

特別研究報告

題目

字句情報に基づく“自然さ”を利用した自動生成ファイルの特定

指導教員

楠本 真二 教授

報告者

土居 真之

平成 30 年 2 月 14 日

大阪大学 基礎工学部 情報科学科

字句情報に基づく“自然さ”を利用した自動生成ファイルの特定

土居 真之

内容梗概

近年、ソースコード解析に関する研究が盛んに行われている。ソースコードの解析において、解析対象のソースファイルの中には自動生成ファイルが含まれていることがある。しかし自動生成ファイルの存在により解析結果が目立たなくなることや解析時間が増加することがあるため、多くの場合自動生成ファイルは除外して解析する必要がある。自動生成ファイルを除外する方法として、ソースコードが自動生成ファイルであるかを目視で判定するという方法がある。しかしこの方法は時間的コストが大きくなってしまふ。他にも自動生成ファイル内に存在する特有のコメント文を文字列検索することにより特定するという方法があるが、この方法に関しても、自動生成ファイル特有のコメント文が消された場合に、自動生成ファイルを自動的に特定できないといった問題がある。

そこで本研究では、自動生成ファイルとしての“自然さ”と人が作成したファイルとしての“自然さ”を比較することで任意の自動生成ファイルを自動的に特定する手法を提案する。ファイルの自然さ、すなわち、自動生成あるいは人が生成したファイルとしてもっともらしい度合いは、確率的言語モデルである N-gram 言語モデルによって数値化を行う。この提案手法を評価するために、4つの自動生成プログラムから生成された自動生成ファイル群を対象に実験を行った。実験の結果、高い精度で自動生成ファイルを特定できた。また、機械学習に基づいた既存の手法と比較した結果、精度が向上していることを確認した。

主な用語

自動生成ファイル, N-gram 言語モデル, ソースコード解析, ソフトウェア保守

目次

1	まえがき	1
2	準備	3
2.1	自動生成ファイル	3
2.2	自動生成ファイル特定の既存手法	5
2.3	言語モデルと自然さ	7
3	提案手法	11
4	実装	12
4.1	字句解析	13
4.2	N-gram 言語モデルの構築と適用	14
5	評価実験	15
5.1	実験対象	15
5.2	評価尺度	15
5.3	実験 1	15
5.4	実験 2	17
5.5	実験 3	17
6	考察	19
6.1	実験 1 の考察	19
6.2	実験 2 の考察	19
6.3	実験 3 の考察	19
7	妥当性への脅威	20
8	あとがき	22
	謝辞	23
	参考文献	24

目次

1	自動生成ファイルの例	3
2	言語モデルの例	8
3	“A cat caught a mouse” におけるトリグラム言語モデルの適用	9
4	提案手法の概要	12
5	字句解析と自然さ計測の例	13
6	誤判定されたファイルの一部	20

表目次

1	自動生成プログラムの種類とファイル数	4
2	データセット作成時におけるノイズデータ含有率	5
3	自動生成プログラムのファイル名生成規則	6
4	ファイル名検索による自動生成ファイル特定の結果	7
5	交差検証による精度評価の結果	16
6	ブートストラップ法による精度評価の結果	16
7	MIX ファイル群の交差検証による精度評価の結果	17
8	別種類の自動生成ファイル言語モデルを適用時の交差検証による精度評価の結果	17

1 まえがき

近年，ソースコード解析に関する研究が盛んに行われている．例えば，ソースコード中に存在する互いに一致または類似したコード片であるコードクローンに関する研究や，ソフトウェア開発履歴データなどを対象に分析を行うリポジトリマイニングに関する研究などが行われている．このソースコード解析において，解析対象のソフトウェアに含まれるソースコードファイルの中には自動生成ファイルが含まれており，ソースコード解析を行う際に自動生成ファイルが弊害となることがある [15, 18]．具体的には，コードクローン検出において，自動生成ファイルから多くのコードクローンが検出されることにより他のコードクローンが目立たなくなってしまう [8, 20]．また，リポジトリマイニングにおいて，ソースコードを追跡する際に自動生成部分の追跡によって解析時間が増加してしまう [14]．したがって，ソースコード解析において自動生成ファイルは除外すべき存在である．

通常，自動生成ファイルにはそれ自身が自動生成ファイルであると明示するためのコメント文が残されている．ゆえに，そのようなソースコードに対しては `grep` コマンドを用いれば特定および除去することができる．しかしソースコードを修正していく過程でそのコメント文が消されてしまう場合がある．その場合 `grep` コマンドによる特定は難しく，またコメント文が消された自動生成ファイルを目視などで特定するのは時間的コストが大きい．したがって，そのようなコメント文が消された場合も含めて自動生成ファイルを自動的に特定することが必要となる．

コメント文が消された自動生成ファイルを自動的に特定するためには，自動生成ファイルにおける何らかの特徴を用いる必要がある．しかし任意の自動生成ファイルに共通する特徴を目視などで発見するのは困難である．そこで本研究では，N-gram 言語モデルを用いて，コメント文の有無にかかわらず自動生成ファイルか否かを自動的に判定する手法を提案する．

N-gram 言語モデルとは，既存のコードの字句や行などの並びからモデルを生成することにより，未知のファイルのモデルに対する自然さ，すなわちモデルらしさを数値として出力する統計的言語モデルである．提案手法では自動生成ファイルだと判明しているソースファイル，及び自動生成でないファイルと判明しているソースファイルをそれぞれ収集する．収集したソースファイルから，自動生成ファイルの言語モデルと自動生成でないファイルの言語モデルを生成する．それらを用いて未知のソースコードの自動生成ファイルとしての自然さと自動生成でないファイルとしての自然さをそれぞれ求め，これらを比較することによって自動生成ファイルかどうかを判定する．

提案手法の評価を行うために，4種類の自動生成ファイルを対象とする評価実験を行った．実験の結果，高い精度で自動生成ファイルを特定できることを確認した．

以降，2章では，準備として自動生成ファイルおよび N-gram 言語モデルについて詳細に述べる．3章では，提案手法について説明する．4章では提案手法の評価実験を行うための実装について説明し，5

章では評価実験の結果を述べる。6章では、実験結果に対し考察を行い、7章では、妥当性の脅威について述べる。最後に、8章で、本研究のまとめおよび今後の課題について述べる。

```

/* Generated By:JavaCC: Do not edit this line. Func2Parser.java */

    . . .

    if ((active0 & 0x400L) != 0L)
        return jjStartNfaWithStates_0(2, 10, 10);
    else if ((active0 & 0x4000L) != 0L)
        return jjStartNfaWithStates_0(2, 14, 10);
    else if ((active0 & 0x40000L) != 0L)
        return jjStartNfaWithStates_0(2, 18, 10);
    else if ((active0 & 0x400000L) != 0L)
        return jjStartNfaWithStates_0(2, 22, 10);
    else if ((active0 & 0x8000000L) != 0L)
        return jjStartNfaWithStates_0(2, 23, 10);
    else if ((active0 & 0x10000000L) != 0L)
        return jjStartNfaWithStates_0(2, 28, 10);
    else if ...

