

ソースコード内の文字出現率を利用したコードクローン判別法

垣谷 広輝† 安藤 優作† 増田 智樹† 平山 雅之† 菊地 奈穂美‡

日本大学†

沖電気工業株式会社‡

1. はじめに

ソフトウェアの保守性を阻害する要因としてソースコード内のコードクローン（以降、クローンと表記）の存在が問題視されている。クローンは、コード内の処理をコピー&ペーストすることで生まれる類似コード断片（コードブロック）である。保守性を向上させるためにはクローンを早期に検出し、内容確認や場合によっては除去するなどの対策を行う必要がある。

従来、クローンは主に構文解析を利用する方法によって検出されてきた。この方法では文法的に誤りのあるコードへの適用が難しいといった問題があった。このため我々は、クローンの検出において構文解析を用いず、コードブロック内の処理の形状面から類似性に着目したクローン検出法(CICD:Code Image Clone Detection method)の開発を進めている。

本論では、CICDの中で利用している文字出現率を用いたコードブロックの類似性判定法について、簡易な実験結果も交えて報告する。

2. CICDの概要

本章ではCICD法のおおよその流れを紹介する。CICD法では次のStep1からStep4の4ステップによってクローンの検出を行う[1]。

Step1:対象ソースコードを処理ブロックに分割。

Step2:各ブロック内の総文字数、総行数を利用したクローン候補ブロックの特定。

Step3:候補ブロックの各行ごとの文字数比較によるクローンブロックを特定。

Step4:特定したクローンブロックをクローンとして確定するためにブロック内の顕著な特徴を抽出して確定診断を行う。

図1はStep3までの処理で特定されるクローンブロックの様子を示す。図1の例では、コードA、コードBの2つのコードについて、それぞれ3つの処理ブロックに分けられる。そして、そのうちの①と⑥、②と④、③と⑤が形状的に

行番号	コードA	文字数	計	コードB	文字数	計
1	for(x=2;k=0;	13		int a=0,i,j;	11	B[0]
2	while(k<=6){	12		for(i=0;i<=4;j++){	18	34(4)
3	printf("%d\n",x);	19	A[0]	a++;	4	
4	x++;	4	5(6)	}	1	
5	k=x;	4			0	
6	}	1		while(a<=20){	13	B[1]
7		0		j=j+i;	6	26(4)
8	int a=0,i,j;	11		a++;	3	
9	for(i=0;i<=4;j++){	18	A[1]	}	1	
10	a=a+1;	6	26(4)		0	B[2]
11	}	1		printf("Hello world!\n");	25	25(1)
12		0		}	1	
13	while(a<=20){	13	A[2]	int b=rand(2);	13	
14	j=j+i;	6	24(4)	switch(b){	10	B[3]
15	a++;	4		case 0:return 1;	14	52(6)
16	}	1		case 1:return 2;	14	
17		0	A[3]	}	1	
18	printf("Hello world!\n");	25	25(1)	}	0	

図1 CICD法によるクローンブロック候補抽出

クローンペア候補	クローンペア
1	①と⑥
2	②と④
3	③と⑤

類似するクローン候補ブロックとして特定できることを示している。

本論ではこのようにして抽出されたクローンブロック候補に対し、各ブロック内で出現する特定文字の出現率を利用して、より高い精度でクローン判定を行うことを目的としている。

3. 特定文字出現頻度を利用したクローン判定

Step4では、クローン候補となるコードブロック内の特定文字の出現頻度を比較し、クローンであるかどうかを高精度で判定する。特定文字には、コードブロック内の処理特徴をあらわす文字を利用する。特定文字は対象言語の特徴から抽出されるものもある。

具体的には、各ブロック内での論理構造や演算処理面の処理特徴を表す文字として、C言語の場合であれば、次のものを候補として考える。

特殊文字[記号]: <, >, +, -, /, %, &=, (,) など
 アルファベット: for, while, if, else, do, switch, case など

例えば、図1で候補として抽出したブロック③と⑤、およびブロック②と④の場合(図2)、それぞれの特徴文字の出現数は表1のようになる。この表をみると、どちらのブロックの組も出現する特殊文字の数に大きな差は見られず、処理内容や論理構造が同じであれば、特定文字の出現数もほぼ同じになると考えることができる。

コードA-ブロック③	コードB-ブロック⑤	コードA-ブロック②	コードB-ブロック④
while(a<=20){ j=j+a+i; a++; }	while(a<=20){ j=j+a+i; a++; }	int a=0,i,j; for(i=0;i<=4;j++){ a=a+1; }	int a=0,i,j; for(i=0;i<=4;j++){ a++; }

図2 特定文字を利用したブロック類似評価

Code clone candidate detection technique with characters' appearance ratio

†Nihon University

‡Oki Electric Industry Co., Ltd.

