



算法设计与分析

实验报告（二）查重

姓 名	熊恪峥
学 号	22920202204622
日 期	2022年3月23日
学 院	信息学院
课程名称	算法设计与分析

实验报告（二）查重

目录

1	问题描述	1
2	实现思路	1
2.1	两点假设	1
2.2	算法	1
2.3	结果分析	2
A	附录：代码实现	4
A.1	矩阵连乘	4
A.2	查重程序	5

1 问题描述

2 实现思路

要求给定两个程序，判断它们的相似性。显然，程序的相似性和代码字符串的相似性无关，而与实际执行逻辑的相似性有关，例如1中有两段代码本身不尽相同的代码，但是执行逻辑完全一致。那么最准确的方式是进行DFA(Data Flow Analysis)和CFA(Control Flow Analysis)，对于相似的程序它们应当能相当准确地反映出相似度。这正是现代IDE对重复代码给出修改建议的方式。但这种方式实现相当复杂，本程序通过对问题进行简化有效地实现了**基于语义**的代码相似性判断。

图 1: 逻辑相同但代码本身差异较大的代码

1	<code>int main()</code>	1	<code>int main()</code>
2	<code>{</code>	2	<code>{</code>
3	<code> bool a = true;</code>	3	<code> int the_flag = 1;</code>
4	<code> if (a)</code>	4	<code> if (the_flag)</code>
5	<code> {</code>	5	<code> {</code>
6	<code> printf("helloworld");</code>	6	<code> puts("helloworld");</code>
7	<code> }</code>	7	<code> }</code>
8	<code> else</code>	8	<code> else</code>
9	<code> {</code>	9	<code> {</code>
10	<code> printf("worldhello");</code>	10	<code> puts("worldhello");</code>
11	<code> }</code>	11	<code> }</code>
12	<code> return 0;</code>	12	<code> return 0;</code>
13	<code>}</code>	13	<code>}</code>

2.1 两点假设

为了简化问题，首先进行以下两个假设

- 程序的抽象语法树(AST, Abstract Syntax Tree)和实际执行逻辑高度相关
- 抽象语法树中的语句节点和表达式节点是所有节点中和实际执行逻辑最相关的两类节点

根据这两点假设，通过从程序编译时的抽象语法树的语句(Statement)节点和表达式(Expression)节点序列中寻找最长公共子序列可以有效地衡量程序的逻辑相似性。定义逻辑相似度 s ，其中 AST_j 是程序代码 j 的抽象语法树的语句节点和表达式节点序列

$$s = \frac{|LCS_{AST_1, AST_2}|}{\max\{|AST_1|, |AST_2|\}}$$

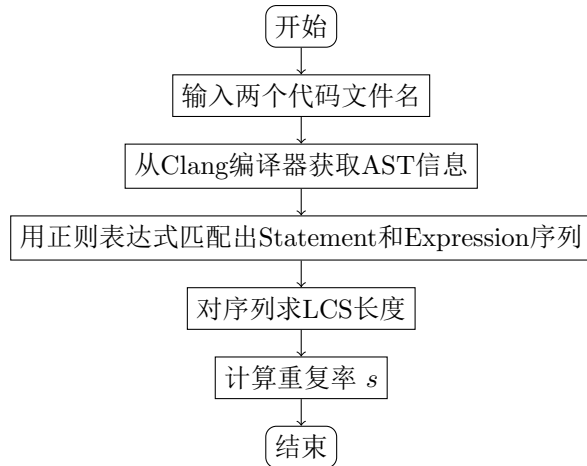
2.2 算法

根据以上分析，首先需要获得AST中Statement和Expression组成的序列。这里借助Clang编译器。通过观察可以发现输出结果中Statement和Expression都符合正则表达式(1)。

$$([a - z A - Z \backslash_][0 - 9 a - z A - Z \backslash_]* Expr)([a - z A - Z \backslash_][0 - 9 a - z A - Z \backslash_]* Stmt) \quad (1)$$

因此可以用正则表达式匹配输出来获得上述序列。

图 2: 查重流程图设计



根据以上分析可以得到算法流程图 2，然后可以实现核心部分如代码 1，完整代码见附录：代码实现节中的代码 3。该实现依赖Clang编译器。

代码 1: 核心部分代码

```

1 def lcs(s1, s2):
2     f = [[0] * len(s1) * 2] * len(s2) * 2
3
4     for i in range(1, len(s1) + 1):
5         for j in range(1, len(s2) + 1):
6             if s1[i - 1] == s2[j - 1]:
7                 f[i][j] = 1 + f[i - 1][j - 1]
8             else:
9                 f[i][j] = max(f[i - 1][j], f[i][j - 1])
10
11     return f[len(s1)][len(s2)]
12
13
14 def duplication_check(file1: str, file2: str):
15     ast1 = str(shell(['clang -cc1 -ast-dump {}'.format(file1)]))
16     ast2 = str(shell(['clang -cc1 -ast-dump {}'.format(file2)]))
17
18     elems1 = list([i[0] if i[0] != '' else i[1] for i in
19                    re.findall(r'([a-zA-Z\_][0-9a-zA-Z\_]*Expr)|([a-zA-Z\_][0-9a-zA-Z\_]*Stmt)', ast1)
20                    ])
21     elems2 = list([i[0] if i[0] != '' else i[1] for i in
22                    re.findall(r'([a-zA-Z\_][0-9a-zA-Z\_]*Expr)|([a-zA-Z\_][0-9a-zA-Z\_]*Stmt)', ast2)
23                    ])
24
25     print("Repeat Rate: {}".format((lcs(elems1, elems2) / max(len(elems1),
26                                                                    len(elems2))) * 100.0))
  
