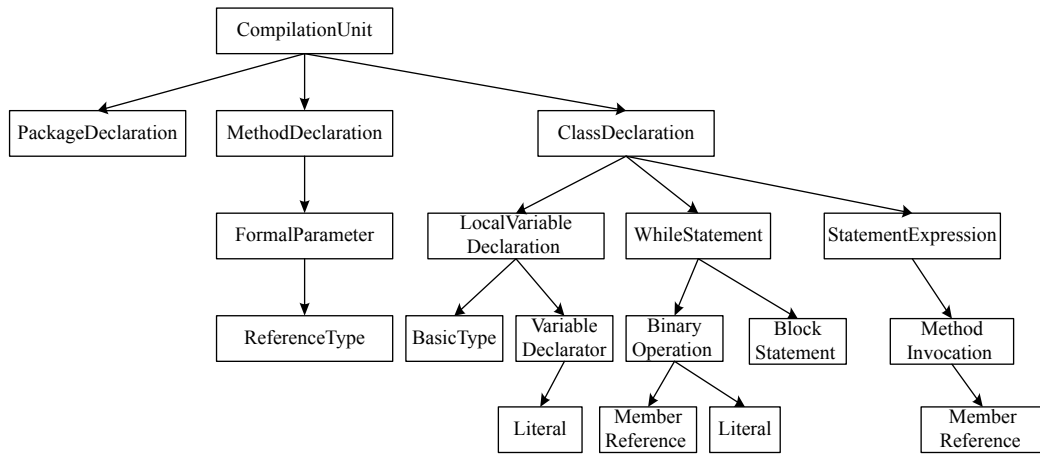


(a) 代码 a 的抽象算法树



(b) 代码 b 的抽象算法树

图 2 示例代码的抽象语法树

表 1 使用到的所有结点类型

类别	结点类型
方法调用结点	MethodInvocation, SuperMethodInvocation
声明结点	PackageDeclaration, InterfaceDeclaration, ClassDeclaration, MethodDeclaration, ConstructorDeclaration, VariableDeclaration
控制流结点	IfStatement, WhileStatement, DoStatement, ForStatement, SwitchStatement, AssertStatement, BreakStatement, ContinueStatement, ReturnStatement, ThrowStatement, TryStatement, SynchronizedStatement, BlockStatement, TryResource, CatchClause, EnhancedForControl
其他结点	FormalParameter, BasicType, CatchClauseParameter, MemberReference, SuperMemberReference, ReferenceType

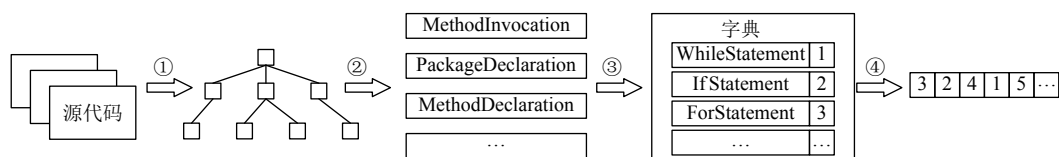


图 3 从源代码到得到整数向量的流程图 (①从源代码中解析出抽象语法树; ②从抽象语法树中提取表征向量; ③根据提取的特征构建字典; ④将表征向量映射为整数向量)

3 卷积神经网络设计与结构

卷积神经网络是一种含卷积计算且具有深度结构的前馈式神经网络,具有表征学习、深度挖掘数据的能力.卷积神经网络已经被证实在图像处理、语音识别、自然语言处理等领域有着不错的性能表现.而且,卷积神经网络的隐含层内的卷积核参数共享、层与层之间的稀疏连接使得卷积神经网络能够以较小的计算量对格点化特征.

3.1 数据预处理

在软件缺陷预测领域,数据分类不均衡问题是普遍存在的,如何处理这种不均衡问题也是具有挑战性的.在获得的数据集中,往往有缺陷的样本数是要少于没有缺陷的样本数,如果直接用这样的样本对模型进行训练,训练出来的模型往往会对样本数量较少的类别的识别能力较弱,在软件缺陷预测时,模型便会偏向没有缺陷的结果.因此,我们需要对数据进行预处理以弥补数据集分类不均衡带来的偏差.

