

# 改进的神经语言模型及在代码提示中的应用

张献 贲可荣

(海军工程大学电子工程学院 武汉 430033)

**摘要** 语言模型旨在刻画文本段的发生概率,作为自然语言处理领域中的一类重要模型,其近年来被广泛应用于不同软件分析任务,例如代码提示。为增强模型对代码特征的学习能力,提出了一种改进的循环神经网络语言模型 CodeNLM。通过分析词向量形式表示的源代码序列,该模型能够捕获代码规律,实现对序列联合概率分布的估计。考虑到现有模型仅学习代码数据,信息利用不充分,提出了附加信息引导策略,通过非代码信息的辅助提高了代码规律的刻画能力。针对语言建模任务特点,提出了节点逐层递增策略,通过优化网络结构改善了信息传递的有效性。在验证实验中,针对 9 个 Java 项目共 203 万行代码,CodeNLM 得到的困惑度指标明显优于对比的  $n$ -gram 类模型和神经语言模型,在代码提示应用中得到的平均准确度(MRR 指标)较对比方法提高 3.4%~24.4%。实验结果表明,CodeNLM 能有效地实现程序语言建模和代码提示任务,并具有较强的长距离信息学习能力。

**关键词** 软件分析,代码提示,自然语言处理,语言模型,循环神经网络

中图法分类号 TP311.5

文献标识码 A

DOI

## A Modified Neural Language Model and Its Application in Code Suggestion

ZHANG Xian BEN Ke-rong

(School of Electronic Engineering, Naval University of Engineering, Wuhan 430033, China)

**Abstract** Language models are designed to characterize the occurrence probabilities of text segments. As a class of important model in the field of natural language processing, it has been widely used in different software analysis tasks in recent years. To enhance the learning ability for code features, this paper proposed a modified recurrent neural network language model called CodeNLM. By analyzing the source code sequences represented in embedding form, the model can capture rules in codes and realize the estimation of the joint probability distribution of the sequences. Considering that the existing models only learn the code data and the information is not fully utilized, this paper proposed an additional information guidance strategy, which can improve the ability of characterizing the code rules through the assistance of non-code information. Aiming at the characteristics of language modeling task, a layer-by-layer incremental nodes setting strategy is proposed, which can optimize the network structure and improve the effectiveness of information transmission. In the verification experiments, for 9 Java projects with 2.03M lines of code, the perplexity index of CodeNLM is obviously better than the contrast  $n$ -gram class models and neural language models. In the code suggestion task, the average accuracy (MRR index) of our model is 3.4%~24.4% higher than the contrast methods. The experimental results show that with a strong long-distance information learning capability, CodeNLM can effectively model programming language and perform code suggestion well.

**Keywords** Software analysis, Code suggestion, Natural language processing, Language model, Recurrent neural network

## 1 引言

近年来,伴随软件的更新演化、开源思想的深入,软件仓库,如 GitHub、Stack Overflow 等,累积了海量的代码和其他元数据。这些日益扩增的软件制品促使着新的、数据驱动的方法分析软件。目前已有一批重要成果引入自然语言处理(natural language processing, NLP)、机器学习等方法解决不同软件分析任务<sup>[1-12]</sup>。可以看出,这些技术的应用正成为近年来软件分析的热点内容。

语言模型(language model)是 NLP 领域中的一类重要模型,其通过语料库学习对文本段的发生概率进行估计,进而刻画出语言规律。目前,该模型已在自动分词、句法分析和

机器翻译等问题中得到了广泛应用<sup>[3]</sup>。近些年,不少研究者开始尝试采用语言模型来分析程序语言/软件代码,应用包括代码提示(code suggestion)<sup>[4]-[6]</sup>、代码补全(code completion)<sup>[7]</sup>、API 推荐(application programming interface recommendation)<sup>[8]</sup>、编程规范(coding convention)提示<sup>[9]</sup>、缺陷检测<sup>[10]</sup>、程序合成(Program Synthesis)<sup>[11]</sup>、伪代码自动生成<sup>[12]</sup>等。这些成果主要采用了经典的  $n$ -gram 语言模型或其改进形式。虽然这类模型构建简单、直接,但由其自身特点引起的维数灾难问题十分显著,致使模型无法有效地学习长距离上下文信息。

近年来,迅速兴起的深度学习方法为克服上述问题带来了希望。由神经网络构建的语言模型在利用词间语义相似性和学习长距离序列特征方面具有天然优势,可以有效克服传

到稿日期: 返修日期:

本文受国家安全重大基础研究计划(613315)项目资助。

张献(1990-),男,博士生,CCF 学生会员,主要研究方向为软件缺陷挖掘、机器学习, E-mail: tomtomzx@foxmail.com。

贲可荣(1963-),男,博士,教授,CCF 高级会员,主要研究方向为软件质量保证、人工智能, E-mail: benkerong08@163.com。

统语言模型的固有缺陷, 目前已在 NLP 领域取得了显著进展<sup>[13]</sup>。为此, 近几年人们也开始重点关注此类方法在程序语言建模方面的应用与拓展<sup>[2]</sup>。例如, White 等<sup>[14]</sup>面向海量代码语料库建立了循环神经网络 (recurrent neural network, RNN) 结构的程序语言模型; Dam 等<sup>[15][16]</sup>提出了一种基于长短期记忆 (long short-term memory, LSTM) 单元的 RNN 模型 DeepSoft; Allamanis 等<sup>[17]</sup>为推荐合适的方法/类名称构建了一种基于神经网络的程序语言模型; 尹春林等<sup>[18]</sup>基于 RNN 语言模型提出了一种源代码主题建模方法并应用于特征定位; 张献等<sup>[19]</sup>利用神经语言模型度量得到的代码熵值进行缺陷预测。这些工作都采用了深度神经网络 (deep neural network, DNN) 架构, 取得的效果相比传统模型均有较为明显的提升。但这些方法仍存在一些不足, 例如在解决实际任务时, 它们都仅以代码数据为分析对象, 而与任务密切相关的非代码信息未得到有效利用, 致使模型对软件的分析还不够充分。