```

図1 自動生成ファイルの例

2 準備

2.1 自動生成ファイル

本節では、自動生成ファイルの定義、収集及び特定方法について述べる。自動生成ファイルとは、プログラムによって自動的に生成されたソースファイルのことである。ソースファイルを自動で生成するプログラムは多数存在し、例として GUI プログラムを生成する GUI ビルダー、あるプログラミング言語を別の言語に変換するトランスレータなどがある。本研究では比較的収集が容易であるパーサジェネレータによって生成されたファイルを対象とする。

また研究の準備段階において、パーサジェネレータ生成ファイルに対して目視による調査を行ったところ、ソースコードにおいて以下のような特徴があることが分かった。

- 変数宣言文が連続して出現することが多い
- 似た条件式が連続して出現することが多い

自動生成ファイルのソースコードの一部を図1に示す。このように、人が書いたソースコードとプログラムによって自動的に生成されたソースコードには視覚的な差異が存在するため、目視によって自動生成ファイルを特定することは可能である。しかし、目視によって自動生成ファイルを特定すると時間

的成本が大きくなってしまふ。

通常、自動生成ファイルにはそれ自身が自動生成ファイルであると明示するためのコメント文が記述されている。自動生成ファイルのコメント文の例を図1の上部に示す。このようなコメント文を文字列検索することにより、自動生成ファイルを自動的に特定することができる。しかし、機能追加やバグ修正などの過程においてコメント文が削除されてしまう場合があるため、このコメント文検索だけでは特定できない自動生成ファイルが存在する。また、コメント文が残っていた場合でも `grep` コマンドの引数として適切な文字列を与えなければ特定することはできない。

2.1.1 自動生成ファイルと自動生成でないファイルの収集

自動生成ファイルを特定するためのデータセットとして、自動生成ファイル群と自動生成でないファイル群を用いる。そこで本研究では、下仲らが作成したデータセットを用いた [19]。その収集方法について述べる。

まず自動生成ファイルの収集方法について説明する。上述したとおり、自動生成ファイルにはそれ自身が自動生成ファイルであると明示するためのコメント文が記述されているため、そのコメント文の一部を検索キーワードとして用いて自動生成ファイルを収集した。下仲らは、GitHub [7] からコメント文検索により収集した [19]。しかし GitHub 上の検索では 1,001 件以上の検索結果を取得することができない。この問題に対して高澤らは、検索結果数が 1,000 件を超えないようなファイルサイズを指定し、ファイルサイズごとに分割して収集することで対応している [21]。下仲らも同様の方法を用いた。また収集の際には、jsoup [11] を用いたウェブスクレイピングを行い、自動で収集を行った。収集した自動生成ファイルは、4つの自動生成プログラムから生成されたものである。4つの自動生成プログラム名と、収集した自動生成ファイル数を表1に示す。

これらの自動生成プログラムはパーサジェネレータであり、Java で記述されたソースファイルを生成する。以降、ANTLR によって生成されたファイル、JavaCC によって生成されたファイル、JFlex によって生成されたファイル、SableCC によって生成されたファイルをそれぞれ **ANTLR 生成ファイル**、**JavaCC 生成ファイル**、**JFlex 生成ファイル**、**SableCC 生成ファイル**と呼ぶ。

次に、自動生成でないファイルの収集方法について説明する。下仲らは、大規模なソースファイルの集合である Apache リポジトリ [4] からソースファイルをランダムに収集した。また、Apache リポジトリの中に自動生成ファイルが含まれている可能性があるため、それらをコメント文検索により特定

表1 自動生成プログラムの種類とファイル数

自動生成プログラム	ANTLR	JavaCC	JFlex	SableCC
ファイル数	9,218	15,033	3,737	16,603

し、除外した。

上記の方法で収集した自動生成ファイル群および自動生成でないファイル群の中には、誤って自動生成ファイルとみなしているもの、もしくは誤って自動生成でないファイルとみなしているもの（これらをノイズデータと呼ぶ）が存在している可能性がある。そこで、ノイズデータを除去するために目視確認を行った。しかし、すべてのソースファイルを目視確認するのは時間的コストが膨大であるため、4種類の自動生成ファイル群、および自動生成でないファイル群に対し、以下の処理を行った。

1. 各 1,000 ファイルずつランダムに抽出する
2. 目視確認によりノイズデータを除去する
3. 除去したソースファイルの数だけ、再度ランダムに抽出する
4. 3. で抽出したソースファイルに対して目視確認によりノイズデータを除去する

上記の 3. および 4. を繰り返し行い、ANTLR 生成ファイル、JavaCC 生成ファイル、JFlex 生成ファイル、SableCC 生成ファイル、自動生成でないファイルがそれぞれ 1,000 ファイルずつ存在するデータセットを作成した。この 1,000 ファイルを収集する過程でファイル群に含まれていたノイズデータの割合を表 2 に示す。

コメント文検索により自動生成でないファイルが自動生成ファイルと判定されていた要因としては、キーワードとして用いたコメント文が文字列リテラルとしてソースコード中に存在していたことがあげられる。

2.2 自動生成ファイル特定の既存手法

本節では、2.1 で述べた自動生成ファイルを特定する既存手法について述べる。既存手法としてはコメント文検索、ファイル名検索、機械学習による検索が挙げられる。

表 2 データセット作成時におけるノイズデータ含有率

ファイル群	チェックした ファイル数	ノイズデータ 含有率
ANTLR 生成ファイル群	1,000	0.0%
JavaCC 生成ファイル群	1,006	0.6%
JFlex 生成ファイル群	1,010	1.0%
SableCC 生成ファイル群	1,004	0.4%
自動生成でないファイル群	1,032	3.0%

