性阅读的“我们谨此再次提醒您,本协议内容中以加粗方式显著标识的条款,请您着重阅读”。

结论 1. 本文提出的多标签分类模型在首级标签和次级标签的平均精确率均在 94%以上,该模型可以实现隐私权声明语句的有效分类。

表 6 10 个首级标签的多标签分类模型评价指标

首级标签	精确率	召回率	F1 值
第一方收集/使用	0.951	0.934	0.942
第三方收集/使用	0.901	0.842	0.871
第三方共享/转让/公开	0.933	0.878	0.905
用户操作	0.892	0.812	0.850
数据安全	0.887	0.852	0.869
Cookie 信息相关	0.982	0.961	0.971
其他通用信息	0.924	0.883	0.903
面向特定人群用户条款	0.974	0.968	0.971
条款更改	0.919	0.841	0.878
无关语义句子	0.835	0.689	0.755
微平均	0.940	0.917	0.928
宏平均	0.950	0.927	0.938

3.3.2 条目完整性

本文使用训练得到的多标签分类模型，对三类应用共 1200 份隐私权声明开展完整性分析，不同类别隐私权声明的完整性分布情况如图 8 所示。根据实验结果，本文有如下结论：

(1) 隐私权声明存在内容缺失的现象严重。如图 8(a)所示，三类应用隐私权声明完整性为 100%的数量仅为 250 篇，占全部隐私权声明的 20.8%。其中，“金融理财”类应用的隐私权声明相对完整，有 120 篇的隐私权声明是完整的，占该类隐私权声明的 30%；“健康运动”类应用的隐私权声明内容缺失最为严重，完整性在 80%的仅有 248 篇，占该类隐私权声明的 30%，低于“金融理财”的 315 篇，78.8%和“旅游出行”的 286 篇，71.5%。由此可见，不同类别的隐私权声明完整性存在差异，“金融理财”类应用对信息处理和保护较其他类型应用更为敏感。

(2) 隐私权声明内容完整性与隐私权声明文本长度在高完整指标区间内存在正相关。如图 8(b)所示，在 40%-100%区间，完整性指标越高，字符数越多，语句数越多。而低完整性的隐私权声明文本数量较少，平均长度指标容易受到个别样本影响。在不同类别的隐私权声明中，最完整的“金融理财”类的文本平均语句数为 246.3 句，平均字符数为 10468 字，最缺失的“健康运动”类的文本平均语句数为 185.8 句，平均字符数为 7928 字。

(3) 在不同类别应用的隐私权声明中，内容缺失现象集中在存在重合的某几类标签上。如图 8(c)所示，隐私权声明通常都有的标签集中在“第一方收集/使用”，在“金融理财”和“旅游出行”类应用中标签缺失最为严重的是“Cookie 信息相关-Cookies”。

结论 2. 在测试的三类应用共 1200 份隐私权声明中，存在内容缺失的有 950 篇，占 79.1%，内容完整性与隐私权声明文本在高完整指标区间内存在正相关，不同类隐私权声明中缺失标签呈现出集中和重合的现象。

3.3.3 目的完整性

本文通过数据要素提取方法，对“金融理财”类应用的隐私权声明进行要素提取，其中数据目的分布情况如图 9(a)所示。在划分目的缺失类型上，由于数据目的与数据类型紧密相关，本文进一步将目的缺失细分为四种类型。其中，“没有缺失”指隐私权声明对全部声明的数据类型都做了使用目的的阐述；“部分缺失”指隐私权声明提到了一部分类型的使用目的，但是仍有一部分类型未声明使用目的；“全部缺失”指隐私权声明中没有对任何数据类型进行使用目的的阐述；“未声明类型”指隐私权声明中未说明收集任何个人信息。对“金融理财”应用的目的缺失检验结果如图 9(b)所示，图中数字表示对应缺失类型的隐私权声明数量。根据实验结果，本文有如下结论：