鉴于此, 本文利用深度学习方法对程序语言进行建模, 依据任务特点提出了一种改进的基于 LSTM 单元的 RNN 语言模型 CodeNLM, 并以代码提示任务为例对模型应用进行了验证。该模型以词化后的源代码序列为输入, 利用 DNN 学习代码的上下文信息、捕获程序特征。为增强模型刻画语言规律的能力, 提出了附加信息引导策略, 即引入任务相关的非代码信息以辅助网络学习与决策。在代码提示应用中, 抽取的附加信息为代码与项目的关系信息。为提高泛化性能, 还提出了节点逐层递增策略, 通过限制神经网络规模优化了结构设置。最后, 本文以实际软件项目为分析对象, 通过对比实验验证了 CodeNLM 在程序语言建模、代码提示及长距离信息学习等方面的有效性。

本文组织结构如下: 第 2 节介绍预备知识及相关成果; 第 3 节详细阐述本文方法; 第 4 节在实际数据集上进行实验与结果分析, 以评价本文方法; 最后一节对全文进行总结并展望下一步工作。

## 2 预备知识与相关工作

### 2.1 语言模型

语言模型旨在刻画语言规律, 其通过语料库学习对序列的联合概率分布进行估计, 是处理语言任务的一类基础模型<sup>[3]</sup>。形式化地, 给定一个长度为  $l$  的语言序列  $S = w_1 w_2 \cdots w_l$ , 其中  $w_i \in \Omega$ ,  $S \in \Omega^*$ ,  $\Omega$  为最小语言单元集合,  $\Omega^*$  为所有语言序列集合, 本文中最小语言单元为代码的词 (token)。那么序列  $S$  的发生概率可用其联合概率表示:

$$P(S) = P(w_1)P(w_2 | w_1) \cdots P(w_l | w_1, w_2, \cdots, w_{l-1}) \\ = \prod_{i=1}^l P(w_i | w_1, \cdots, w_{i-1}) \quad (1)$$

给定一个语料库  $C \subseteq \Omega^*$ , 语言模型以极大似然估计方法估算参数化的序列概率  $P(S, \theta)$ , 其中  $S \in C$ ,  $\theta$  为待估计的

模型参数。估计  $P(S, \theta)$  的方法有多种, 目前最为常用的是  $n$ -gram 模型和神经语言模型 (neural language model)。

#### (1) $n$ -gram 语言模型

$n$ -gram 模型运用统计方法估计序列的联合概率分布, 是语言模型中的一个经典工作。为简化计算, 该模型给定了一个马尔科夫假设: 当前词的出现概率仅仅与前面  $n-1$  个词相关, 因此式(1)被近似表示为:

$$P(S) = \prod_{i=1}^l P(w_i | w_1, \cdots, w_{i-1}) \approx \prod_{i=1}^l P(w_i | w_{i-(n-1)}, \cdots, w_{i-1}) \quad (2)$$

对于式(2), 常采用频数计数法估算其中的条件概率:

$$P(w_i | w_{i-(n-1)}, \cdots, w_{i-1}) = \frac{\text{count}(w_{i-(n-1)}, \cdots, w_{i-1}, w_i)}{\text{count}(w_{i-(n-1)}, \cdots, w_{i-1})} \quad (3)$$

虽然  $n$ -gram 模型的概率计算过程较为简单, 但其难以学习长距离信息。这主要是由于模型将语言单元视为孤立的符号, 无法有效地利用词间的语义关系, 造成维数灾难问题。例如, 考虑一个大小为 2000 的词汇表, 令序列最大长度为 10, 那么不同的语言序列就约有  $10^{33}$  种可能, 因此 10-gram 模型在式(3)的频数统计中, 将造成数据稀疏和内存爆炸。为此, 实际应用中常选择 5-gram 规模以下模型。

#### (2) 神经语言模型

神经语言模型借助神经网络估计序列的联合概率分布, 通过利用词间语义相似性克服传统模型带来的维数灾难问题。该模型由 Bengio 等人<sup>[20]</sup>于 2001 年正式提出, 目前已在机器翻译、信息检索、对话系统等诸多 NLP 任务中起到关键作用。神经语言模型一方面利用词嵌入 (word embeddings) 技术<sup>[21]</sup>将离散的语言符号映射为连续的一组矢量, 语义相似的词将拥有相似的矢量, 增强了模型对词间语义关系的学习能力; 另一方面借助高度非线性的神经网络拟合序列的发生概率, 可以更灵活地刻画语言规律。此外, RNN 等深层网络的兴起, 还使得模型在长程依赖 (long-term dependency) 学习方面具有显著优势, 其中代表性的工作包括 Mikolov 等人首次提出的 RNN 语言模型<sup>[22]</sup>及其在词向量 (word vector) 生成方面的应用<sup>[23]</sup>、Zaremba 等人提出的正则化 LSTM 网络语言模型<sup>[24]</sup>等 (参见综述<sup>[25]</sup>)。

近年来, 已有一批学者尝试采用语言模型处理不同软件分析任务<sup>[2]</sup>。Hindle 等<sup>[4]</sup>运用  $n$ -gram 模型首次对代码的自然性 (naturalness) 进行了分析, 发现程序具有显著的重复性、可预测性。基于此工作, Ray 等<sup>[10]</sup>分析了缺陷代码的自然性, 并成功将  $n$ -gram 模型应用于缺陷检测问题。为增强局域预测能力, Tu 等<sup>[5]</sup>提出了一种附加“cache”组件的  $n$ -gram 模型。Nguyen 等<sup>[6]</sup>通过综合多种代码信息构建了一种统计语义语言模型 SLAMC, 并应用于代码提示任务。为增强传统  $n$ -gram 模型对长序列规律的刻画, Yang 等<sup>[7]</sup>提出了一种基于中间表示的改进模型, 较好地实现了代码补全问题, 即能一次性推荐较长的代码段。Bielik 等<sup>[11]</sup>定义了一种领域特定语言, 以此构建了字符级的语言模型并完成程序合成任务。为规范代

码风格, Allamanis 等<sup>[9]</sup>利用  $n$ -gram 模型实现了编程规范提示。Oda 等<sup>[12]</sup>基于  $n$ -gram 模型和统计机器翻译框架学习伪代码自动生成任务。