2.2.1 ファイル名検索による自動生成ファイルの特定

コメント文検索以外の自動生成ファイル特定手法として、ファイル名による検索がある。通常、自動生成ファイルのファイル名には、自動生成プログラムごとに定められている生成規則がある。例えば SableCC では、以下のようなものが定められている [3]。

- AXxx.java
- TXxx.java

ただし、X は大文字の任意のアルファベット、x は小文字の任意のアルファベットを表す。2.1.1 で作成したデータセットを用いて、ファイル名検索による自動生成ファイルの特定を行った。4 つの自動生成プログラムごとのファイル名生成規則を表 3 に示す。

表中の *ClassName* は Java のクラス名を表す。正規表現を用いて、これらの生成規則にマッチするものを自動生成ファイルとして特定する。また、これらの生成規則にマッチしないものを自動生成でないファイルとみなす。

ファイル名検索による自動生成ファイル特定の結果を表 4 示す。ANTLR 生成ファイル群、JavaCC 生成ファイル群、SableCC 生成ファイル群においては 1,000 ファイル中約 800 ファイルの自動生成ファイルをそれぞれ特定できているのに対し、JFlex 生成ファイル群は 79 ファイルと極端に少なくなって

表 3 自動生成プログラムのファイル名生成規則

自動生成プログラム	生成規則
ANTLR	[<i>ClassName</i>]Lexer.java, [<i>ClassName</i>]Parser.java, [<i>ClassName</i>]Listener.java, [<i>ClassName</i>]BaseListener.java
JavaCC	JJT[<i>ClassName</i>]State.java, [<i>ClassName</i>]Constants.java, Node.java, ParseException.java, SimpleCharStream.java, SimpleNode.java, Token.java, TokenMgrError.java, ASTPerl.java, ASTPython.java, [<i>ClassName</i>]TreeConstants.java, [<i>ClassName</i>]Visitor.java
JFlex	Yylex.java, [<i>ClassName</i>].java
SableCC	Lexer.java, LexerException.java, Parser.java, ParserException.java, DepthFirstAdapter.java, Analysis.java, Switch.java, Switchable.java, TXxx.java, Token.java, AXxx.java, Start.java, Node.java, State.java

いることが分かる。これは、JFlex のファイル名生成規則の数が少ないことが要因と考えられる。このことから、ファイル名生成規則だけでは自動生成ファイルを特定するのに十分ではないことが分かる。

2.2.2 機械学習を利用した自動生成ファイルの特定

コメント文検索とファイル名による検索以外の自動生成ファイル特定手法として、機械学習を用いて特定する手法がある [19]。これは自動生成ファイルと自動生成でないファイルの構文情報の出現回数を学習データを用いてモデルを構築し、このモデルによって自動生成ファイルか否かを判定する手法である。この手法を 2.1.1 で述べたデータセットに対して 4 種類の機械学習アルゴリズム (Decision Tree, Naive Bayes, Random Forest, SVM) を用いた場合の適合率の上限と下限は 99.9%, 82.7%, 再現率の上限と下限は 99.9%, 78.4% である。またサイズの小さいファイルに対して適用した場合、精度が悪くなるといった課題がある。したがって機械学習を利用した手法は自動生成ファイルを特定するのに十分ではないことが分かる。

2.2.3 既存手法における課題

コメント文検索では、コメント文が削除されている場合うまく機能せず、ファイル名検索では、十分な量のファイルを特定できないといった問題が存在する。また機械学習を利用した手法ではサイズの小さいファイルに適用した場合、精度が悪くなるという問題が存在する。一方で、人が書いたソースコードとプログラムによって自動的に生成されたソースコードには字句の並びに視覚的な差異が存在するため、目視では容易に特定することが可能である。このソースコードの視覚的な差異は、そこに存在している構文の違いによるものであるため、字句の並びにより自然さを求める N-gram 言語モデル [2] を利用することにより、高い精度で自動生成されたソースコードを特定できると考えた。

2.3 言語モデルと自然さ

本節では、言語モデルと言語モデルによって求められる自然さについて述べる。まず言語モデル $P(W)$ とは、文字列 W の言語らしさを確率として与えるモデルであり、主に自然言語解析に用いられ

表 4 ファイル名検索による自動生成ファイル特定の結果

自動生成プログラム	特定した自動生成ファイル数	特定した自動生成でないファイル数
ANTLR	894	986
JavaCC	767	989
JFlex	79	1,000
SableCC	867	965

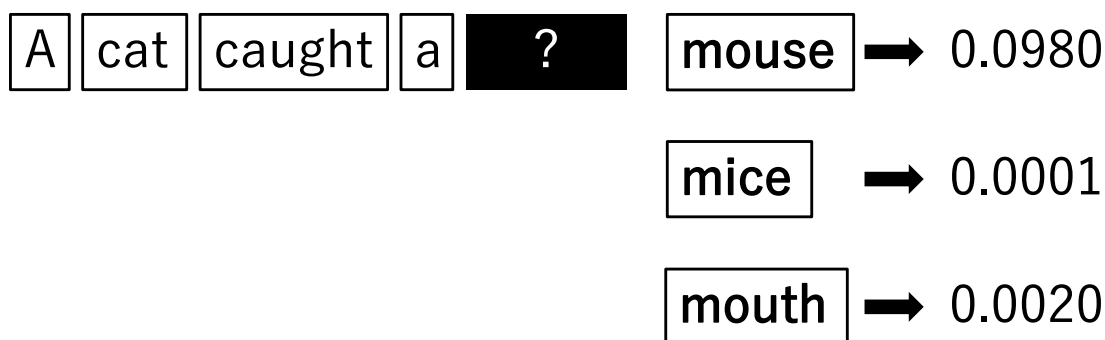


図2 言語モデルの例

ている [17]. この言語モデルによって求まる言語らしさを自然さと呼ぶ. すなわち文字列 W が言語としてもっともらしければ自然さは高い値をとり, 言語もっともらしくなければ自然さは低い値をとる. 例として “A cat caught a mouse” という文における自然さを図2に示す. “A cat caught a” という単語列に対し, 次に mouse がくると自然な文になるため高い自然さを得られる. 一方で mice のように文法に合わない単語や, mouth のように意味が不自然な文だと自然さは mouse の場合と比べて自然さは小さくなる.