一般而言,针对数据集中存在的样本分类不均衡问题,通常采用过采样或欠采样的方法来使得数据集中各类的样本数量达到均衡.欠采样是从多数样本中随机选择和少数样本一样数量的样本,以此构建分类平衡的数据样本,但这样会大大减少数据的训练样本,造成信息的损失,会造成模型的欠拟合,因此我们选择过采样的方法来处理分类不均衡问题.由于AST数值向量并非是特征向量,所以类似像SMOTE^[5]等基于样本间的欧式距离来生成新样本的分类不均衡处理方法并不一定适用.因此,在这里,我们采取简单的随机过采样的方法来对数据进行预处理,在少数样本即有缺陷的样本中随机选择进行复制,使得有缺陷的样本数量和没有缺陷的样本数量保持一致,以此保证数据类别间的均衡.具体的算法如算法1所示.

算法1. 随机过采样算法

输入: 分类不均衡的数据集 D

输出: 分类均衡的数据集 D'

- (1) 初始化已复制样本集合 C 为空集.
- (2) 将数据集 D 复制, 组成数据集 D' .
- (3) 在数据集 D 中筛选出有缺陷的样本集 d .
- (4) 在数据集 d 中随机选择样本 a , 如果 a 不在已复制样本集合 C 中, 将样本 a 加入数据集 D' 和已复制样本集合 C ; 如果已复制样本集合 C 的大小和 d 一样, 则清空 C .
- (5) 重复步骤 (4) 直至数据集 D' 中无缺陷样本和有缺陷样本数量保持一致.

3.2 网络结构

首先介绍所要构建的网络中用到的基础卷积块 Inception 块^[6], 这个卷积模块的结构如图4所示.

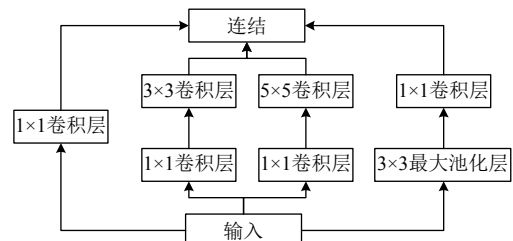


图4 Inception 块

如图4所示, Inception 块含有4条并行的网络线路.前三条线路所使用的卷积核大小分别是 1×1 、 3×3 、 5×5 , 用来抽取不同空间尺寸下的信息, 第2, 3层会先对输入做 1×1 卷积操作以减少输入通道数, 以降低模型复杂度. 第4条线路会先使用 3×3 最大池化层, 然后接一个 1×1 卷积层来改变通道数. 并且, 根据4条线路的具体情况, 设置合适的填充来使得输入和输出的高和宽一致. 最后将每条线路的输出在通道维上连结, 以此完成卷积的功能, 将结果输入到下一层. 相较于普通的卷积层, Inception 块因为使用了3个卷积层和1个池化层, 能够更深层次的挖掘数据中的特征, 以此帮助模型进行更好的学习分类.

为了能够更深层次地挖掘出潜藏在向量中的语法语义特征, 本文基于 GoogLeNet^[7] 设计了一个卷积神经网络, GoogLeNet 最初设计出来是用来进行图像处理的, 在 ImageNet 图像识别挑战赛中大放异彩, GoogLeNet 串联了多个 Inception 块来对图像进行深度挖掘, 以此进行更好的分类. 本文基于 GoogLeNet 设计一个卷积神经网络, 具体的网络结构如图5所示, 除了最后的输出层使用 Sigmoid 函数作为激活函数外, 其他层的激活函数均使用 ReLU 函数, 并且根据实际情况调整网络中各个层的参数, 网络整体分为主要分为以下3个部分:

