



ソースコードの盗作判定システム SA-Plag

研究駆動コース 石川琉聖

SA-Plag とは？

競技プログラミングのソースコードから学習した
AI盗作判定システムとそのWeb API

製作背景と課題

プログラミングコンテストと呼ばれる、短時間でプログラムを記述する競技が世界中で開催されています。

参加者はアルゴリズム能力を問う問題を解決するプログラムを記述し、コンテストサイトに提出します。コンテストサイトはそのプログラムを採点し、この結果に応じて順位が決定します。



しかし、多くのコンテストはオンラインで開催されているため、競技中の不正行為が問題となっています。特に正答プログラムを複数人で共有する、ソースコードの盗作は検出することが非常に難しいです。これは、2つのソースコードの類似度を正確に測定する尺度は存在しないためです。

そこで、本研究ではこのような盗作を検出し、プログラミングコンテストでの公平性を確保することを目的とします。



また、本研究の成果は大学の宿題での盗作判定やコードクローンの検出にも活用できます。

提案手法

盗作判定の仕組みは4つのフェーズから成ります。

1. Fetch
2. Tokenize
3. Vectorization
4. Learning

1. 競技プログラミングサイトから C++ で書かれたソースコードを約 3000 個取得します。

2. ソースコードを整形しコメントを削除した後、スペースや括弧などで区切ります。そして各トークンに対して型、変数など14種類のラベルを設定します。



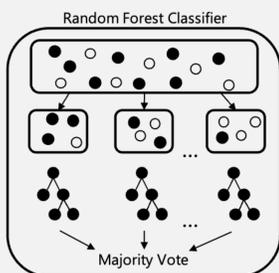
3. 作成したラベルの出現数を計算し、ベクトル化します。ベクトルにはソースコードの文字数とトークン数も加えます。



その後、各ソースコードから作成したベクトルを盗作ペア、非盗作ペアごとにベクトルを結合します。

型	変数	演算子	数詞	終止符	型	変数	演算子	数詞	終止符	関数...
1	1	1	1	1	1	1	1	1	1	0
3	2	1	0	2	3	2	1	0	2	2
...

4. 機械学習アルゴリズムである Random Forest Classifier を用いて盗作データセットを学習し、盗作判定モデルを作成します。



評価実験と結果

提案手法の有用性を評価することを目的として、先行研究である JPlag, Moss, PMD, SIM の4つの盗作判定ソフトウェア^{※1}と性能を比較しました。

実験は 280 個のデータセットに対し各ソフトウェアで盗作判定を行います。盗作データセットは以下の表に示す 14 種類の盗作パターンをもとに作成しました。

盗作パターン	内容
Visual code formatting	インデント、スペース、改行などの変更
Comments modification	コメントの変更、追加、削除
Translation of program parts	変数名、コメント、出力の自然言語を翻訳
Identifier rename	変数名、定数名、関数名の識別子の名前を変更
Changing constant values	定数、最終変数、enum値などの値を変更
Reordering independent lines of code	並び順に意味がないコード行の並び替え
Adding redundant lines of code	使われていないまたは何もしていないコードを追加
Splitting up lines of code	1行のコード、1つのクラスや関数を複数に分割
Merging lines of code	複数のコードやクラス、関数を一つに結合
Changing of statement specification	オペランドの順序を変更、論理演算を等価なものに変更
Replacing control structures with equivalents	制御構造体の置き換え
Simplifying the code	不要なコード行の削除
Changing the logic	大まかなロジックをコピー
Combining copied and original code	部分的なパーツをコピー

その結果、SA-Plagは各指標で高スコアを取得できました。

software	accuracy	recall	precision	F1
👑 SA-Plag	0.961	0.971	0.951	0.961
JPlag	0.918	0.864	0.968	0.913
Moss	0.950	0.957	0.944	0.950
PMD	0.732	0.500	0.933	0.651
SIM	0.832	0.829	0.835	0.832

※1 Martins, Vitor et al. (2014). Plagiarism detection: A tool survey and comparison. OpenAccess Series in Informatics.

一般公開

アーキテクチャ

本研究で作成したシステムは Web API として一般公開しました。ユーザは盗作ソースコードペアを Webサーバに送信すると、盗作可能性の有無とその確信度を取得できます。



Repository & Website



<https://github.com/xryuseix/SA-Plag>



<https://xryuseix.github.io/apps/sa-plag>

外部発表

情報通信システムセキュリティ研究会 (ICSS) にて論文発表。
「プログラミングコンテストにおけるソースコードの盗作検知手法の実装と評価」
石川 琉聖, 服部 祐一, 井上 博之, 猪俣 敦夫