

基于 CNN 的 Android 恶意软件识别方法

柏强盛 周 丽

四川大学锦城学院计算机与软件学院 四川 成都 611731

【摘要】随着互联网时代快速发展,移动智能终端设备快速增长,恶意软件也随之大量出现,给用户的财产和隐私安全带来巨大的危害。为了识别 Android 恶意软件,本文将深度学习中的卷积神经网络与 Android 系统权限相结合,基于 LeNet-5 对 APK 中提取出的 AndroidManifest.xml 文件中的系统权限转化特征图进行处理,通过 AndroidManifest.xml 文件中的系统权限来对恶意软件进行识别。提出了一种新的 Android 恶意软件识别方法,实验结果表明,识别准确率达到 87.07%。

【关键词】CNN; LeNet-5; Android; 权限; 恶意软件识别

1 引言

随着移动互联网技术的快速发展与普及^[1],更多的移动智能终端设备例如智能手机成为人们日常工作和生活中的重要工具。

目前 Google 发明的 Android 系统凭借其简单效率高、高扩展性等优点,现已占据了市场的主导地位^[2]。在一份 2020 年 2 月至 2021 年 2 月移动操作系统全球市场份额调查表中可以看出,安卓操作系统占市场份额的 71.9%,所以使用安卓系统手机的用户自然而然地也成为了恶意软件攻击的受害者。Android 开放的生态环境为应用程序的编写提供了便利,但也使恶意软件可利用的漏洞增多^[3]。据统计,Android 平台上恶意程序的数量总体一直呈上升趋势^[4]。根据公开访谈网站上公布的研究结果,79%的恶意软件是专门针对 Android 系统的。而在中国这样的情况更加严重,2017 年的中国移动互联网发展安全报告显示,国家互联网应急中心采集对 205 万个样本进行抓取分析,结果显示其中有 99.9%的样本是针对安卓平台进行攻击的。2012 年以来,移动终端的恶性样本从几十万个上升到几千万个,移动终端(尤其是 Android 终端)的网络安全问题日益突出。每天有 70 万左右手机用户会被恶意软件攻击。为了

应对源源不断的恶意软件,必须要有一些措施来防范恶意软件,让更多的用户避免被恶意软件攻击。

为了识别 Android 恶意软件,本文基于 CNN 的深度学习方法进行 Android 恶意软件识别。

2 卷积神经网络

图像识别是区分不同类别的图像,卷积神经网络(Convolutional Neural Network, CNN)是完成图像识别任务的最佳算法之一^[5]。卷积神经网络主要由输入层、卷积层、池化层和全连接层组成^[6]。卷积层的作用是通过卷积运算提取图像特征,池化层是通过下采样简化网络,全连接层用于图片分类。卷积运算是卷积神经网络的核心,也是卷积神经网络和传统神经网络的不同之处。

3 LeNet-5 模型

卷积神经网络 LeNet-5 算法是由 Y. Le Cun 提出的一个多层的神经网络,对于二维图像的特征提取有着出色的表现^[7]。该模型在手写数字识别中取得了很好的识别效果。LeNet-5 模型共有 8 层(包括输入层和输出层),图 1 展示了 LeNet-5 模型的整体框架结构。

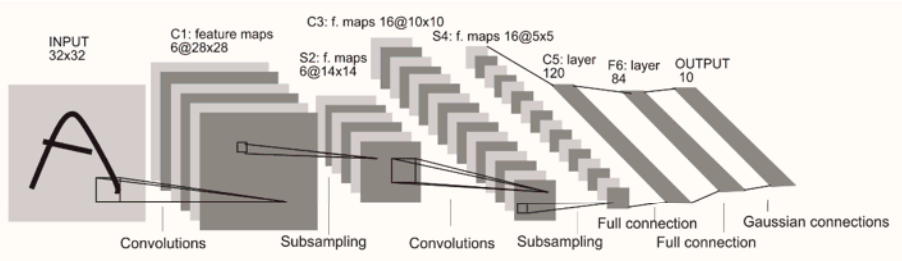


图 1: LeNet-5 网络结构图

LeNet-5 模型将 32*32 的图像像素矩阵作为输入,模型组成结构如图所示。C 表示该网络层为卷积层, S 表示该网络层为下采样层, F 层表示该网络层为全连接层^[8]。

C1 层有 6 个 28*28 大小的神经元矩阵的特征图。

S2 网络层是由 6 个特征图组成的下采样层,特征图是 14*14 大小的神经元矩阵。S2 层的每个神经元与 C1 层中 2*2 大小的感受野内神经元相接。

C3 网络层共有由 16 特征图组成,特征图是 10*10 大小的神经元矩阵。C3 层中的每个神经元与 S2 层的特征图的中 5*5 大小的感受野内神经元相接。

4 数据集以及初始化

通过反编译和解压缩等手段从 APK 中提取出 AndroidManifest.xml 文件。AndroidManifest.xml 文件中共有 200 种不同的权限。我们将 1 个权限看作软件