(1) 输入层: 输入层主要是一个嵌入层^[8], 嵌入层的主要作用是将输入的整数向量中的整数元素转换成整数向量, 使得向量可以进行卷积操作. 嵌入层有两个重要参数: 嵌入层字典的大小 (num_embeddings) 和每个产出向量的大小 (embedding_dim). 这里, 本文将 num_embeddings 设置为 2.3 节构建的字典中所含有的

特征的数量,将 embedding_dim 设置为 2.3 节中通过映射得到的整数向量的长度.将长度为 n 的整数向量输入到嵌入层,嵌入层将给出一个 $n \times n$ 的矩阵向量.并且,为了提高内存利用率和训练模型的速度,本文选择分批进行训练,设置每次训练样本个数(批尺寸, Batch Size)为 16,即一次输入 16 个样本进行训练.

(2) 卷积部分:卷积部分是网络的主体部分,共由 5 个模块组成.模块与模块之间使用步幅为 2 的 3×3 最大池化层来减小输出高度.第 1 个模块包含 3 层的 3×3 卷积层;第 2 个模块使用 2 个卷积层,首先接一个 64 通道的 1×1 卷积层,然后接了一个将通道数扩大 3 倍的 3×3 卷积层;第 3 个模块串联了 2 个完整的 Inception 块;第 4 模块串联了 5 个 Inception 块;第 5 模块串联了 2 个 Inception 块.通过多层的不同空间尺寸的卷积操作,来深度挖掘数据中的特征,从而进行性能更好稳定性更高的学习分类.

(3) 输出层:输出层主要是根据之前卷积层输出的结果来输出分类结果.首先使用一个全局平均池化层来将每个通道的高和宽都变成 1,然后接上一个全连接层,输出通道数为标签类别数,最后,连结一个 Sigmoid 函数构建逻辑回归分类器来计算程序代码的缺陷概率,从而得到分类结果.

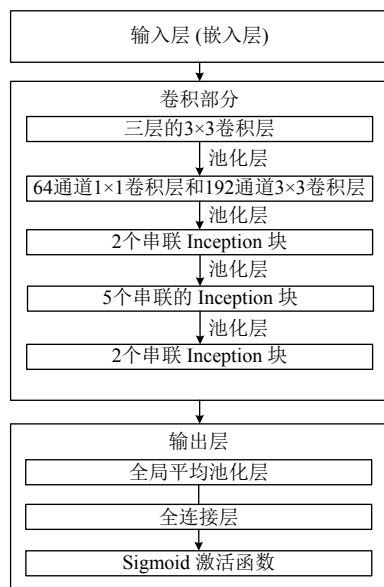


图 5 网络结构图

3.3 模型优化

之前,我们在数据预处理时采用随机过采样的方

法来解决数据分类不均衡问题,提升了模型的泛化能力,但是这样也有一定的过拟合的风险,因此我们选择使用丢弃法 (Dropout)^[9],通过随机丢弃一部分神经元来避免过拟合.在训练过程中,随机丢弃一部分隐藏层神经元,即所有神经元都有可能被清零,这样就减少了神经元之间的依赖性,输出层的计算也无法过度依赖任何一个隐含层神经元,从而在训练模型时起到正则化的作用,用来应对过拟合.在测试模型时,为了拿到更加准确的结果,我们不使用丢弃法.

另外,在训练模型的过程中,为了得到最优的模型参数,我们需要根据损失函数的梯度不断地对参数进行迭代,这里我们选择使用 Adam^[10] 优化器来更新参数. Adam 算法结合了 AdaGrad 和 RMSProp 两种优化算法的优点^[10],能够自动调整学习率,并且参数的更新不受梯度的伸缩变换影响. Adam 算法能够从梯度均值及梯度平方两个角度进行自适应地调节,综合考虑梯度的一阶矩估计和二阶矩估计计算出更新步长,而不是直接由当前梯度决定.