为增强代码特征学习能力并克服  $n$ -gram 类模型固有缺陷, 近几年人们也开始尝试运用神经语言模型。White 等<sup>[14]</sup>利用 RNN 语言模型对大规模语料库进行了挖掘, 通过多种实验验证了深度学习模型在学习软件代码规律上优于  $n$ -gram 类模型。为提高模型对长序列依赖关系的学习能力, Dam 等<sup>[15][16]</sup>使用了 LSTM 单元构建 RNN 并提出了一种程序语言模型 DeepSoft, 在 Java 语料库上该模型性能优于原始 RNN 模型。Allamanis 等<sup>[17]</sup>提出了两种神经语言模型, 然后结合词嵌入表示方法, 较好地完成了方法命名问题 (method naming problem), 即为方法或类推荐功能性描述的名称, 用以提高代码可读性和维护性。为增强对代码语境的学习, 尹春林等<sup>[18]</sup>基于 RNN 语言模型提出了一种源代码主题建模方法, 其利用神经网络学习代码及描述文本的语义向量, 通过计算向量间的相似性进行特征定位, 并取得了较优的查准率。张献等<sup>[19]</sup>面向软件缺陷预测提出了一个框架 DefectLearner, 其首先利用 RNN 语言模型对代码模块进行熵值度量, 然后基于度量结果改善预测性能。

## 2.2 代码提示

代码提示作为一项重要技术常集成于软件开发环境, 其通过即时推荐候选代码辅助编程任务。准确的提示结果不仅可以提高开发效率, 还能保证软件质量, 因此探索更有效的代码提示方法一直是软件分析领域的一个重要研究内容。图 1 展示了 2 个代码提示示例 (示例由本文方法在 Log4J 项目中得到, 见 4.1 节)。

```
public class JDBCTest
{.....
    try {
        Class.forName("sun.jdbc.odbc.JdbcOdbcDriver");
    }
    catch (Exception e)
    {
        e.printStackTrace();
        System.out.println(e. ? )
    }
    .....

    try {
        Category log = Category.getInstance("main");
        log.debug("Debug 1");
        Thread.sleep(500);
        log.info("info 1");
        Thread.sleep(500);
        log.warn("warn 1");
        Thread.sleep( ? )
    }
}
```

提示候选项:  
 1. toString ✓  
 2. getMessage  
 3. getClass  
 4. getName  
 5. printStackTrace

提示候选项:  
 1. 500 ✓  
 2. 200  
 3. 100  
 4. 10  
 5. 50

图1 代码提示示例 (由 CodeNLM 得到)

Fig1 An example of code suggestion (acquired by CodeNLM)

代码提示即是为当前语句推荐后续代码, 其可视为文本

序列预测问题。回顾到式(1), 语言模型在计算序列联合概率的同时, 需要估计条件概率  $P(w_i | w_1, \dots, w_{i-1})$ , 这些条件概率恰好可用于文本序列预测。因此自然地, 可利用语言模型进行代码提示工作。Hindle 等<sup>[4]</sup>最早进行了尝试, 其应用  $n$ -gram 模型于 Java 代码提示, 实验结果表明该方法可增强 Eclipse 开发环境的提示效果。基于此工作, Tu 等<sup>[5]</sup>提出了一种 cache 语言模型, 以增强局域预测能力, 并取得了优于基本  $n$ -gram 模型的代码提示性能。Nguyen 等<sup>[6]</sup>提出了一种语言模型 SLAMC, 通过综合词的属性、数据类型、从属范围等信息, 增强了程序语义理解能力, 并在代码提示任务取得了较优的准确性。Bielik 等<sup>[26]</sup>提出了一种新的概率生成模型 PGOH (probabilistic higher order grammar), 其较  $n$ -gram 模型能更多综合上下文信息, 在对抽象语法树节点的代码提示任务中取得了较优结果。在[22]工作的基础上, White 等<sup>[14]</sup>构建了几种 RNN 结构的语言模型, 通过实验对比, 这些方法在代码提示任务上, 较  $n$ -gram 类模型有显著优势。本文针对任务特点构建了一个基于 LSTM 单元的改进 RNN 语言模型 CodeNLM, 该模型可以有效综合历史信息、捕获代码规律, 进而为当前代码段提供符合语法、贴近语义的候选项, 辅助完成编程工作。下文将对本文方法进行详细阐述。

## 3 改进的神经语言模型

面向软件分析任务, 本文提出了一种附加信息引导的神经语言模型 CodeNLM, 其整体架构如图 2 所示。该模型首先依据实际任务收集目标软件项目, 创建语料库; 提取出任务相关的样本附加信息, 即一些非代码信息, 如项目名称、软件类型、软件重要性等级等 (依据任务的不同, 提取不同附加信息); 之后针对源代码文件进行必要的预处理工作; 然后将代码序列及其附加信息一并输入给模型的核心部分, 即一个基于 LSTM 单元的 DNN, 用于捕获代码上下文信息、学习程序语言特征; 最终, 模型由一个 softmax 多分类器预测输入序列的下一词, 以得到的概率分布来刻画整段代码的发生规律。

如图 2 所示, CodeNLM 的 DNN 部分主要包括输入层、数据表示层、特征学习层、决策层和输出层。其中, 数据表示层利用分布式表示方法<sup>[21]</sup>将代码的每个词映射为一个稠密的实值矢量, 语境相似的词将被赋予为相似的矢量; 特征学习层则针对这些矢量序列学习代码语义特征和词序规律; 最后决策层通过综合网络信息预测序列的下一词, 以此刻画代码序列的联合概率分布。为增强模型对软件任务的理解, 本文将任务相关的样本附加信息融入网络。具体地, 首先将附加信息量化为实值向量, 然后分别同词向量输入和 softmax 层输入相融合, 以辅助网络学习与决策。