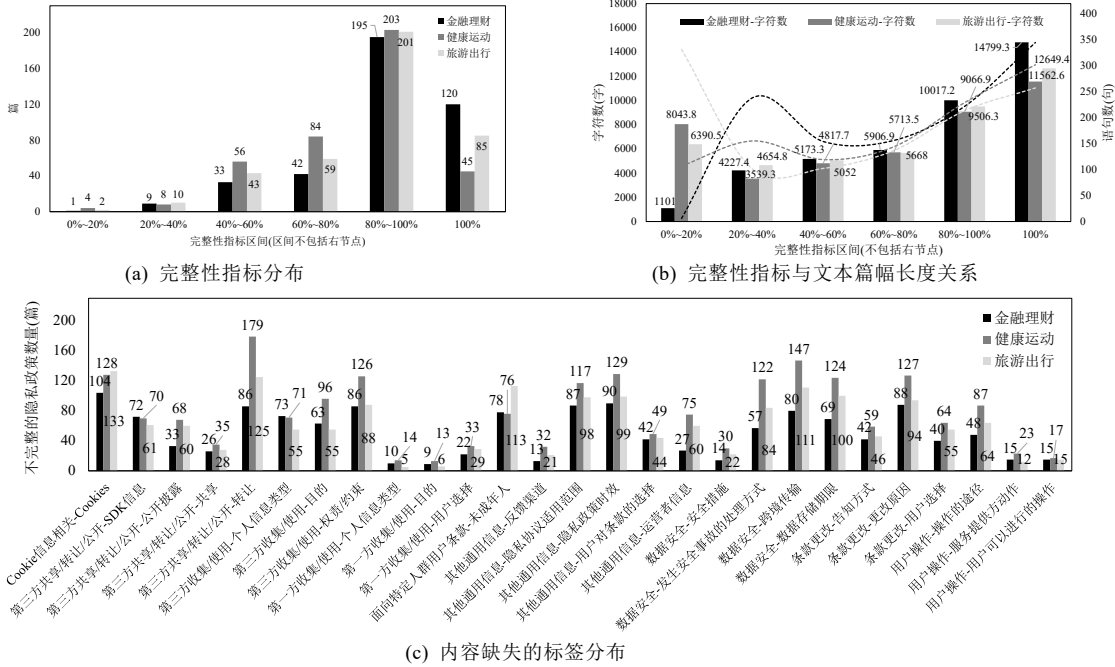


图 8 三类应用的隐私权声明条目完整性分布情况

(1) 隐私权声明的数据目的类型分散，常见目的的表述笼统模糊。超过 200 篇隐私权声明中均出现的目的类型只有一种“提供服务”，该目的也不指向确切的服务类型，其余的数据目的在不同隐私权声明中出现的频率均低于 100 次。

(2) 隐私权声明的目的缺失现象严重，收集信息而不告知使用目的的现象较为普遍。测试数据集中，63%（253 篇）的隐私权声明都存在目的缺失现象，数据目的表述完整的隐私权声明仅占约三分之一，其中包括了 2%（6 篇）的未声明任何数据类型的隐私权声明。由此可见，移动应用在声明获取用户信息过程时，更倾向于遗漏或隐瞒真实的使用目的。造成这一现象的原因是完整的隐私权声明的撰写需要配备专业法律知识，同时缺乏强制、明确的隐私权声明撰写规范标准，这造成了许多隐私权声明撰写人员会有意或无意地遗漏关键信息。

结论 3. 在金融理财类应用的 400 份应用中展开测试，存在目的缺失的有 253 篇，占 63%。隐私权声明更倾向声明数据类型和数据动作，而忽略了对使用目的做出解释。