4 实验与结果分析

4.1 数据集

为了评估训练出来的模型的性能,本文从 Promise^[11] 上下载了 5 个工程,总共 11 个项目,组成 6 组软件缺陷预测任务,用于模型的测试.在同一软件工程中,将旧版本的工程项目作为训练集,将新版本的工程项目作为测试集,根据测试结果来评估模型的预测能力.例如,对于 Camel 工程,我们将 Camel 1.4 版本的工程代码用来训练模型,然后用 Camel 1.6 版本的代码用来测试模型.表 2 列出了测试时所使用的软件项目的基本信息.

表 2 测试使用的工程信息

工程	版本	平均代码数	平均缺陷率(%)
Camel	1.4, 1.6	918	18.1
Lucene	2.0, 2.2	209	55.7
Synapse	1.0, 1.1, 1.2	212	25.5
Poi	2.5, 3.0	413	64.0
Xalan	2.6, 2.7	844	47.3

另外,在数据集中,每个项目不仅含有工程的源代码,还统计了源代码的静态代码度量和缺陷注释,度量方法主要是针对面向对象编程的静态代码度量,具体的指标内容如表 3 所示.这些指标可以用于其他的软件缺陷预测方法,来和本文模型进行比较.

表3 数据集中所使用的20个静态代码度量

简写	全称
LOC	Lines of code
DIT	Depth of inheritance tree
NOC	Number of children
RFC	Response for a class
CBO	Coupling between object classes
LCOM	Lack of cohesion in methods
LCOM3	Lack of cohesion in methods
NPM	Number of public methods
DAM	Data access metric
MOA	Measure of aggregation
MFA	Measure of function abstraction
IC	Inheritance coupling
CAM	Cohesion among methods of class
CBM	Coupling between methods
AMC	Average method complexity
Ca	Afferent couplings
Ce	Efferent couplings
Avg(CC)	Average McCabe
Max(CC)	Maximum McCabe
WMC	Weighted methods per class

4.2 评估指标

本文采用 AUC^[12] 和 F1-measure^[13] 这两个指标来评估模型性能, AUC 主要用来评估模型的识别缺陷的能力, 而 F1-measure 主要用来评估模型的稳定性。

在二分类的学习预测过程中, 根据分类结果可以将其分为4类: (1) 若一个实例为正类且被预测为正类, 即为真正类 (True Positive, *TP*); (2) 若一个实例为正类但被预测负类, 即为假负类 (False Negative, *FN*); (3) 若一个实例为负类但被预测为正类, 即为假正类 (False Positive, *FP*); (4) 若一个实例为负类且被预测为负类, 即为真负类 (True Negative, *TN*)。基于这4个数据, 可以得到击中概率 (*TPR*) 和虚报概率 (*FPR*), 其计算公式如式 (1) 和式 (2) 所示:

$$TPR = \frac{TP}{TP + FN} \quad (1)$$

$$FPR = \frac{FP}{TN + FP} \quad (2)$$

然后以 *FPR* 为横轴, *TPR* 为纵轴, 就可以绘制出 ROC 曲线, 而 ROC 曲线下的面积就是 AUC。根据 AUC 的定义, 识别能力更好的模型对应着更高的 *TPR* 和更低的 *FPR*, 所以有 AUC 值越大的预测方法越好。

F1-measure 是精确率 (*P*) 和召回率 (*R*) 的调和平均, 其中精确率和召回率的计算公式如式 (3) 和式 (4) 所示:

$$P = \frac{TP}{TP + FP} \quad (3)$$

$$R = \frac{TP}{TP + FN} \quad (4)$$

通常情况下, 我们希望精确率越高越好, 召回率也越高越好, 但事实上这两个指标在某些情况下是矛盾的, 而 F1-measure 则综合考虑了这两个指标。F1-measure 的计算公式如式 (5) 所示。

$$F_1 = \frac{2 \cdot P \cdot R}{P + R} \quad (5)$$

另外, 用于评估软件缺陷预测模型的指标还有很多, 例如 MCC^[14] 和 G-mean^[15], MCC 考虑所有的有缺陷数据和无缺陷数据来综合评估预测结果和真实结果之间的关系, G-mean 是在数据不平衡时十分有参考价值的一个指标, 但因为 AUC 和 F1-measure 综合评估了模型的准确性和稳定性, 具有广泛的代表意义。