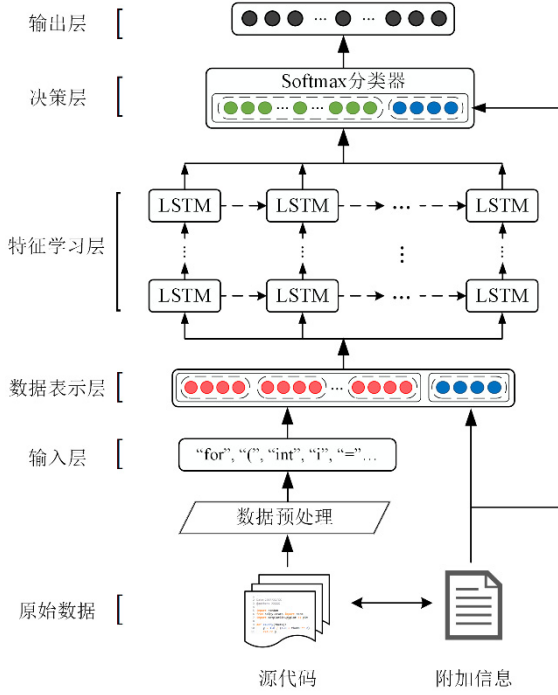


图2 CodeNLM 框架

Fig2 The framework of CodeNLM

下面对 CodeNLM 的主要构成进行详细介绍。

### 3.1 数据预处理

在训练神经网络之前，需要对语料库数据进行必要的预处理工作，具体方法如下所示：

- 步骤1：提取源代码文件。从语料库中提取源代码文件作为模型训练和验证的原始数据。
- 步骤2：生成有效代码。去除源代码文件中的注释和空行，保留有效代码行。
- 步骤3：代码词化（tokenization）。按照语法规则将源代码进行词化处理，即生成词序列，例如“for, '(', 'int', 'i', '=', '0', ';'...”。
- 步骤4：生成词汇表。合并源代码文件，统计独立的词及其词频，创建词汇表。
- 步骤5：替换少见词。将模型处理的词汇总数设为  $V$ ，即保留词频数前  $V$  的词，其余用<RARE>标志替换。
- 步骤6：生成实验数据集。生成固定长度的代码词序列，创建十折交叉验证数据集。

### 3.2 网络结构及改进策略

CodeNLM 的核心结构是一个附加信息引导的多层 RNN。在样本附加信息的辅助下，网络利用 LSTM 单元对长距离信息的建模优势来学习程序语言特征、捕捉代码序列的语义信息。

LSTM 单元由 Schmidhuber 等<sup>[27]</sup>提出，通过使用内存门控机制防止梯度消失问题，克服了早期 RNN 难以记忆长程依赖的问题。LSTM 单元的数学模型为：

$$\begin{aligned} i_t &= \sigma(W_{xi}x_t + W_{hi}h_{t-1} + b_i) \\ f_t &= \sigma(W_{xf}x_t + W_{hf}h_{t-1} + b_f) \\ o_t &= \sigma(W_{xo}x_t + W_{ho}h_{t-1} + b_o) \\ c_t &= f_t \odot c_{t-1} + i_t \odot \tanh(W_{xc}x_t + W_{hc}h_{t-1} + b_c) \\ h_t &= o_t \odot \tanh(c_t) \end{aligned}$$

其中  $\sigma(x) = 1/(1 + e^{-x})$  和  $\tanh(x) = (e^x - e^{-x}) / (e^x + e^{-x})$  为激活函数； $W$  和  $b$  分别表示权重矩阵和偏置向量； $\odot$  表示按元素相乘，即 Hadamard 积； $i_t, f_t, o_t$  分别表示输入门、遗忘门和输出门操作； $x_t, c_t, h_t$  分别表示网络输入、记忆单元状态和隐含层状态； $t$  表示展开时间。

考虑到现有的语言模型在数据学习和网络结构设计方面仍具有一定盲目性，本文提出两种改进策略：附加信息引导策略（策略1）和节点逐层递增策略（策略2），以增强模型对软件任务的理解、改善程序语言建模能力。

#### （1）附加信息引导策略

现有语言模型大都从源代码数据中提取知识，而有助于解决软件任务的部分样本信息（如与项目的对应关系、软件类型、软件重要性等级）难以在代码中体现，导致模型未能对软件进行充分分析。为克服这一问题，本文提出将任务相关的非代码信息融入网络，引导模型更准确地刻画语言特征，提高原有网络对数据学习的针对性。需要指出，应依据具体应用任务，抽取不同附加信息。下文将以代码提示任务为例进行阐述。

具体地，该策略首先将抽取的附加信息以实值向量形式表示，然后分别同词向量输入和 softmax 层输入相融合，用于辅助网络学习与决策。其形式化描述为：

$$\begin{aligned} x_t^1 &= w_t \\ x_t^{softmax} &= h_t^{N_{layer}} \end{aligned} \longrightarrow \begin{aligned} \tilde{x}_t^1 &= [w_t \quad p_t] \\ \tilde{x}_t^{softmax} &= [h_t^{N_{layer}} \quad p_t] \end{aligned} \quad (4)$$

其中  $x$  表示原网络输入， $\tilde{x}$  表示改进后的网络输入， $w$  表示词向量， $p$  表示向量化的样本附加信息， $h$  表示 LSTM 隐含状态， $N_{layer}$  表示隐含层层数， $t$  表示展开时间。

对于不同应用任务，可以为模型引入不同附加信息。本文以代码提示应用为例进行分析。代码提示旨在为编程人员提供准确的代码候选项，以提高开发效率。由于不同的软件项目拥有不同的应用背景及开发过程，致使反映在代码产品上，存在统计差异。为使模型更精准地捕获代码规律、更准确地完成代码提示，本文将样本与项目的关系信息融入网络，如图3所示。针对代码提示应用，式(4)中的附加信息  $p$  即为量化后的项目名称/ID，其量化使用词嵌入技术。具体实现中，通过路径识别出代码文件对应的项目 ID，然后将代码 token 的 ID 序列及其对应的项目 ID 序列输入给神经网络，这一过程是自动化的。



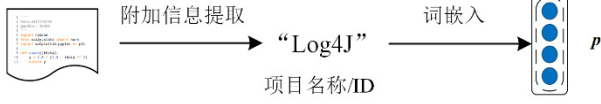


图3 样本附加信息的提取与量化

Fig3 Extraction and quantification of additional information

### (2) 节点逐层递增策略

在神经网络结构设计方面，目前缺乏经验性总结，导致实验盲目、低效。本文通过分析神经语言模型特点，提出一种简单有效的网络规模正则化方法：节点逐层递增策略，即由输入层到特征学习层再到输出层的网络节点数按递增模式设置：

$$\dim_{in} < \dim_h^1 < \dim_h^2 < \dots < \dim_h^{N_{layer}} < \dim_{out}$$

其中， $\dim_{in}$  表示输入层节点数，即词嵌入维度  $\dim_w$ ； $\dim_h^l$  表示隐含层节点数， $N_{layer}$  表示隐含层层数； $\dim_{out}$  表示输出层节点数，即 softmax 分类类别数。