2.3.1 N-gram 言語モデル

言語モデルの代表例として, N-gram 言語モデルがある. N-gram 言語モデルは, 文を単語に分解し, その単語の並びに対して自然さを計算する言語モデルである. なお N-gram 言語モデルの N はオーダーと呼び, 自然さを計算する際にどれだけの単語に対して確率を定義するかの長さを表す. すなわち N-gram 言語モデルでは, ある時点での単語の生起確率はその直前の $N-1$ 個の単語にのみ依存すると仮定し, 直前の $N-1$ 個の単語列から次の単語への遷移確率の計算を行うのが N-gram 言語モデルである. 自然さは式 1,2 で計算される.

なお式 1,2 における W は入力となるテキストデータであり, 出力として自然さ $P(W)$ を得る. また W を単語に分解した場合の i 番目の単語が w_i であり, $c(w)$ はその単語列の出現回数である.

$$P(W) = \prod_{i=1}^{|W|+1} p(w_i | w_0 \dots w_{i-1}) \quad (1)$$

$$p(w_i | w_0 \dots w_{i-1}) = \frac{c(w_{i-n+1} \dots w_i)}{c(w_{i-n+1} \dots w_{i-1})} \quad (2)$$

特にオーダーが $N=1,2,3$ のときの N-gram 言語モデルはそれぞれユニグラム言語モデル, バイグラ

A cat caught a mouse

$$\begin{aligned}
 \text{自然さ} &= p(A) \\
 &\times p(\text{cat} \mid A) \\
 &\times p(\text{caught} \mid A \text{ cat}) \\
 &\times p(a \mid \text{cat caught}) \\
 &\times p(\text{mouse} \mid \text{caught a})
 \end{aligned}$$

図3 “A cat caught a mouse”におけるトリグラム言語モデルの適用

ム言語モデル，トリグラム言語モデルという．例としてトリグラム言語モデルによる“A cat caught a mouse”という文の自然さの計算過程を図3に示す．

2.3.2 ゼロ頻度問題とスムージング

N-gram 言語モデルでは，文に対する確率を単語単位に分解してその積を計算するが，そのうちどこか1カ所にでも学習用のテキストにおいて出現しないような単語の並びがあると，文全体に対する確率が0となる．このように言語モデルにおいて確率が0となる問題をゼロ頻度問題と言う．N-gram 言語モデルでは，オーダーを大きくするとゼロ頻度問題が発生しやすくなり，逆にオーダーを小さくすると精度が落ちることが知られている [2]．このゼロ頻度問題を避けるために言語モデルにはスムージングを行う．

スムージングの手法には加算スムージング，PPM など様々な手法があるが，最も良い性能をもつスムージング手法の一つとして Kneser-Ney スムージング [12] がある．Kneser-Ney スムージングでは高次の N-gram の確率を計算する際にそれよりも低次の N-gram の確率を用いることによってスムージングを行う．

2.3.3 言語モデルのプログラム解析への応用

この N-gram 言語モデルによって求められる自然さは近年プログラム解析の分野にも応用されている．実際にプログラムの字句を単語に見立て，プログラム全体を文とみなしてモデルを生成することで，コード補完 [10] やバグ検出 [16] が行えることが知られている．

本研究ではこの N-gram 言語モデルを用いて自動生成ファイルと自動生成でないファイルの判定を自動的に行う.

3 提案手法

本研究では、N-gram 言語モデルを利用し、自動生成ファイルを自動的に特定する手法を提案する。提案手法の概要を図 4 に示す。本研究における提案手法の入力は、学習データ、すなわち自動生成ファイル群及び自動生成でないファイル群と、テストデータ、すなわち自動生成ファイルか判定したいファイルである。また出力はテストデータが自動生成ファイルか否かの判定結果である。

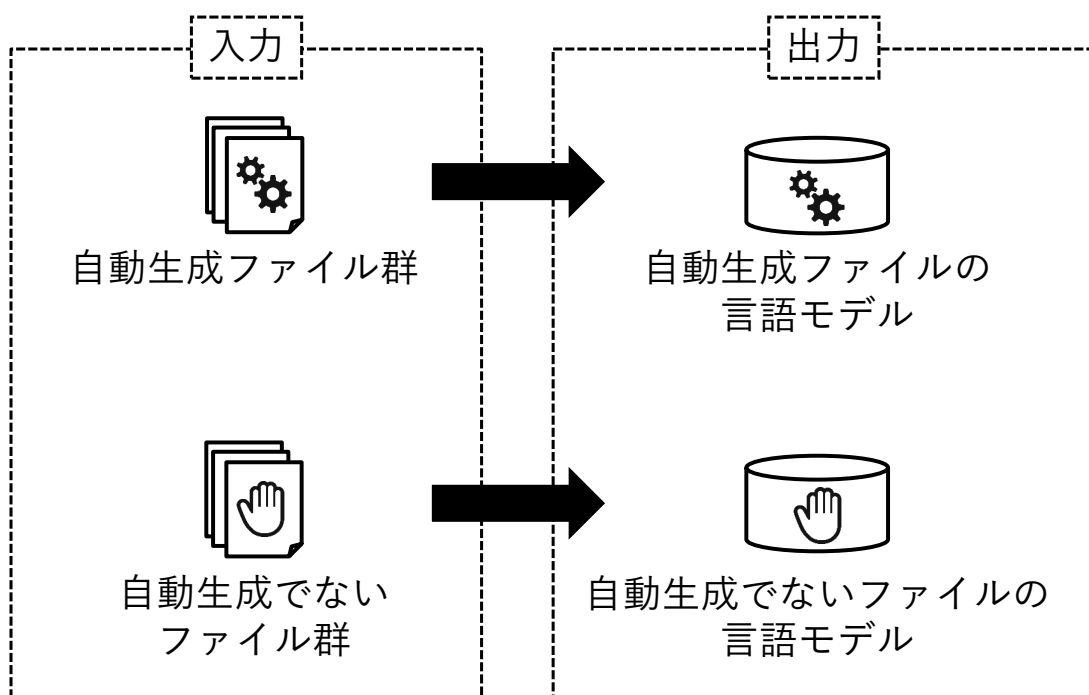
提案手法は次の 2 ステップから構成される。Step1 では、学習データから自動生成ファイルと自動生成でないファイルそれぞれの言語モデルを構築する。Step2 では、Step1 で構築した各言語モデルを用いてテストデータの自然さを計測し、それぞれの自然さを比較する。以降、各ステップについて詳細に説明する。