表1 特定文字の出現数評価結果

		<	+	=	i	f	w
ブロック	A③	1	4	2	2	0	1
	B⑤	1	4	2	2	0	1
ブロック	A②	1	3	4	4	1	0
	B④	1	4	3	4	1	0

4. 特定文字の選択実験

前章に示したように特定文字を利用して、クローンブロックの処理特徴を評価する場合、特定文字としてどの文字を採用するかが重要となる。このため、ここでは、実際のC言語のオープンソースを評価材料として、

- a. どのような文字が個々のブロック内に含まれるか、
- b. その文字の中でどの文字を評価する特定文字にすべきか

を決定するために実験を行った。

4.1 実験方法

評価材料は下記のオープンソースの2ファイルから一部を無作為に抜粋したものを使用する。

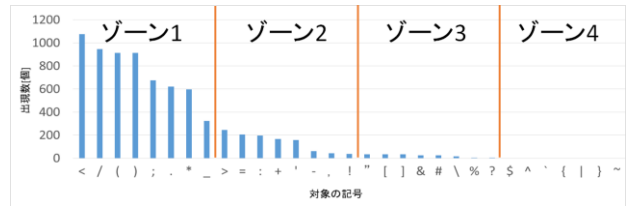
- ・Borland C++ Compiler 5.5 stdlib
- ・indent ver2.2.9 intl

実験ではこれらのファイルをCICDのStep1を適用し処理ブロックに分割し、各ブロック内にどのような文字が含まれるかをカウントした。この際、ブロックはif, whileなどの論理構造をもつ複雑な処理ブロックと、それらを持たない単純構造の処理ブロックの2タイプに分けて評価を行った。

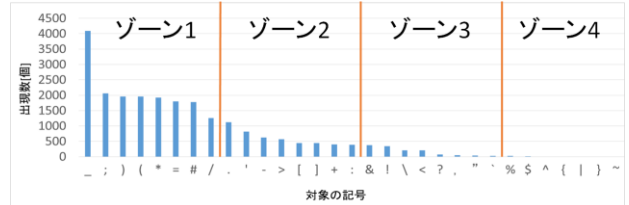
4.2 評価結果

(a) 特殊文字[記号]

図3(1)はstdlibの複雑ブロック内の記号出現数を、図3(2)はintl内のブロック全体における記号の出現数を示したものである。この2つの図はそれぞれのブロック内での記号の出現数を頻度順に表示しており、出現頻度別に上位から四分の一ずつにゾーンで区切っている。この中で第1ゾーンはどのようなブロックでも比較的多く出現する記号であるため処理特徴を表現する文字としては適していないと考える。また、第3, 4ゾーンは出現が非常に希なものであり、個々の処理ブロックに含まれない場合もあるため処理特徴を評価する記号としては適していないと考えられる。これらを考えるとブロックの特徴を表現する文字としては第2ゾーンに含まれる文字を採用することが良いと考えられる。このゾーンには、処理内容の加算+や減算-比較>などを表現する文字が含まれており、これらをブロックの処理特徴を評価する文字として利用することが有効と考えられる。



(1) stdlib 複雑ブロック

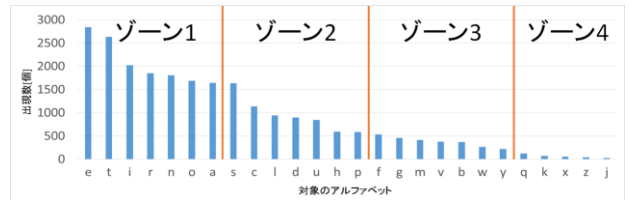


(2) intl 全体

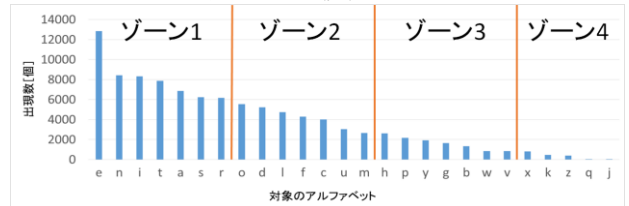
図3 記号の出現頻度

(b) アルファベット

図4(1)(2)にアルファベットの出現頻度を示す。記号と同様に4つのゾーンに分けると、この場合も第2ゾーンに含まれる文字を評価用の特殊文字に採用するのが適切と考えられる。このゾーンには、for, while, caseなどに含まれるf, h, cなどが入っており、これらの文字は処理の特徴をとらえる文字として利用することが有効と考えられる。



(1) stdlib 複雑ブロック



(2) intl 全体

図4 アルファベット(小文字)の出現頻度

5. まとめ

本論ではクローンブロックに含まれる特定文字によるクローンの判定精度向上法のアイデアを紹介した。また、簡易実験により判定時に利用すべき特定文字候補の絞り込みについて検討した。今後、この結果をさらに精査し、特定文字による判別の精度を上げていきたい。

参考文献

[1] 垣谷 広輝, 安藤 優作, 平山 雅之, 菊地 奈穂美, "ソースコードを構成する処理ブロックの特徴に着目したコードクローン推定技術", 情報処理学会研究報告ソフトウェア工学研究会, Vol. 2013-SE-182(29), 2013.