这一经验策略有悖于常见的 DNN 结构设计方法，即整体上网络节点逐渐递减或节点先增后减的设计模式，例如经典的 LeNet5<sup>[28]</sup>、AlexNet<sup>[29]</sup>和 GoogLeNet<sup>[30]</sup>等 DNN。这样的差异主要是因为面向任务的不同。上述这些经典模型主要面向计算机视觉任务，例如在图像识别问题中，考虑处理像素  $256 \times 256$  的图像，那么模型的输入维度为 65,536；输出层为分类类别，其节点一般在 10~5,000 之间，远小于输入维度。因此，网络通常需要逐层压缩数据，以实现特征的抽象。而本文面向的是语言建模任务，模型的输入通常为词向量，维度一般在 100~1,000 之间；输出层对应词汇表中的词，维度一般在 2,000~20,000 之间，远大于输入维度。这正是语言建模任务的一个重要特点。针对这一特点，本文使用各层节点逐层递增的模式设计网络，此方法不仅有助于信息的逐层传递，还能有效限制网络规模、降低过拟合风险，从而提高模型泛化性能和训练速度。

### 3.3 模型训练

深度学习模型在训练过程中常出现过拟合、梯度爆炸或梯度弥散等问题。为此，本文采用以下几种方法进行优化网络训练：

- 1) Dropout 方法<sup>[24]</sup>：在模型训练中，对网络节点按照一定概率将其暂时丢弃，以阻止神经元的共适应问题，是一种有效的抗过拟合方法。
- 2) 梯度剪裁 (gradient clipping)<sup>[31]</sup>：通过限制梯度范数来防止梯度爆炸，方法简单有效。
- 3) 自适应学习率：通过动态学习率优化模型训练、提高泛化性能。本文令梯度下降算法的学习率  $l_{rate}$  随训练轮数动态调整：

$$l_{rate} = \begin{cases} 1, & i < \lambda \\ \eta^{i-\lambda}, & \lambda \leq i \leq N_{epoch} \end{cases}$$

其中  $i, N_{epoch} \in \mathbf{Z}^+$  分别表示训练轮数和最大训练轮数， $\lambda \in [1, N_{epoch}]$  和  $\eta \in (0, 1)$  是调节学习率曲线的参数。

### 3.4 模型评估

语言模型的评价常采用交叉熵 (cross-entropy) 或困惑度 (perplexity) 作为度量指标<sup>[3]</sup>。交叉熵用于衡量估计模型结果与真实概率分布之间的差异情况，其值越小说明模型表现越好。对于给定的语料库  $C = \{S_1, S_2, \dots, S_N\}$  ( $C \subseteq \Omega^*$ ) 和语言模型  $M$ ，则模型在该语料库上的交叉熵值可由下式计算：

$$\begin{aligned} CE_M(C) &= -\frac{1}{l_C} \log_2 P_M(C) \\ &= -\frac{1}{l_C} \sum_{j=1}^N \log_2 P_M(S_j) \\ &= -\frac{1}{l_C} \sum_{j=1}^N \sum_{i=1}^{l_j} \log_2 P_M(w_{j,i} | w_{j,1}, \dots, w_{j,i-1}) \end{aligned}$$

其中  $S_j \in C$  表示语料库中的序列， $N$  表示序列总数， $l_j$  表示序列  $S_j$  的长度， $l_C = \sum_{j=1}^N l_j$  表示语料库上序列的总长度， $P_M(S_j)$  表示模型  $M$  求得的序列  $S_j$  的发生概率， $P_M(C) = \prod_{j=1}^N P(S_j)$  表示语料库的发生概率。

困惑度指标是交叉熵的指数形式，同样地其值越小说明模型表现越好。对于给定的语料库  $C$  和语言模型  $M$ ，则模型在该语料库上的困惑度值可由下式得到：

$$PP_M(C) = 2^{CE_M(C)}$$

其中  $CE_M \in [0, +\infty)$ ， $PP_M \in [1, +\infty)$ 。

## 4 实验与结果分析

本节针对实际项目数据，通过实验分析和结果对比验证 CodeNLM 模型的有效性，其中 4.1 节介绍实验数据集，4.2 节说明实验环境及模型参数，4.3 节将从语言模型性能、代码提示准确度和长距离信息学习能力三个方面综合评价本文方法。

### 4.1 实验数据集

为便于模型对比，本文选用文献[4]-[6]和[16]均使用的 Java 公开数据集，其信息如表 1 所示。该数据集包含 9 个 GitHub 开源项目共 9060 个 java 文件，总代码行数约 203 万，有效代码行约 114 万。

表1 实验数据集

Table1 Information of the evaluated dataset				
项目名	版本	文件数	代码行	有效代码行
Ant	20110123	1,193	254,457	128,266
Batik	20110118	1,657	367,293	195,467
Cassandra	20110122	545	135,992	97,029
Log4J	20101119	353	68,528	34,479
Lucene	20100319	2,279	429,957	266,573
Maven2	20101118	386	61,622	38,593
Maven3	20110122	847	114,527	70,462
Xalan-J	20091212	973	349,837	171,619
Xerces	20110111	827	257,572	138,841
总计	—	9,060	2,039,785	1,141,329

#### 4.2 实验环境与模型参数

实验运行环境为 Windows 7 64 位操作系统、Intel i7-6700K CPU、16G 内存及 1 块 GeForce GTX 1070 显卡，使用的深度学习工具为 TensorFlow 1.0.0<sup>1</sup>，采用的代码词化工具来自文献[10]。

实验中 CodeNLM 的超参数设置如下：词汇表大小  $V=5000$ ；词向量维度  $dim_w=500$ ；项目属性的嵌入维度  $dim_p=5$ ；LSTM 层数  $N_{layer}=2$ ；LSTM 层的隐含节点  $dim_h$  分别为 800 和 1200；dropout 方法的节点保留概率  $p_{keep}=0.6$ ；按时间反向传播展开步长  $N_{step}=35$ ；批处理的大小  $N_{batch}=30$ ；训练轮数  $N_{epoch}=15$ ；允许的最大梯度范数  $grad_{max}=5$ ；LSTM 单元遗忘门偏置的初始值为 0.0，网络其他参数的初始值按  $[-0.08, 0.08]$  的均匀分布随机设置；模型优化使用梯度下降算法，自适应学习率  $l_{rate}$  参数  $\lambda=10$ 、 $\eta=0.5$ 。不难看出，CodeNLM 的“输入层-隐含层-输出层”（即“词嵌入层-LSTM 层-softmax 层”）维度依次为 505, 800, 1205 和 5000，各层节点逐层递增。