Step1: 言語モデルの構築

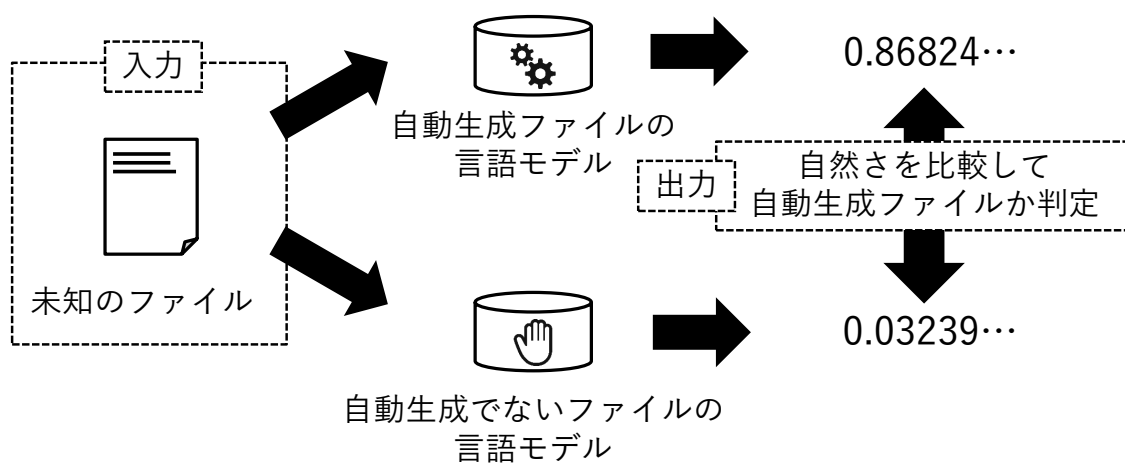
Step1 では、学習データから字句の並びを取得し、言語モデルを構築する。Step1 における入力は学習データであり、出力は学習データの各ファイル群より生成された自動生成ファイルの言語モデルと自動生成でないファイルの言語モデルである。まず学習データの各ファイル群に対して字句解析を行い字句の並びを取得する。この字句の並びから N-gram 言語モデルを構築することによって自動生成ファイルの言語モデルと自動生成でないファイルの言語モデルを得る。

Step2: 言語モデルの適用

Step2 では、自動生成ファイルであるか判定したいソースファイルに対して Step1 で作成した 2 つの言語モデルを適用する。Step2 における入力はテストデータであり、出力はテストデータが自動生成ファイルか否かの判定結果である。まずテストデータを Step1 と同様に字句解析を行い、その結果を言語モデルに適用する。その結果得られた自然さを比較することで自動生成ファイルか否かの判定を行う。すなわち、自動生成ファイルとしての自然さの方が高ければ自動生成ファイル、そうでないならば自動生成でないと判定する。