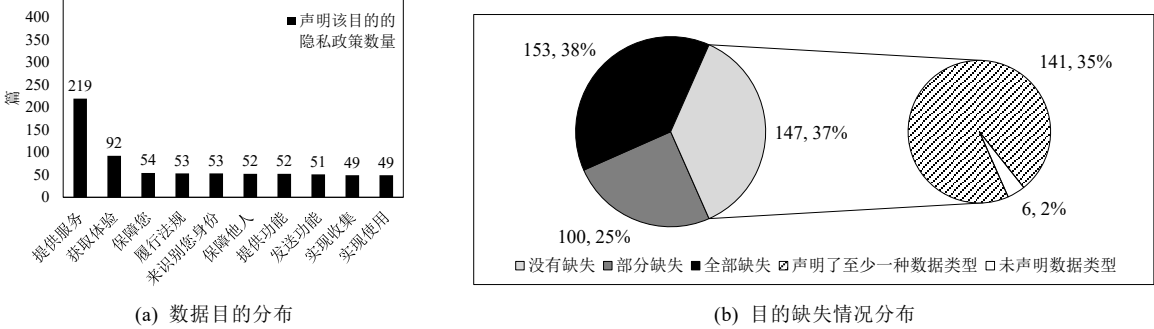


图 9 数据目的分布与目的缺失情况分布

3.3.4 表述清晰性

本文对 400 篇“金融理财”类应用的隐私权声明展开表述模糊分析，实验结果如图 10 所示。根据实验结果，本文有如下结论：

(1) 隐私权声明中存在普遍的表述模糊问题。如图 10(a)所示，在测试样本中，完全清晰的隐私权声明数量仅有 24 篇，占总样本的 6%。隐私权声明的模糊性语句集中在 1-5 句和 6-10 句区间内。超过 25 篇隐私权声明中存在 25 句以上的模糊语句。

(2) 隐私权声明的模糊性表述主要集中在可能性表达、数量关系涵盖与抽象实体类型上。如图 10(b)所示，(i)只有 31 篇隐私权声明中出现了条件约束，2 篇隐私权声明中出现了泛化性约束，原因是隐私权声明中该类模糊性表述往往和否定词连接在一起使用，例如“多数情况我们不会收集您的额外信息”，尽管存在模糊词，但没有隐藏收集使用用户信息，因此不构成完整的模糊表述；(ii)而另外三种表述非常频繁，超过 336 篇，84% 的隐私权声明存在至少一类该种模糊，原因是隐私权声明相比于一般文本语料而言，会更频繁地出现数据类型列举和抽象实体表述，这增加了出现模糊表述的概率，然而从法律法规角度而言，这些模糊性表述损害了用户的知情权。

结论 4. 在金融理财类应用的 400 份应用中展开测试，存在模糊表述的隐私权声明占 94%，模糊性表述集中在可能型表达(346 篇，86.5%)、数量关系涵盖(337 篇，84.3%)和抽象实体(338 篇，84.5%)类型上。

