

0 摘要

针对大规模线上班级作业给出反馈, 评估最终作业和给陷入困难学生提示是非常困难的, 这些班级学生成千甚至上百万个。我们提出神经网络编码代码, 从嵌入预设条件空间映射到嵌入后置条件空间, 利用线性映射作为特征提出规模化反馈的算法。

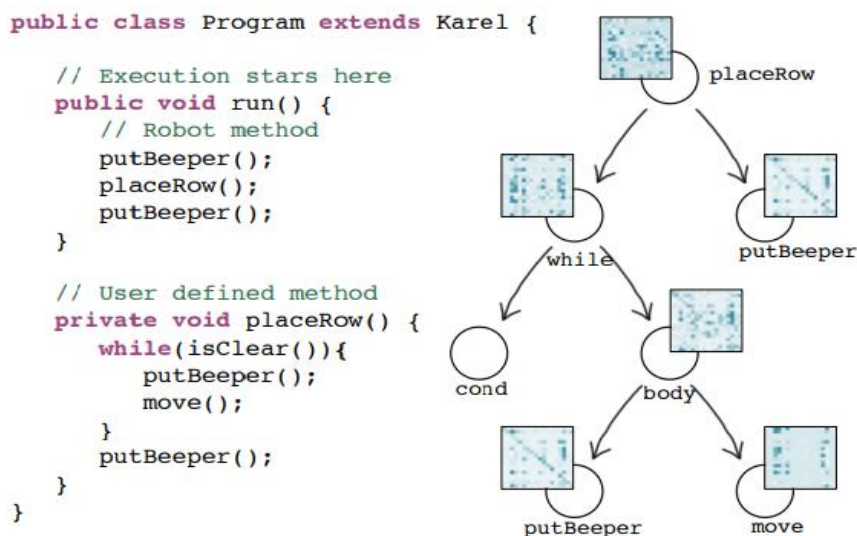
1 介绍

直接对代码中的内容应用机器学习算法很难。代码表示, 应用标准统计方法不能直接生成语法树, 树的编辑距离矩阵也区分度不够, 难于对代码提出准确反馈。因为有相似语法树的代码, 可能有完全不同的意思, 需要不同的评论。虽然单元测试常用于检测答案是否正确, 但者并不适用于对学生提供帮助, 给出及时的答案, 单元测试也不能对文体内容给出反馈。

本文有两个目标, 第一个是自动给出学生提交代码的特征嵌入向量, 捕捉功能及样式元素, 这在传统监督机器学习系统中很容易使用。第二个是利用特征学习如何给学生自动反馈。受最近在 NLP 和视觉领域深度学习学习特征研究的启发, 本文提出新的神经网络, 联合优化代码嵌入向量及特征空间记忆状态。图 1 是示例代码和相应矩阵嵌入。

为了收集数据, 我们利用程序是可以执行这一事实, 基于特定输入 (预先条件) 评估代码, 看输出状态 (后置条件)。对代码和它的组成部门, 我们收集一些输入输出映射。这些数据作为训练集, 我们得到代码共享表示。为了衡量代码嵌入, 我们测试增强教师反馈的能力。利用真实的数据训练, Code.org Hour of Code, 应用到斯坦福编程方法论作业集。这些代码结构完整, 使用者不能自己定义变量。

我们的贡献有以下几点。首先, 我们提出计算代码特征的方法, 捕捉功能及文体元素。我们的模型同时将代码的前置及后置条件嵌入特征空间, 代码看作空间的线性映射。第二, 我们的代码特征是有用的, 自动为大规模课程学生生成指导反馈。第三, 我们的方法在大规模数据集是有效的。机器学习学习中学习代码的嵌入研究很多, 但嵌入能生成反馈对计算机科学教育作用重大。



3 嵌入霍尔三元组

核心问题是把代码表示成固定维度空间，可以直接应用于传统监督学习算法。虽然代码有很多维度，像空间和时间复杂度等，我们首先聚焦最基本的方面--功能。捕捉代码功能忽视里其他有用的方面，本文在后面根据语法树子树建模子程序功能，以此获取样式元素。给定程序 A 和输入 P ，学习 A 的特征，这对预测程序输出有用。

本文提出利用训练集 (P, A, Q) 三元组学习代码特征。如果 f_P 和 f_Q 是输入 P 和输出 Q 的 m 维非线性特征表示，有如下等式。

$$f_Q = M_A \cdot f_P.$$

相关系数矩阵 M 是代码 A 的特征代表，是代码嵌入矩阵。我们需要学习特征空间映射 f 和线性映射 M ，使所有输入输出上式都成立。

学习输入输出的非线性嵌入函数 f ，可以获得丰富的非线性关系，就像核方法允许非线性决策边界一样。

3.1 神经网络编码及解码状态

假设输入有基准编码器转为 d 维函数，定义为 P 。同样的，输出也有基编码器 Q 。

前提是输入及输出空间共享一个基编码器，解码器参数也一样。受非线性自动编码启发，编码器是将输入条件 P 编码为非线性特征表示 f_P 。

$$f_P = \phi(W^{enc} \cdot P + b^{enc}),$$

特征表示 f_P 用来解码：

$$\begin{aligned}\hat{Q} &= \psi(W^{dec} \cdot f_Q + b^{dec}), \\ &= \psi(W^{dec} \cdot M_A \cdot f_P + b^{dec}).\end{aligned}$$

3.2 代码嵌入的非参数模型

为了编码代码嵌入矩阵，提出非参数模型，每个代码与自己嵌入矩阵相关。

$$\Theta = \{W^{dec}, W^{enc}, b^{enc}, b^{dec}\} \cup \{M_i : i = 1, \dots, m\}.$$

为了学习这些参数，最小化下面三项的和：（1）预测损失，衡量给定输入预测代码输出的质量（2）编码损失：衡量以重构给定输入的编码和解码参数质量。（3）正则化项，本文选择 L2 正则。

$$L(\Theta) = \frac{1}{n} \sum_{i=1}^n \ell^{pred}(Q_i, \hat{Q}_i(P_i, A_i; \Theta)) \\ + \frac{1}{n} \sum_{i=1}^n \ell^{auto}(P_i, \hat{P}_i(P_i, \Theta)) + \frac{\lambda}{2} \mathcal{R}(\Theta),$$

这个最优化问题是试图找到状态空间好的共享表示，找到输入及输出的非线性编码使代码在新的特征空间看上去线性的。不考虑动量的小批量随机梯度下降计算梯度实现联合最优化。超参数随机搜索，学习率按 **Adagrad** 设定。学习状态空间自动编码器，对每个代码实现向量值的岭回归，抽取输入到输出特征的映射矩阵。

4 反馈传播

第一阶段，算法挑选示例代码的子程序，实现有限集标注。第二部分，人工标注作为标签，预测无标签代码的反馈。反馈有 **N** 种情况，这样就将反馈转化为 **N** 类分类问题。

4.1 递归嵌入包含结构信息。

利用嵌入矩阵预测反馈，同时考虑语法树所有组成子树，基于递归神经网络提出模型。

语法树 **j** 节点子树用矩阵表示，这由 **j** 的叶子的子树表示及 **j** 的子树嵌入矩阵组合而成。**K-means** 给出最终反馈。

$$a^{(j)} = \phi \left(\sum_{i=1}^{a_\omega} W_i^\omega \cdot a^{(c_i[j])} + b^\omega + \mu M_j^{NPM} \right)$$

总体上，实验结果优于 **RNN, bag of trees** 和单元测试。