#### 4.3 实验结果与对比

本小节从语言模型性能、代码提示准确度和长距离信息学习能力三个方面评价本文方法。

##### 4.3.1 语言模型性能评估与对比

针对实验数据集，采用十折交叉验证方法验证 CodeNLM 的有效性，具体结果如图 4 所示。本文方法在测试集上取得的最优困惑度为 2.650、最劣值为 2.886、均值为 2.802，模型平均耗时约 11,102s。

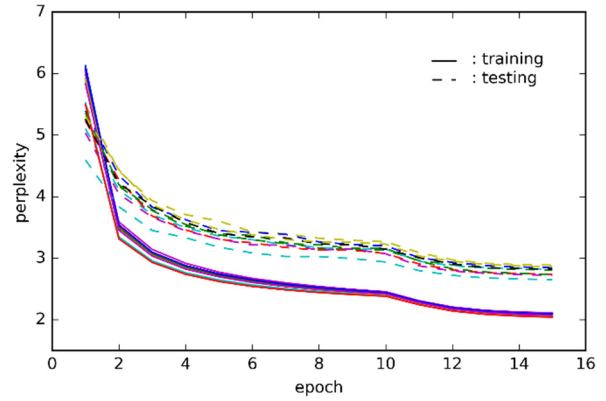


图4 损失函数曲线（十折交叉验证）

Fig4 Loss curves (10-fold cross validation)

由图 4 可以看出，模型经过 15 轮数据训练，困惑度指标逐渐降低并趋于收敛，其中在第 10 轮曲线明显变化的原因是由于学习率的下降。通过学习率的自适应调整，全局搜索和局部搜索实现平衡，模型得到充分训练。此外，10 次实验的损失曲线十分接近，测试集与训练集的困惑度相差不大，说明 CodeNLM 具有良好的稳定性和泛化能力。

表 2 展示了不同语言模型在此数据集上的性能，进一步通过模型对比验证 CodeNLM 的有效性。文献[4]和[5]分别采用了  $n$ -gram 模型和 cache  $n$ -gram 模型进行代码分析，此两种模型获得的最优交叉熵介于 3~4 之间<sup>[5]</sup>，即困惑度均大于 8（ $2^3=8$ ）。文献[15]面向软件代码提出了一种基于 LSTM 单元深度语言模型 DeepSoft，在相同数据集上，其求得的最优困惑度为 4.72<sup>[16]</sup>。文献[24]提出了一种有效的正则化 LSTM 网络，是语言建模领域的一项经典工作并被多次引用，其不同规模网络在表 1 数据集上的最优困惑度为 3.044（具体参数设置见文献[24]）。本文方法得到的指标值为 2.802，明显优于上述多种模型结果，表明了 CodeNLM 在刻画代码复杂特征、学习程序语言规律方面具有优越性。其中 CodeNLM\*表示未使用策略 1（即附加信息引导策略）；CodeNLM\*\*表示未使用策略 1 和策略 2（即节点逐层递增策略），其词向量和隐含层维度均为 800，其他方面较 CodeNLM 不变。这两个版本模型得到的困惑度指标分别为 2.881 和 2.946，可以看出本文提出的两种策略都起到了明显的正向促进作用，提高了语言模型的学习能力。但另一方面也可以看到，深度学习模型带来了过高的计算成本，其中本文 CodeNLM 的训练时间达到了 185 分钟，远大于  $n$ -gram 类模型（见表 2，其中  $n$ -gram 类模型数据来自文献[5]），这种显著的开销对比也可在文献[11]中得到印证。因此，如何构建低功耗的 DNN 仍是当前一个重要挑战和研究方向<sup>[2]</sup>。

<sup>1</sup> <https://www.tensorflow.org/>

表2 不同模型的困惑度指标值

Table2 Perplexity values of the investigated models

模型	<i>n</i> -gram 类模型		DeepSoft [16]	Regularized LSTMs [24]			本文方法		
	<i>n</i> -gram [4]	Cache LM [5]		Small	Medium	Large	CodeNLM**	CodeNLM*	CodeNLM
困惑度	12~14	8~9	4.720	3.739	3.044	3.152	2.946	2.881	<b>2.802</b>
训练时间	0.9 min	2.9 min	—	79 min	457 min	1750 min	150 min	180 min	185 min

注：\*表示未使用策略 1；\*\*表示未使用策略 1 和策略 2。

#### 4.3.2 代码提示任务结果与对比

本文选用两类常用指标衡量模型的代码提示性能<sup>[4]-[6]</sup>：

平均排序倒数（mean reciprocal rank, MRR）和 Top-*n* 准确率。指标的计算方法如下：

$$\text{MRR} = \frac{1}{|T|} \sum_{i=1}^{|T|} \frac{1}{\text{rank}_i}$$

$$\text{Top-}n = \frac{1}{|T|} \sum_{i=1}^{|T|} r_i$$

其中

$$r_i = \begin{cases} 1, & \text{rank}_i \leq n \\ 0, & \text{rank}_i > n \end{cases}$$

式中 *T* 表示测试的 token 总数， $\text{rank}_i$  表示正确结果  $t_i$  获得的推荐名次。Top-*n* 准确率衡量了正确结果出现在前 *n* 个候选项中的可能性；MRR 指标则表示正确结果获得的平均排名的倒数，例如 MRR=0.5 表明平均意义上正确结果出现在第二推荐位置。

表 3 罗列了不同语言模型在表 1 数据集上的代码提示结果，进一步评价 CodeNLM 的有效性。

表3 不同模型在代码提示任务中的性能对比

Table3 Performance of the investigated models for code suggestion

性能 指标	<i>n</i> -gram类模型			Regularized LSTMs [24]			本文方法		
	<i>n</i> -gram [4]	Cache LM [5]	SLAMC [6]	Small	Medium	Large	CodeNLM**	CodeNLM*	CodeNLM
平均	MRR	52.8%	67.6%	—	68.5%	73.8%	72.9%	75.4%	<b>77.2%</b>
	Top-1	46.3%	—	65.6%	58.9%	64.9%	64.0%	66.9%	<b>68.5%</b>
	Top-5	56.8%	—	78.2%	80.3%	84.9%	83.8%	86.0%	<b>88.0%</b>