(a) Step 1



(b) Step 2

図4 提案手法の概要

4 実装

本章では提案手法の実装方法について述べる。

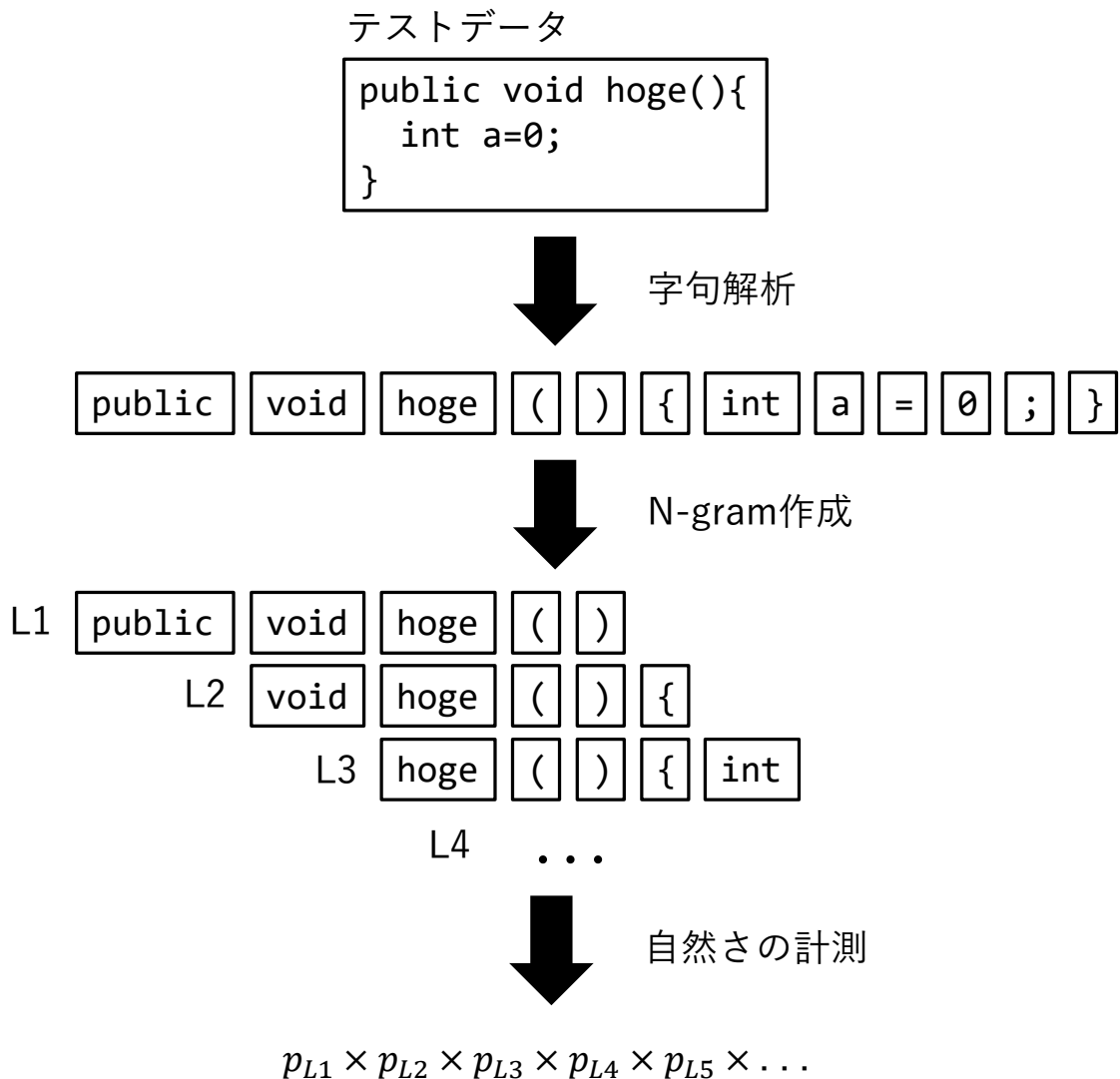


図5 字句解析と自然さ計測の例

4.1 字句解析

ソースファイルの字句解析には Eclipse JDT 4.5.2 [5] を用いて AST 解析を行った結果を用いる。Eclipse JDT 4.5.2 では 92 種類の AST ノードが定義されているが、そのうち 3 つはコメント文に関するノードであるので、本研究ではコメント文に関するノードを除いた 89 種類のノードからなる AST を構成する字句列を抽出して、その字句の並びから自然さの計測を行う。字句解析及び自然さの計測の例を図 5 に示す。

4.2 N-gram 言語モデルの構築と適用

N-gram 言語モデルの構築と適用には KenLM[9] を用いる。KenLM におけるスムージングの手法は 2.3.2 で述べた Kneser-Ney スムージングを用いている。

5 評価実験

本章では、提案手法を評価するために行った実験と、その実験結果について述べる。

提案手法の評価を行うために、3つの実験を行った。その概要を以下に示す。

実験 1 ある単一種類の自動生成ファイル群で構築した言語モデルを用いて提案手法の精度を評価した。

実験 2 複数種類の自動生成ファイルを用いた自動生成ファイル群で構築した言語モデルを用いた場合の精度を評価した。

実験 3 ある自動生成プログラムの自動生成ファイル群から構築した言語モデルを用いて別種類の自動生成ファイル判定した場合の精度を評価した。

5.1 実験対象

学習データとして、2.1.1 で述べたデータセットを用いた。加えて、4種類の自動生成ファイル群計4,000ファイルの中から1,000ファイルをランダムで抽出した。このソースファイル群を **MIX** ファイル群と定義する。

5.2 評価尺度

実験の評価尺度として、適合率と再現率を用いる。以下、それぞれの尺度の定義について説明する。

適合率 自然さを比較して自動生成ファイルと判定されたファイルのうち、実際に自動生成ファイルであるものの割合

再現率 実際に自動生成ファイルであるもののうち、自然さを比較して自動生成ファイルと判定されたものの割合

5.3 実験 1

まず収集した各自動生成ファイル群から言語モデルを構築する。同様に自動生成でないファイル群から言語モデルを構築し、計5種類の言語モデルを得た。なお本実験では N-gram 言語モデルのオーダーは5として構築した。構築された計5種類の言語モデルの性能を評価するために交差検証とブートストラップ法による検証を行った。

交差検証では、まずデータセットを N 個のブロックにランダムに分割する。分割したブロックのうち、 $N - 1$ 個のブロックを学習データとし、残りの1個のブロックをテストデータとして評価を行う。この処理を、すべてのブロックが1度テストデータとなるようブロックを変化させながら N 回行い、それらの精度の平均をとる。なお本実験におけるブロック数 N は、最も実際の精度に近いとされてい

る $N = 10$ として言語モデルの評価実験を行った [13]. またその結果を下仲らによる機械学習の 4 種類
 のアルゴリズム (Decision Tree, Naive Bayes, Random Forest, SVM) を用いた手法と比較した [19].
 その結果を表 5 に示す.

提案手法では全ての場合において適合率, 再現率ともに 99% を超える精度になっていることが分かる.
 また機械学習による判定結果と比較すると, すべての場合において精度が同じもしくは提案手法の
 方が精度が高くなっていることが分かる.

交差検証は実際の精度に近い精度をだすとされているが, データセットからブロックを作成する際
 に, データの偏りが生じて実験結果の分散が大きくなるといった問題がある [1]. その問題を回避する
 ためブートストラップ法による検証も行った. ブートストラップ法では, まずデータセットからデータ
 セットと同じ数だけのファイルをランダムに抽出し, 抽出したファイルで新しいデータセットを作成す
 ることを繰り返して複数のデータセットを得る. この複数のデータセットにおける精度の平均を計算す
 ることで, 元のデータセットの評価を行う. ブートストラップ法における必要な繰り返し回数は 50 回
 から 2,000 回で十分とされているため, 本実験では 2,000 個のデータセットを作成し評価を行った [6].
 その結果を表 6 に示す.

提案手法では適合率, 再現率ともに 99% を超える精度になっていることが分かる. また機械学習を
 用いた手法による判定結果と比較すると, すべての場合において提案手法の方が精度が高くなっている
 ことが分かる.

表 5 交差検証による精度評価の結果

自動生成プログラム		ANTLR		JavaCC		JFlex		SableCC	
		適合率	再現率	適合率	再現率	適合率	再現率	適合率	再現率
提案手法		100.0%	100.0%	99.6%	99.3%	99.8%	99.7%	100.0%	99.9%
下仲らの 手法 [19]	Decision Tree	99.9%	99.9%	97.0%	97.0%	99.4%	99.4%	96.3%	96.2%
	Naive Bayes	98.8%	98.8%	88.0%	85.7%	99.5%	99.5%	82.7%	78.4%
	Random Forest	99.9%	99.9%	97.3%	97.3%	99.7%	99.7%	96.1%	96.1%
	SVM	99.8%	99.8%	95.7%	95.7%	99.6%	99.6%	86.8%	84.5%

表 6 ブートストラップ法による精度評価の結果

自動生成プログラム		ANTLR		JavaCC		JFlex		SableCC	
		適合率	再現率	適合率	再現率	適合率	再現率	適合率	再現率
提案手法		99.99%	99.98%	99.86%	99.75%	99.97%	99.86%	100.0%	99.92%