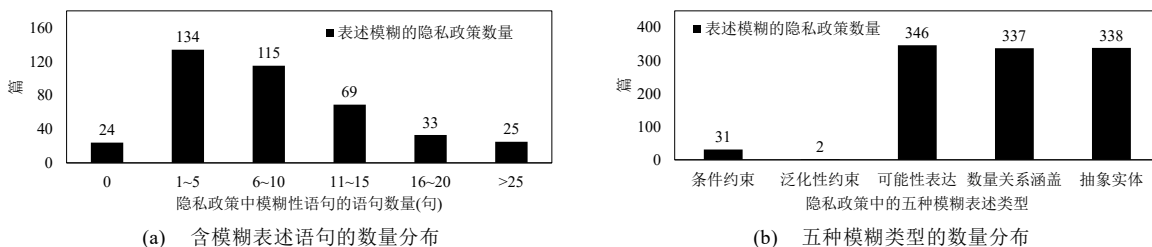


图 10 隐私权声明模糊表述分布情况

3.4 现象讨论

通过实验，本文认为造成当下隐私权声明不合规现象严重的原因主要有以下三点：

(1) 当前仍缺乏标准统一的隐私权声明规范。(i) 现如今我国已经颁布了多项法律法规，对隐私权声明提出了内容方面的不同要求，然而，目前仍存在法条较多、不同要求粒度不同的问题。本文在文本标签类别上进行了统一规范的尝试，发现除去法律法规关注到的标签类别以外，仍存在一些未被法律法规关注到，而在隐私权声明文本中普遍存在的标签，如术语的“定义”、“免责声明”等。(ii) 对于目的缺失和表述模糊的问题当前法律法规中也没有可以量化的定义，更多的是从定性的角度进行描述，较多使用如《安全规范》中“知情”“清晰”“透明”等词语，而较少使用如《认定方法》中提到的“点击少于 4 次”“未提供简体中文版”等可量化的描述。(iii) 现有的法律法规体系大多是面向移动应用整体的行为进行约束，而非隐私权声明文档的专有规范。

因此，本文建议监管者加快隐私权声明内容规范的制定与推行工作，为隐私权声明的内容合规在顶层设计上提供标准规范支持。

(2) 时下开发者欠缺隐私权声明撰写需要的背景知识。(i) 当下对移动应用行为进行约束的条例文件已经超过十余部，且可以预见的是，随着国家在数据立法的不断推进上，这些条例数量将会越来越多。能够完全理解，并将之付诸于隐私权声明撰写的合规自检工作，需要开发者具备较强的相关领域知识。(ii) 通过实验我们发现，时下不同应用中大多采用《安全规范》中给出的隐私权声明模板，结合自身应用进行适配和修改。然而，应用功能开发与隐私权声明撰写是两种截然不同的任务。如果由应用开发人员来撰写隐私权声

明,则可能会在表述过程中由于语言习惯问题而较多使用“等”来涵盖信息类型,违背“逐一列出”原则。如果分离两种任务则存在两部门人员的交流问题,拉高移动应用的开发成本。

因此,本文建议开发者可以配备专业隐私权声明撰写人员,或是尽快推出基于程序代码分析的隐私权声明自动生成工具,以减少开发人员可能造成的违规隐患。

(3)目前隐私权声明用户友好程度较低,用户对内容的关心程度也不高。在本文实验过程中发现,当下隐私权声明普遍存在篇幅过长、内容繁杂问题。尽管一些移动应用已经通过一些方式改进了用户阅读体验,如提示要点、使用列表形式等,然而大部分用户很少会去阅读隐私权声明的内容,也很少关心隐私权声明的更新情况。在使用便利性和隐私保密性上,用户更倾向于选择前者而放弃后者。当隐私权声明出现内容违规现象时,用户无法及时发现其中存在的问题,在一定程度上也滋生了隐瞒要权、模糊要权的违规现象发生。

因此,本文建议一方面可以改善隐私权声明阅读的用户友好程度,鼓励、帮助用户了解隐私权声明的核心内容,另一方面用户可以提高隐私忧患意识,对陌生应用在敏感信息获取和使用上提高警惕。

4 总 结

移动应用中的隐私权声明是用户了解应用具体行为的桥梁,应当向用户清晰有效地披露数据收集和使用细节。隐私权声明合规性是移动应用隐私合规中首先考虑的关键问题,本文提出了一套隐私权声明合规性检测方法,该方法从隐私权声明获取出发,收集构建了含100篇、共86万字的隐私权声明分类数据集,基于这一数据集,本文构建了一个多标签分类模型,可以自动地实现隐私权声明语句解析,进而实现内容完整性分析、目的缺失检验和表述模糊分析。实验结果发现,本文提出的多标签分类模型在首级标签和次级标签的平均精确率均在94%以上。通过对三类应用共1200份隐私权声明的统计分析,本文发现仅有20.8%的隐私权声明是完整的。在400份金融理财类别的隐私权声明中,63%存在目的缺失现象,94%存在至少一句表述模糊的语句。移动应用中的隐私权声明违规问题十分严重。