4.3 实验方法

为了能够正确估计模型对于软件缺陷预测的性能, 将本文提出的模型与以下3种方法进行比较:

(1) 静态代码度量+逻辑回归 (LR): 以数据集中提供的20个静态代码度量作为代码特征, 并用逻辑回归的方法进行分类

(2) 深度置信网络+逻辑回归 (DBN)^[2]: 使用深度置信网络从源代码中提取特征, 然后使用逻辑回归的方法进行分类

(3) 卷积神经网络+逻辑回归 (CNN)^[4]: 利用单层卷积神经网络对源代码进行特征提取, 然后使用逻辑回归分类器得到分类结果

对于传统软件缺陷预测算法, 因为使用的是20个静态代码度量所构成的特征向量, 所以在数据预处理时, 可以使用 SMOTE 方法进行过采样来处理数据集分类不平衡问题; 而对于 DBN、CNN 和本文模型, 只能简单地采用随机过采样的方法来对数据进行预处理。

本文使用 Python 环境以及深度学习框架 PyTorch 来实现本文提出的模型, 所有的实验均在一台带有 NVIDIA GTX 1080 的 Linux 服务器上运行。此外, 因为随机过采样和丢弃法都具有一定的随机性, 因此实验中每个方法都执行 10 次, 取平均值来进行模型性能的比较。

4.4 实验结果及分析

本文采用 AUC 和 F1-measure 来比较 4 种方法在

6组预测任务上的性能.表4和表5分别记录了这4种方法关于AUC和F1-measure的实验结果,每次测试任务的表现最好的已在表格中加粗.

表4 4种方法关于6项测试任务的AUC

训练集	测试集	LR	DBN	CNN	本文模型
Camel 1.4	Camel 1.6	0.599	0.641	0.687	0.709
Lucene2.0	Lucene2.2	0.628	0.626	0.635	0.641
Synapse1.0	Synapse1.1	0.600	0.639	0.594	0.646
Synapse1.1	Synapse1.2	0.637	0.697	0.622	0.674
Poi2.5	Poi3.0	0.665	0.651	0.710	0.718
Xalan2.6	Xalan2.7	0.651	0.683	0.675	0.674
平均值		0.630	0.656	0.654	0.677

表5 4种方法关于6项测试任务的F1-measure

训练集	测试集	LR	DBN	CNN	本文模型
Camel 1.4	Camel 1.6	0.355	0.342	0.489	0.513
Lucene2.0	Lucene2.2	0.611	0.666	0.709	0.724
Synapse1.0	Synapse1.1	0.444	0.457	0.456	0.461
Synapse1.1	Synapse1.2	0.471	0.493	0.487	0.473
Poi2.5	Poi3.0	0.707	0.716	0.756	0.776
Xalan2.6	Xalan2.7	0.633	0.642	0.654	0.681
平均值		0.537	0.553	0.592	0.605

表3和表4分别列出了4种方法关于每个测试任务的AUC值和F1-measure. AUC评估了模型分类的准确性,而F1-measure评估了模型的稳定性.从表3和表4中我们可以看到,总体而言,本文提出的模型在软件缺陷预测性能方面和模型稳定性明显优于LR、DBN和CNN的.而本文模型的AUC和F1-measure的均值也都高于其他方法,这也证实了本文提出模型的合理性和可行性.此外,从两张表中我们可以看出,相较于传统的软件缺陷预测方法,应用深度学习方法在软件缺陷预测性能和模型稳定性上都得到一定的提高.这也证实了,在软件缺陷预测性能方面,深度学习方法优于传统的机器学习方法.

综上所述,针对传统软件缺陷预测方法中对源代码语义特征挖掘不足的问题,本文测试实验结果表明,在软件缺陷预测领域,相比于传统的预测方法,应用深度学习方法得到了一定的提高.而本文也根据前人的工作,提出了用多层卷积神经网络对基于抽象语法树得到的表征向量进行分类学习,有效提高了缺陷预测的准确性.