5.4 実験 2

実験 1 では、単一種類の自動生成ファイルで構築した言語モデルの性能を評価した。実験 2 では、MIX ファイル群によって自動生成ファイルを特定できるかを確認するため交差検証を行った。また実験 1 同様、下仲らの手法との比較を行った。この結果を表 7 に示す。

提案手法では適合率、再現率ともに 98% を超える精度になっていることが分かる。また機械学習を用いた手法による判定結果と比較すると、実験 1 同様にすべての場合において提案手法の方が精度が高くなっている。

5.5 実験 3

実験 1 では同一種類の自動生成ファイルに対して交差検証を行っているが、実験 3 では、ある自動生成ファイル群を用いて構築した言語モデルで、別種類の自動生成ファイルを特定できるかどうかを確認した。そのため、言語モデルを適用したのは自動生成ファイルのみであり、自動生成でないファイルには適用していない。ある自動生成ファイルから構築した言語モデルを別種類の自動生成ファイルへ適用した結果を表 8 に示す。表中の数値は *Recall* を表している。表 8 から全ての場合において別種類の

表 7 MIX ファイル群の交差検証による精度評価の結果

自動生成プログラム		MIX	
		適合率	再現率
提案手法		98.7%	99.7%
下仲らの手法 [19]	Decision Tree	95.3%	95.3%
	Naive Bayes	85.3%	80.6%
	Random Forest	96.8%	96.8%
	SVM	85.2%	79.1%

表 8 別種類の自動生成ファイル言語モデルを適用時の交差検証による精度評価の結果

自動生成プログラム	ANTLR	JavaCC	JFlex	SableCC
ANTLR	-	8.7%	0.0%	0.0%
JavaCC	17.1%	-	13.2%	7.7%
JFlex	0.7%	38.8%	-	0.0%
SableCC	0.0%	4.8%	0.0%	-

自動生成ファイルで構築した言語モデルで自動生成ファイルの検出は難しいことが分かる。このことから、各自動生成ファイルの字句の並びは大きく異なると考えられる。

6 考察

本章では、5で述べた評価実験、および提案手法における有用性について考察を行う。

6.1 実験1の考察

交差検証の結果において、誤判定されたファイルを調べた。誤判定されたファイルの一部を図6に示す。その結果、自動生成ファイルと誤判定された自動生成でないファイルの特徴として、以下のことがあげられる。

- case文が多い
- シグネチャが類似したメソッドが多数ある

これらは、2章で述べたように自動生成ファイルの特徴であるので、それらが誤検出の要因であると考えられる。また自動生成でないファイルと誤判定された自動生成ファイルに関しては、図6(b)で示しているような目視によっても自動生成か判断のつきにくいファイル、すなわち自動生成でないようなファイルであった。

また、下仲らの手法ではサイズが小さいファイルに対して適用した際に誤検出されてしまう傾向があったが、提案手法においてはその傾向はあまりみられなかった。そのため、提案手法はファイルサイズが小さい場合においても有用であると考えられる。

6.2 実験2の考察

交差検証の結果において、誤判定されたファイルを調べたところ、実験1で誤判定されたファイルが多く含まれていた。したがって自動生成ファイルの特徴をもつ自動生成でないファイルは誤検出される傾向があると考えられる。また実験1と比較して適合率の値が低下した。これは複数種類の自動生成ファイル群を混ぜたことによって、自動生成ファイルの種類ごとに存在する特有の特徴が目立たなくなったためと考えられる。

6.3 実験3の考察

交差検証の結果において、ある自動生成ファイル群を用いて構築した学習モデルを別種類の自動生成ファイル群に適用した場合、全く特定することができなかったものもある。したがって、別種類の自動生成ファイル間における学習モデルの適用は難しく、実験2の結果より学習モデルに対象の自動生成ファイルを含める必要性があると考えられる。


```

. . .
case '%b':
    sb.append("%%b");
    break;
case '%f':
    sb.append("%%f");
    break;
case '%n':
    sb.append("%%n");
    break;
case '%r':
    sb.append("%%r");
    break;
case '%t':
    sb.append("%%t");
    break;
. . .

```

(a) 自動生成ファイルと誤判定された自動生成でないファイル

```

. . .
private void initTermBuffer() {
    if (termBuffer == null) {
        if (termText == null) {
            termBuffer = new char[MIN_BUFFER_SIZE];
            termLength = 0;
        } else {
            int length = termText.length();
            if (length < MIN_BUFFER_SIZE) length = MIN_BUFFER_SIZE;
            termBuffer = new char[length];
            termLength = termText.length();
            termText.getChars(0, termText.length(), termBuffer, 0);
            termText = null;
        }
    } else if (termText != null)
        termText = null;
}
. . .

```

(b) 自動生成でないファイルと誤判定された自動生成ファイル

図6 誤判定されたファイルの一部

7 妥当性への脅威

本章では、評価実験に含まれる妥当性への脅威について述べる。

本研究では、自動生成でないコードを Apache リポジトリから収集した。その際、自動生成でないことを目視確認をしているが、この目視確認が間違っている可能性がある。そのため、本来自動生成ファイルであるものを自動生成でないファイルとみなしている可能性がある。

本実験で対象とした自動生成ファイルは、Java で記述された 4 種類の自動生成ファイルである。そのため、他の種類の自動生成ファイルを用いた場合や、他の言語で記述された自動生成ファイルを用いた場合には、本実験とは異なる結果が得られる可能性がある。

8 あとがき

本研究では、N-gram 言語モデルを用いて自動生成ファイルを自動的に特定する手法を提案した。提案手法では、まず自動生成ファイル特有のコメント文の有無にかかわらず、自動生成ファイルか否かを判定するために、自動生成ファイルと自動生成でないファイルそれぞれのソースコードの字句の並びを抽出する。次にそれらを学習データとして自動生成ファイルの言語モデルと自動生成でないファイルの N-gram 言語モデルを構築し、この言語モデルを自動生成か判定したいファイルの自然さを計測することで自動生成ファイルか判定する。