的 1 个特征，若某个软件拥有某种权限令该特征值为 1，若某个软件没有某种权限则令该特征值为 0。分别把良性软件和不良软件的权限统计制作成 csv 文件。

经过上述处理之后，每个软件在 csv 文件中就是 1



图 2: 数据样本示例图

根据 LeNet-5 处理数据的方式，将 10*20 像素大小的特征图拉伸为 32*32 像素大小。

5 训练模型配置

模型配置如下：

卷积层 1：输入通道数：1，输出通道数：6，卷积核大小：3*3，

卷积步长：1，填充大小：1

激活函数：ReLU

池化层 1：平均池化，池化核大小：2*2，池化步长：2

卷积层 2：输入通道数：6，输出通道数：12，卷积核大小：3*3，

卷积步长：1，填充大小：1

激活函数：ReLU

池化层 2：最大池化，池化核大小：2*2，池化步长：2

全连接层 1：输入大小：12*8*8=768，输出大小：

256

激活函数：ReLU

全连接层 2：输入大小：256，输出大小：128

激活函数：ReLU

全连接层 3：输入大小：128，输出大小：2

损失函数：交叉熵函数

优化器：SGD 优化器

学习率：0.07

动量：0.9

权重衰减参数：0.005

分类器：SoftMax

6 模型的训练和测试

6.1 模型的训练方式

在输入层输入 32*32 像素大小的图片，经过第一层卷积，通道数由 1 增加至 6。由于图像是二值图，所以在池化层 1 选择平均池化的方式，经过池化层 1，图像的大小由 32*32 减小至 16*16。进入第二层卷积，通道数由 6 增加至 12。经过池化层 2，图像大小由 16*16 减小至 8*8。经过三层全连接层，最后输出分类结果。

6.2 4 折交叉验证

交叉验证是机器学习中建立模型和检验模型参数的常用方法，通常用于评估模型的有效性。交叉验证是重复使用数据，从样本中切割数据，将其组合成不同的训练集和测试集，使用训练集训练模型，测试集评估模型。

数据集有 2210 个样本，其中正样本 1035 个负样

本 1175 个。把其中 2094 数据集分为 4 份：

行 200 列的数据，再把这个向量转化为一个软件的特征图，特征图的大小为 10*20 个像素，特征图中 1 个像素点的值就是软件对应 1 个权限的特征值。

部分样本如图 2：

第一份 524 个样本，其中 259 个正样本，265 个负样本，记为 A

第二份 524 个样本，其中 259 个正样本，265 个负样本，记为 B

第三份 524 个样本，其中 259 个正样本，265 个负样本，记为 C

第四份 522 个样本，其中 258 个正样本，264 个负样本，记为 D

把剩余的 116 个负样本做测试集，记为 E。

训练时采用 K 折交叉验证的方式：每次训练时选择 A、B、C、D 中三份数据来训练，另外一份来验证，做 4 次，即 4 折交叉验证

K 折交叉验证的作用是充分利用有限的找到合适的模型参数，防止过度拟合。

6.3 训练结果

训练结果如表 1：

训练集和验证集	训练集损失	训练集结果 (准确率)	测试集结果 (准确率)
A, B, C 为训练集, D 为验证集	0.1640	94.21%	90.23%
A, B, D 为训练集, C 为验证集	0.2217	91.72%	91.03%
A, C, D 为训练集, B 为验证集	0.3083	91.78%	90.27%
B, C, D 为训练集, A 为验证集	0.2029	92.48%	84.73%

表 1: 训练结果表

通过表 1 可以看出，模型的预测的准确率较高，且训练时的损失值较为稳定。说明该模型在训练集上取得了较好的学习效果，在验证集上的表现也较为稳定。

6.4 测试结果

选择验证集结果最好的模型对测试集进行测试，结果如表 2：

测试集	测试集结果 (准确率)
E	87.07%

表 2: 测试结果表

通过测试结果看出，模型的泛化性很好，识别的准确率达到 87.07%。证明了该方法的可行性。

结束语

通过上述内容看出, 本文提出了一种基于卷积神经网络的 Android 恶意软件识别的方法。方法主要的思路是提取 APK 中的 AndroidManifest.xml 文件, 通过文件中标注权限来制作 CSV 文件, CSV 文件中的数据可以转化为软件的权限图, 再通过卷积神经网络对权限图进行训练, 验证集进行验证, 以达到识别恶意软件的目的。测试集上的识别准确率也检验了该方法的可行性。

【参考文献】

[1] 谢聪敏. 基于神经网络的安卓恶意软件检测设计 [J]. 电子设计工程, 2020, 28(09):50-53. 10.14022/j.issn1674-6236.2020.09.011.

[2] 张伟, 徐洋, 张思聪, 徐贵勇. 基于 gru 的 android 恶意软件检测 [J]. 电子技术与软件工程, 2021(04):52-53.

[3] 超凡, 杨智, 杜学绘, 孙彦. 基于深度神经网络的 android 恶意软件检测方法 [J]. 网络与信息安全学报, 2020, 6(05):67-79.

[4] 刘亚姝, 王志海, 李经纬, 赵烜, 文伟平. 基于卡方检验的 android 恶意应用检测方法 [J]. 北

京理工大学学报, 2019, 39(03):290-294. 10.15918/j.tbit1001-0645.2019.03.011.

[5] HuH, Yang Y. A combined GLQP, and DBN-DRF for face recognition in unconstrained environments[C]//2nd International Conference on Control, Automation and Artificial Intelligence (CAAI 2017), 2017.

[6] 孟佳娜, 吕品, 于玉海, 郑志坤. 基于 cnn 的方面级跨领域情感分析研究 [J/OL]. 计算机工程与应用:1-11[2021-04-28 15:04]

[7] 李瑞辰, 姚宇峰, 蒋元华. 基于 lenet-5 的编组站站内机车车号识别系统的研究 [J]. 铁路计算机应用, 2019, 28(09):16-20.

[8] 王博, 朱兆旻, 唐天兵. 一种改进的 lenet-5 卷积神经网络算法 [J]. 大众科技, 2020, 22(10):1-3.