References:

- [1] Standing committee of the National People's Congress. Personal Information Protection Law of the People's Republic of China. National People's Congress,2021.
- [2] National Standardization Administration. Information Security Technology - Personal Information Security Specification: GB/T 35273-2020. National Standard Full Text Disclosure System,2020.
- [3] Personal Information Protection Task Force on Apps. Self evaluation guide for illegal collection and use of personal information on apps. Personal Information Protection Task Force on Apps,2019.
- [4] Cyberspace Administration of China, Ministry of Industry and Information Technology PRC, Ministry of Public Security PRC, and State Administration for Market Regulation. Identification Method for Illegal Collection and Use of Personal Information by App. China NetEase,2019.
- [5] Costante E, Sun YH, Petkovi M, Hartog J. A machine learning solution to assess privacy policy completeness. Proceedings of the 2012 ACM Workshop on Privacy in the Electronic Society, 2012: 91-96. [doi: 10.1145/2381966.2381979]
- [6] Bhatia J, Breaux T D. Semantic incompleteness in privacy policy goals. 2018 IEEE 26th International Requirements Engineering Conference (RE), 2018: 159-169. [doi: 10.1109/RE.2018.00025]
- [7] Liu S, Zhao BY, Guo RJ, Meng GZ, Zhang F, Zhang MS. Have You been Properly Notified? Automatic Compliance Analysis of Privacy Policy Text with GDPR Article 13. Proceedings of the Web Conference 2021, 2021: 2154-2164. [doi: 10.1145/3442381.3450022]
- [8] Bhatia J, Breaux T D. A data purpose case study of privacy policies. 2017 IEEE 25th International Requirements Engineering Conference (RE), 2017: 394-399. [doi: 10.1109/RE.2017.56]
- [9] Bhatia J, Breaux T D, Reidenberg J R, Norton T B. A theory of vagueness and privacy risk perception. 2016 IEEE 24th