4 種類の自動生成ファイルを収集し、それらを対象に評価実験を行った。その結果、全ての場合で適合率、再現率ともに 99% 以上と、高い精度で自動生成ファイルを特定できていることを確認した。また 4 種類の自動生成ファイルを混ぜた場合の評価実験も行った。その結果に関しても、適合率、再現率ともに 98% 以上と、高い精度で自動生成ファイルを特定できていることを確認した。加えて他種類の自動生成ファイルを特定できるかの評価実験を行った。しかしながら他種類の自動生成ファイルの特定は行えないことを確認した。

今後は任意の種類、および任意の言語の自動生成ファイルにおいて、高精度で特定できる手法とするため改善していく予定がある。

謝辞

本研究を行うにあたり，理解あるご指導を賜り，暖かく励まして頂きました楠本真二教授に心より感謝申し上げます。

本研究の全過程を通し，終始熱心かつ丁寧なご指導を頂きました，肥後芳樹准教授に深く感謝申し上げます。

本研究に関して，有益かつ的確なご助言を頂きました，松本真佑助教に深く感謝申し上げます。

本研究を進めるにあたり，適切なお助言および多大なるご助力を頂きました大阪大学大学院情報科学研究科コンピュータサイエンス専攻博士前期課程 2 年の下仲健斗氏に心より感謝申し上げます。

本研究を進めるにあたり，日常の中で相談に乗って声をかけて頂き，研究の基礎を一から指導して頂きました大阪大学大学院情報科学研究科コンピュータサイエンス専攻博士前期課程 1 年の有馬諒氏に心より感謝申し上げます。

本研究を進めるにあたり，様々な形で励まし，協力を頂きましたその他の楠本研究室の皆様のご協力に心より感謝申し上げます。

最後に，本研究に至るまでに，講義，演習，実験等でお世話になりました大阪大学基礎工学部情報科学学科の諸先生方に，この場を借りて心から御礼申し上げます。

参考文献

- [1] Yoshua Bengio and Yves Grandvalet. No unbiased estimator of the variance of k-fold cross-validation. *J. Mach. Learn. Res.*, Vol. 5, pp. 1089–1105, December 2004.
- [2] Peter F. Brown, Peter V. deSouza, Robert L. Mercer, Vincent J. Della Pietra, and Jenifer C. Lai. Class-based n-gram models of natural language. *Comput. Linguist.*, Vol. 18, No. 4, pp. 467–479, December 1992.
- [3] Pinaki Chakraborty. Object-oriented compilers: A review. *IUP Journal of Information Technology*, Vol. 13, No. 1, p. 36, 2017.
- [4] Apache Source code repository. <http://svn.apache.org/repos/asf/>.
- [5] Eclipse Java development tools (JDT). <http://www.eclipse.org/jdt/>.
- [6] Bradley Efron. Bootstrap methods: another look at the jackknife. In *Breakthroughs in statistics*, pp. 569–593. Springer, 1992.
- [7] GitHub. <http://github.com/>.
- [8] Nils Göde and Rainer Koschke. Frequency and risks of changes to clones. In *Proceedings of the 33rd International Conference on Software Engineering, ICSE '11*, pp. 311–320, New York, NY, USA, 2011. ACM.
- [9] Kenneth Heafield. Kenlm: Faster and smaller language model queries. In *Proceedings of the Sixth Workshop on Statistical Machine Translation, WMT '11*, pp. 187–197, Stroudsburg, PA, USA, 2011. Association for Computational Linguistics.
- [10] Abram Hindle, Earl T. Barr, Mark Gabel, Zhendong Su, and Premkumar Devanbu. On the naturalness of software. Vol. 59, pp. 122–131, New York, NY, USA, April 2016. ACM.
- [11] jsoup: Java HTML Parser. <http://jsoup.org/>.
- [12] Reinhard Kneser and Hermann Ney. Improved backing-off for m-gram language modeling. In *In Proceedings of the IEEE International Conference on Acoustics, Speech and Signal Processing*, Vol. I, pp. 181–184, May 1995.
- [13] Ron Kohavi. A study of cross-validation and bootstrap for accuracy estimation and model selection. In *Proceedings of the 14th International Joint Conference on Artificial Intelligence - Volume 2, IJCAI'95*, pp. 1137–1143, San Francisco, CA, USA, 1995. Morgan Kaufmann Publishers Inc.
- [14] Alexander C MacLean, Landon J Pratt, Jonathan L Krein, and Charles D Knutson. Trends that affect temporal analysis using sourceforge data. *Proceedings of the 5th International*

- Workshop on Public Data about Software Development (WoPDaSD' 10), p. 6, 2010.
- [15] Pam McDonald, Dan Strickland, and Charles Wildman. Estimating the effective size of autogenerated code in a large software project. In *Proceedings of the 17th International Forum on COCOMO and Software Cost Modeling*. Citeseer, 2002.
- [16] Baishakhi Ray, Vincent Hellendoorn, Saheel Godhane, Zhaopeng Tu, Alberto Bacchelli, and Premkumar Devanbu. On the "naturalness" of buggy code. pp. 428–439, 2016.
- [17] Andreas Stolcke. Srilm – an extensible language modeling toolkit. In *ICSLP*, pp. 901–904, 2002.
- [18] Shinji Uchida, Akito Monden, Naoki Ohsugi, Toshihiro Kamiya, Ken-Ichi Matsumoto, and Hideo Kudo. Software analysis by code clones in open source software. *Journal of Computer Information Systems*, Vol. 45, No. 3, pp. 1–11, 2005.
- [19] 下仲健斗, 鷺見創一, 肥後芳樹, 楠本真二. 機械学習を利用した構文情報に基づく自動生成ファイルの特定. *情報処理学会論文誌*, Vol. 58, No. 4, pp. 1–10, 4 2017.
- [20] 大田崇史, 井垣宏, 堀田圭祐, 肥後芳樹, 楠本真二ほか. ソフトウェア開発におけるコピーアンドペーストによって生じたコード片に対する調査. *研究報告ソフトウェア工学 (SE)*, Vol. 2014, No. 22, pp. 1–6, 2014.
- [21] 高澤亮平, 坂本一憲, 鷺崎弘宜, 深澤良彰. Repositoryprobe: リポジトリマイニングのためのデータセット作成支援ツール. *コンピュータ ソフトウェア*, Vol. 32, No. 4, pp. 103–114, 2015.