注：\*表示未使用策略 1；\*\*表示未使用策略 1 和策略 2。

由表 3 可以看出，神经语言模型整体较传统 *n*-gram 类模型有显著优越性，3 种准确度指标均有明显提升（为避免二次实验引入偏差，*n*-gram 类模型结果使用了原文数据，其中文献[6]结果是此类模型得到的最佳结果）。这是由于 LSTM 模型在数据表示、非线性特征变换、长距离信息学习等方面具有突出优势，其中 *n*-gram 类模型可综合的 context 信息长度远小于神经类模型的（4.3.3 节将重点分析）。在全部 9 个项目上，本文方法均得到了最优性能，其中 MRR 指标均值为 77.2%，较其他模型（绝对）增长 3.4%~24.4%；Top-1 指标均值为 68.5%，增长 2.9%~22.2%；Top-5 指标均值为 88.0%，增长 3.1%~31.2%。表 3 还给出了 CodeNLM\*和 CodeNLM\*\*模型结果。以 MRR 指标为例，它们分别得到了 75.4%和 76.2%，同样优于其他方法，但劣于使用了改进策略的 CodeNLM。

通过此次对比实验可以得出，CodeNLM 能够捕获代码

复杂特征、学习程序语言规律，并较好地完成了代码提示工作，较对比模型存在显著优势。同时，实验结果也验证了本文提出的 2 种改进策略能够提升模型性能，适用于程序语言建模及相关软件分析任务。

#### 4.3.3 长距离信息学习能力分析

学习长距离信息一直是语言建模任务的一个难题。作为软件的核心构成，源代码中蕴含了复杂的语法及语义信息。因此，对于语言模型来说，拥有良好的长程依赖学习能力至关重要。本小节考察 CodeNLM 此方面的性能。

图 5 展示了由本文方法得到的序列下一词的估计概率分布，即给定长度为 *l* 的序列  $S = w_1 w_2 \cdots w_l$ ，模型得到的  $P(x = w_{l+1} | w_1, w_2, \cdots, w_l)$  分布，这里 *x* 表示模型的预测结果， $w_{l+1}$  表示序列 *S* 的（真实）下一词， $w_1 w_2 \cdots w_l$  称作  $w_{l+1}$  的“历史”。

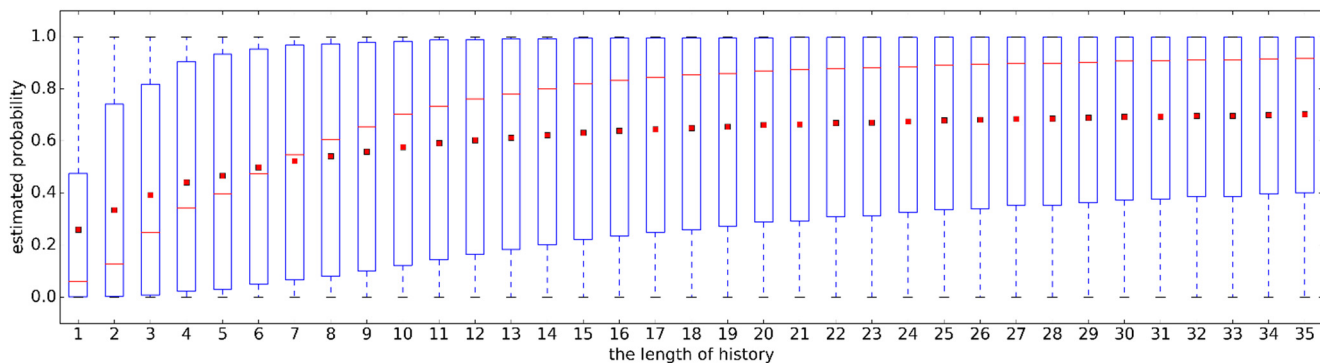


图5 序列下一词的估计概率

Fig5 The estimated probability of the next real token of a given sequence

由图 5 可以看出, 随着给定历史信息的增加, CodeNLM 对序列下一词  $w_{t+1}$  的预测愈加准确, 得到的估计概率分布呈上升趋势。当输入给模型的序列长度为 1 时,  $w_{t+1}$  获得的估计概率均值仅为 0.268; 当历史长度大于 18 后, 估计概率均值超过了 0.650 并趋于平缓; 当序列长度达到 35 时, 模型赋予  $w_{t+1}$  的概率均值为 0.702。其中, 估计概率的中位数变化更为明显, 其随给定历史信息的增长, 由 0.060 升高到 0.917。而对于传统的  $n$ -gram 语言模型, 由于自身固有的维数灾难问题, 其难以处理长文本数据。在文献[4]-[6]中, 此类模型在处理代码提示任务时分析的历史长度不超过 4。因此从上述分析不难得出, CodeNLM 具有优越的长距离信息学习能力, 能有效综合历史输入, 进而捕获语言规律。

**结束语** 本文面向程序语言提出了一种附加信息引导的 RNN 语言模型 CodeNLM。该模型以源代码的词序列为分析对象, 利用多层 LSTM 网络进行特征学习, 克服了传统模型固有的维数灾难问题, 增强了长距离信息学习能力。在神经网络的构建上, 提出了两种改进策略: (1) 将任务相关的非代码信息融入网络, 引导模型更准确地刻画代码特征, 改善了原有方法对数据学习的盲目性 (以代码提示任务为例, 引入的附加信息为代码与项目的关系信息); (2) 针对语言模型特点, 令网络节点逐层递增, 该策略能够限制模型规模、降低过拟合风险。

选取了含 2.03M 代码行的 Java 数据集用于模型验证。在代码提示任务 (token 级) 中, CodeNLM 的预测性能优于对比的  $n$ -gram 类模型和神经语言模型, Top-1、Top-5 和 MRR 指标得到明显提升 (2.9%~31.2%)。实验结果表明, 本文方法能更精准地刻画语言规律、捕获代码特征, 较好地实现了程序语言建模和代码提示任务。下一步工作将考

虑不同方面附加信息的提取与融合, 并应用模型于缺陷检测等其他软件分析问题。

## 参考文献

- [1] Wang Qian-Xiang, Zhang jian, Xie Tal, et al. Software analysis: techniques, applications and trends/CCF 2015-2016 China Computer Science and Technology Development Report[M]. Beijing: China Machine Press, 2016: 55-113. (in Chinese)  
(王千祥, 张健, 谢涛, 等. 软件分析: 技术、应用与趋势/ CCF 2015-2016 中国计算机科学技术发展报告[M]. 北京: 机械工业出版社, 2016: 55-113.)
- [2] Zhang Xian, Ben Kerong. Application of deep learning methods in software analysis [J]. Computer Engineering and Science, 2017(12):2260-2268. (in Chinese)  
(张献, 贲可荣. 深度学习方法在软件分析中的应用[J]. 计算机工程与科学, 2017(12):2260-2268.)
- [3] Jurafsky D, James H M. Speech and language processing (2nd ed) [M]. Upper Saddle River: Pearson/Prentice Hall, 2009.
- [4] Hindle A, Barr E T, Su Z, et al. On the naturalness of software[C]. //Proceedings of the 34th International Conference on Software Engineering. IEEE, 2012:837-847.
- [5] Tu Z, Su Z, Devanbu P. On the localness of software[C]. //Proceedings of the 22nd ACM SIGSOFT International Symposium on Foundations of Software Engineering. ACM, 2014: 269-280.
- [6] Nguyen T T, Nguyen A T, Nguyen H A, et al. A statistical