- International Requirements Engineering Conference (RE), 2016: 26-35. [doi: 10.1109/RE.2016.20]
- [10] Shvartzshnaider Y, Apthorpe N, Feamster N, Nissenbaum H. Going against the (appropriate) flow: A contextual integrity approach to privacy policy analysis. *Proceedings of the AAAI Conference on Human Computation and Crowdsourcing*, 2019, 7(1): 162-170. [doi: 10.1609/hcomp.v7i1.5266]
- [11] Andow B, Mahmud S Y, Wang WY, Whitaker J, Enck W, Reaves B, Singh K, Xie T. Policylint: investigating internal privacy policy contradictions on Google play. *28th {USENIX} Security Symposium ({USENIX} Security 19)*, 2019: 585-602.
- [12] Slavin R, Wang XY, Hosseini M B, Hester J, Krishnan R, Bhatia J, Breaux T D, Niu JW. Toward a framework for detecting privacy policy violations in android application code. *Proceedings of the 38th International Conference on Software Engineering*. 2016: 25-36.
- [13] Yu L, Luo XP, Qian CX, Wang S, Leung H. Enhancing the description-to-behavior fidelity in android apps with privacy policy. *IEEE Transactions on Software Engineering*, 2017, 44(9): 834-854. [doi: 10.1109/TSE.2017.2730198]
- [14] Yu L, Luo XP, Chen JC, Zhou H, Zhang T, Chang H, Leung H. PPChecker: Towards Accessing the Trustworthiness of Android Apps' Privacy Policies. *IEEE Transactions on Software Engineering*, 2021, 47(2): 221-242. [doi: 10.1109/TSE.2018.2886875]
- [15] Zimmeck S, Wang ZQ, Zou LY, Iyengar R, Liu B, Schaub F, Wilson S, Sadeh N, Bellovin S M, Reidenberg J. Automated analysis of privacy requirements for mobile apps. *2016 AAAI Fall Symposium Series*, 2016.
- [16] Wang XY, Qin X, Hosseini M B, Slavin R, Breaux T D, Niu JW. Guileak: Tracing privacy policy claims on user input data for android applications. *Proceedings of the 40th International Conference on Software Engineering*, 2018: 37-47. [doi: 10.1145/3180155.3180196]
- [17] Andow B, Mahmud S Y, Whitaker J, Enck W, Reaves B, Singh K, Egelman S. Actions speak louder than words: Entity-sensitive privacy policy and data flow analysis with polichack. *29th {USENIX} Security Symposium ({USENIX} Security 20)*, 2020: 985-1002.
- [18] Fan M, Yu L, Chen S, Zhou H, Luo XP, Li SY, Liu Y, Liu J, Liu T. An empirical evaluation of GDPR compliance violations in Android mHealth apps. *2020 IEEE 31st International Symposium on Software Reliability Engineering (ISSRE)*, 2020: 253-264. [doi: 10.1109/ISSRE5003.2020.00032]
- [19] Bui D, Yao Y, Shin K G, Choi J M, Shin J. Consistency Analysis of Data-Usage Purposes in Mobile Apps. *Proceedings of the 2021 ACM SIGSAC Conference on Computer and Communications Security*, 2021: 2824-2843. [doi: 10.1145/3460120.3484536]
- [20] Tesfay W B, Hofmann P, Nakamura T, Kiyomoto S, Serna J. I read but don't agree: Privacy policy benchmarking using machine learning and the EU GDPR. *Companion Proceedings of the The Web Conference 2018*, 2018: 163-166. [doi: 10.1145/3184558.3186969]
- [21] Liu S, Zhang F, Zhao BY, Guo RJ, Chen T, Zhang MS. APPCorp: a corpus for Android privacy policy document structure analysis. *Frontiers of Computer Science*, 2023, 17(3). [doi: 10.1007/s11704-022-1627-2]
- [22] Harkous H, Fawaz K, Lebrete R, Schaub F, Shin K G, Aberer K. Polisis: Automated analysis and presentation of privacy policies using deep learning. *27th {USENIX} security symposium ({USENIX} security 18)*, 2018: 531-548.
- [23] Wilson S, Schaub F, Dara A A, Liu F, Cherivirala S, Leon P G, Andersen M S, Zimmeck S, Sathyendra K M, Russell N C, Norton T B, Hovy E, Reidenberg J, Sadeh N. The creation and analysis of a website privacy policy corpus. *Proceedings of the 54th Annual Meeting of the Association for Computational Linguistics*, 2016: 1330-1340. [doi: 10.18653/v1/P16-1126]
- [24] Che WX, Feng YL, Qin LB, Liu T. N-LTP: An Open-source Neural Language Technology Platform for Chinese. In *Proceedings of the 2021 Conference on Empirical Methods in Natural Language Processing: System Demonstrations*, 2021: 41-49.
- [25] Liu F, Fella N L, Liao K. Modeling language vagueness in privacy policies using deep neural networks. *2016 AAAI Fall Symposium Series*, 2016.
- [26] Sun C, Qiu XP, Xu YG, Huang XJ. How to Fine-Tune BERT for Text Classification? *Chinese Computational Linguistics*, 2019:194-206. [doi: 10.1007/978-3-030-32381-3_16]

附中文参考文献:

- [1] 中华人民共和国人民代表大会常务委员会.中华人民共和国个人信息保护法.中国人大网,2021.
- [2] 国家标准化管理委员会.信息安全技术-个人信息安全规范: GB/T 35273-2020.国家标准全文公开系统,2020.
- [3] App专项治理工作组. App 违法违规收集使用个人信息自评估指南.App专项治理工作组,2019.
- [4] 国家互联网信息办公室秘书局,工业和信息化部办公厅,公安部办公厅,市场监管总局办公厅. App 违法违规收集使用个人信息行为认定方法.中国网信网,2019.