```

使用以上程序对图 1 中的代码进行查重，输出的重复率是81.25%，可以看出比起直接比较代码文面，该算法有效地反映出了底层逻辑的相似性，排除了修改变量名等传统降重方法造成的干扰。

2.3 结果分析

根据上述测试结果可以得出结论，通过AST得节点序列可以有效地刻画程序的相似度。然而这种方法的

思路依然是使用更高层级的“形式相似性”近似“逻辑相似性”，它考虑了一定程度的语义信息。因此，可以使用更为高级的降重技巧规避，例如

- 改变语句顺序
- 将递归结构改为非递归结构
- 将一部分代码移动到一个子过程中

因此这种方式相比于正规的程序静态分析手段而言还是有不足的。然而这种方法实现计算简单，计算量较少，在实际使用中有一定的优势。

A 附录：代码实现

A.1 矩阵连乘

代码 2: 矩阵连乘

```
1 p = list([5, 10, 3, 12, 5, 50, 6])
2 N = 6
3
4 m = list([ list([0x7fffffff for i in range(0, N + 1)]) for j in range(0, N + 1)])
5 s = list([ list([0x7fffffff for i in range(0, N + 1)]) for j in range(0, N + 1)])
6
7
8 def pretty_print(i, j):
9     if i == j:
10         print('A{}'.format(i), end=' ')
11     else:
12         print("(", end=' ')
13         pretty_print(i, s[i][j])
14         pretty_print(s[i][j] + 1, j)
15         print(")", end=' ')
16
17
18 for i in range(1, N + 1):
19     m[i][i] = 0
20
21 for i in range(N, 0, -1):
22     for j in range(i, N + 1):
23         if i == j:
24             m[i][j] = 0
25         else:
26             for k in range(i, j):
27                 val = m[i][k] + m[k + 1][j] + p[i - 1] * p[j] * p[k]
28                 if m[i][j] > val:
29                     m[i][j] = val
30                     s[i][j] = k
31
32 for i in range(1, N + 1):
33     for j in range(1, N + 1):
34         print("inf" if m[i][j] == 0x7fffffff else m[i][j], end=' ' if j != N else '\n')
35
36 for i in range(1, N + 1):
37     for j in range(1, N + 1):
38         print("None" if s[i][j] == 0x7fffffff else s[i][j], end=' ' if j != N else '\n')
39
40 pretty_print(1, N)
```

A.2 查重程序

注意：该实现依赖 *Clang* 编译器

代码 3: 查重程序

```
1 import subprocess
2 from typing import Final
3 import re
4 import argparse
5
6
7 def shell(command):
8     try:
9         return subprocess.check_output(command, shell=True, stderr=subprocess.STDOUT).stdout
10    except subprocess.CalledProcessError as exc:
11        return exc.output
12
13
14 def lcs(s1, s2):
15     f = [[0] * len(s1) * 2] * len(s2) * 2
16
17     for i in range(1, len(s1) + 1):
18         for j in range(1, len(s2) + 1):
19             if s1[i - 1] == s2[j - 1]:
20                 f[i][j] = 1 + f[i - 1][j - 1]
21             else:
22                 f[i][j] = max(f[i - 1][j], f[i][j - 1])
23
24     return f[len(s1)][len(s2)]
25
26
27 def duplication_check(file1: str, file2: str):
28     ast1 = str(shell(['clang -cc1 -ast-dump {}'.format(file1)]))
29     ast2 = str(shell(['clang -cc1 -ast-dump {}'.format(file2)]))
30
31     elems1 = list([i[0] if i[0] != '' else i[1] for i in
32                    re.findall(r'([a-zA-Z\_][0-9a-zA-Z\_]*Expr)|([a-zA-Z\_][0-9a-zA-Z\_]*Stmt)', ast1)
33                    ])
34     elems2 = list([i[0] if i[0] != '' else i[1] for i in
35                    re.findall(r'([a-zA-Z\_][0-9a-zA-Z\_]*Expr)|([a-zA-Z\_][0-9a-zA-Z\_]*Stmt)', ast2)
36                    ])
37
38
39     print("Repeat Rate: {}".format((lcs(elems1, elems2) / max(len(elems1),
40                                                                    len(elems2))) * 100.0))
41
42
43 if __name__ == "__main__":
44     parser = argparse.ArgumentParser(prog="dupcheck",
45                                     usage='%(prog)s [options] file1 file2',
46                                     description="Duplication checker")
47     parser.add_argument("file1")
48     parser.add_argument("file2")
49
50     args = parser.parse_args()
51
52     duplication_check(str(args.file1), str(args.file2))
```