- semantic language model for source code[C]. //Proceedings of the 9th Joint Meeting on Foundations of Software Engineering. ACM, 2013: 532-542.
- [7] Yang Y, Jiang Y, Gu M, et al. A language model for statements of software code[C]. //Proceedings of the 32nd IEEE/ACM International Conference on Automated Software Engineering, 2017: 682-687.
- [8] Nguyen A T, Hilton M, Codoban M, et al. API code recommendation using statistical learning from fine-grained changes[C]. //Proceedings of the 24th ACM SIGSOFT International Symposium on Foundations of Software Engineering. ACM, 2016: 511-522.
- [9] Allamanis M, Barr E T, Bird C, et al. Learning natural coding conventions[C]. //Proceedings of the 22nd ACM SIGSOFT International Symposium on Foundations of Software Engineering. ACM, 2014: 281-293.
- [10] Ray B, Hellendoorn V, Godhane S, et al. On the naturalness of buggy code[C]. //Proceedings of the 38th International Conference on Software Engineering. ACM, 2016: 428-439.
- [11] Bielik P, Raychev V, Vechev M. Program synthesis for character level language modeling[C]. //Proceedings of the 5th International Conference on Learning Representations, 2017.
- [12] Oda Y, Fudaba H, Neubig G, et al. Learning to generate pseudo-code from source code using statistical machine translation[C]. //Proceedings of the 30th IEEE/ACM International Conference on Automated Software Engineering. IEEE, 2015: 574-584.
- [13] Hirschberg J, Manning C D. Advances in natural language processing[J]. Science, 2015, 349 (6245) : 261-266.
- [14] White M, Vendome C, Linares-Vásquez M, et al. Toward deep learning software repositories[C]. //Proceedings of the 12th IEEE/ACM Working Conference on Mining Software Repositories, 2015: 334-345.
- [15] Dam H K, Tran T, Grundy J, et al. DeepSoft: a vision for a deep model of software[C]. //Proceedings of the 24th ACM SIGSOFT International Symposium on Foundations of Software Engineering. ACM, 2016: 944-947.
- [16] Dam H K, Tran T, Pham T. A deep language model for software code[C]. //Workshop on naturalness of software, collocated with the 24th International Symposium on Foundations of Software Engineering. ACM, 2016.
- [17] Allamanis M, Barr E T, Bird C, et al. Suggesting accurate method and class names[C]. //Proceedings of the 10th Joint Meeting on Foundations of Software Engineering. ACM, 2015: 38-49.
- [18] Yin Chun-Lin, Wang Wei, Li Tong, et al. Using RNNLM to conduct topic oriented feature location method[J]. Journal of Frontiers of Computer Science and Technology, 2017, 11 (10) : 1599-1608. (in Chinese)  
(尹春林, 王炜, 李彤, 等. 利用 RNNLM 面向主题的特征定位方法[J]. 计算机科学与探索, 2017, 11 (10) : 1599-1608.)
- [19] Zhang Xian, Ben Kerong and Zeng Jie. Cross-entropy: a new metric for software defect prediction[C]. //Proceedings of the 18th IEEE International Conference on Software Quality, Reliability and Security, pp. 111-122, 2018. <https://github.com/TOM-ZXian/-Cross-entropy-metric-of-code-for-defect-prediction>.
- [20] Bengio Y, Ducharme R, Vincent P. A neural probabilistic language model[C]. //Proceedings of the 15th Annual Conference on Neural Information Processing Systems, 2001:932-938.
- [21] Yann LeCun, Yoshua Bengio, Geoffrey Hinton. Deep Learning[J]. Nature, 2015, 512 (7553) : 436-444.
- [22] Mikolov T, Karafiát M, Burget L, et al. Recurrent neural network based language model[C]. //Proceedings of the Annual Conference of the International Speech Communication Association, 2010: 1045-1048.
- [23] Mikolov T, Yih W, Zweig G. Linguistic regularities in continuous space word representations[C]. //Proceedings of the Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies. 2013: 746-751.
- [24] Zaremba W, Sutskever I, Vinyals O. Recurrent neural network regularization[J]. arXiv preprint arXiv:1409.2329,

2014.

- [25] Salehinejad H, Baarbe J, Sankar S, et al. Recent Advances in Recurrent Neural Networks[J]. arXiv preprint arXiv:1801.01078, 2018.
- [26] Bielik P, Raychev V, Vechev M. PHOG: probabilistic model for code[C]. // Proceedings of the 33rd International Conference on Machine Learning. 2016: 2933-2942.
- [27] Hochreiter S, Schmidhuber J. Long short-term memory[J]. Neural Computation, 1997, 9(8): 1735-1780.
- [28] LeCun Y, Bottou L, Bengio Y, et al. Gradient-based learning applied to document recognition[J]. Proceedings of the IEEE, 1998, 86(11): 2278-2324.
- [29] Krizhevsky A, Sutskever I, Hinton G E. Imagenet classification with deep convolutional neural networks[C]. //Proceedings of the 26th Annual Conference on Neural Information Processing Systems. 2012: 1097-1105.
- [30] Szegedy C, Liu W, Jia Y, et al. Going deeper with convolutions[C]. //Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition. 2015: 1-9.
- [31] Pascanu R, Mikolov T, Bengio Y. On the difficulty of training recurrent neural networks[C]. //Proceedings of the 30th International Conference on Machine Learning. 2013: 1310-1318.



**Zhang Xian**, born in 1990. PhD

candidate. Student member of China Computer Federation. His research interests include software analysis, software defect mining, and machine learning.



**Ben Kerong**, born in 1963, PhD,

professor, PhD supervisor. Senior member of China Computer Federation. His main research interests include software quality assurance and artificial intelligence.