5 结束语

本文针对传统软件缺陷预测方法应用静态代码度

量而忽视代码语义的缺点,从代码的抽象语法树中提取出向量,再利用卷积神经网络深度挖掘数据的能力挖掘代码中的语法语义特征,从而对软件缺陷进行学习分类.并且,通过与LR、DBN、MLP方法的实验比较,由AUC和F1_measure两个指标我们可以看出本文提出的模型在软件缺陷预测性能上得到了一定的提高.然而,关于数据集分类不均衡、模型优化等问题,本文的处理方法相对粗糙,这也是未来需要继续研究的方向.

参考文献

- 刘童. 基于机器学习算法的软件缺陷预测技术研究 [硕士学位论文]. 武汉: 华中师范大学, 2018.
- Wang S, Liu TY, Tan L. Automatically learning semantic features for defect prediction. 2016 IEEE/ACM 38th International Conference on Software Engineering. Austin, TX, USA. 2016. 297–308.
- Iqbal A, Aftab S. A classification framework for software defect prediction using multi-filter feature selection technique and MLP. International Journal of Modern Education and Computer Science, 2020, 12(1): 18–25. [doi: 10.5815/ijmecs.2020.01.03]
- Li J, He PJ, Zhu JM, et al. Software defect prediction via convolutional neural network. 2017 IEEE International Conference on Software Quality, Reliability and Security. Prague, Czech Republic. 2017. 318–328.
- Chawla NV, Bowyer KW, Hall LO, et al. SMOTE: Synthetic minority over-sampling technique. Journal of Artificial Intelligence Research, 2002, 16(1): 321–357.
- Szegedy C, Ioffe S, Vanhoucke V, et al. Inception-v4, Inception-ResNet and the impact of residual connections on learning. Proceedings of the 31st AAAI Conference on Artificial Intelligence. San Francisco, CA, USA. 2017. 4278–4284.
- Tang PJ, Wang HL, Kwong S. G-MS2F: GoogLeNet based multi-stage feature fusion of deep CNN for scene recognition. Neurocomputing, 2017, 225: 188–197. [doi: 10.1016/j.neucom.2016.11.023]
- Mikolov T, Sutskever I, Chen K, et al. Distributed representations of words and phrases and their compositionality. Proceedings of the 26th International Conference on Neural Information Processing Systems. Red Hook, NY, USA. 2013. 3111–3119.
- Ba LJ, Frey B. Adaptive dropout for training deep neural

- networks. Proceedings of the 26th International Conference on Neural Information Processing Systems. Red Hook, NY, USA. 2013. 3084–3092.
- 10 Kingma DP, Ba J. Adam: A method for stochastic optimization. arXiv preprint arXiv: 1412.6980, 2014.
- 11 Cukic B. Guest editor's introduction: The promise of public software engineering data repositories. IEEE Software, 2005, 22(6): 20–22. [doi: [10.1109/MS.2005.153](https://doi.org/10.1109/MS.2005.153)]
- 12 Fawcett T. An introduction to ROC analysis. Pattern Recognition Letters, 2006, 27(8): 861–874. [doi: [10.1016/j.patrec.2005.10.010](https://doi.org/10.1016/j.patrec.2005.10.010)]
- 13 Powers DMW. Evaluation: From precision, recall and F-measure to ROC, informedness, markedness & correlation. Journal of Machine Learning Technologies, 2011, 2: 37–63.
- 14 Olsson PQ, Cotton WR. Balanced and unbalanced circulations in a primitive equation simulation of a midlatitude MCC. Part II: Analysis of balance. Journal of Atmospheric Sciences, 1997, 54(4): 479–497.
- 15 Guo HP, Liu HB, Wu CA, *et al.* Logistic discrimination based on G-mean and F-measure for imbalanced problem. Journal of Intelligent & Fuzzy Systems, 2016, 31(3): 1155–1166.

www.c-s-a.org.cn

www.c-s-a.org.cn