

データマイニング手法を用いたプロダクトメトリクス類似性に基づくソフトウェアインスペクション方法の提案

M2012MM028 西川 雅彦

指導教員 青山 幹雄

1. はじめに

ソフトウェアの再利用によって開発が進行すると再利用元で発見できなかった Fault が再利用先に継承される。再利用元と再利用先のプロダクトメトリクス類似性に注目し、データマイニング手法を用いたインスペクション方法を提案する。

2. 研究課題

個人スキルに依存したインスペクションではなく、データマイニングツールを活用した数学的な分析によって Fault が発生する要因を特定、予測を行うことで精度の高いソフトウェアインスペクションを提案する。

3. 関連研究

- (1)ソフトウェア欠陥の有無を推定する研究がある[1].
- (2)インスペクションの工数最適化の研究がある[2].

4. アプローチ

データマイニング手法を用いてプロダクトの類似度で Fault 発生を予測する。仮説検証として、仮説 1 は、再利用ソフトウェアとの類似度、仮説 2 では、再利用元と再利用先とのプロセスの類似度を明らかにする。

5. データマイニングに基づく分析方法の提案

データベース化された再利用元と再利用先のメトリクスをデータマイニング手法によって類似度比較し、プロセス類似バグと構造類似バグに分類することでインスペクション精度の向上策を提案する(図 1)。

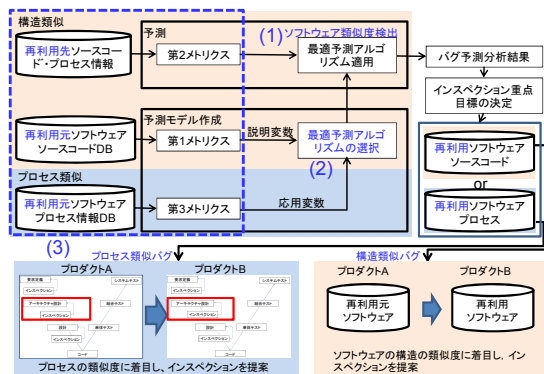


図 1 データマイニングに基づく分析方法の提案

(1)ソフトウェアの類似度の検出

仮説 1 構造類似バグに基づき、Fault が発生したメトリク

スの類似度に着目し、Weka [3]を利用して潜在する再利用先の Fault を明らかにする。ソフトウェアの類似度予測値を評価し、Rating 比較することでソフトウェアの類似度を検証する。

(2)最適アルゴリズムの選択

Weka に実装されている複数の予測アルゴリズムを数種類選択し、再利用元と再利用先の類似度の予測値を評価し、Rating を比較することで最も精度が高いアルゴリズムを選択する(表 1)。

表 1 最適アルゴリズムの選択方法

予想結果	内容
Precision (適合率)	Trueと分類された中で、実際にTrueであるモジュールの割合
Recall (再現率)	全体のTrueモジュールの中で正しくTrueと分類された割合
F-measure (F値)	適合率と再現率の調和平均
ROC area	ROC曲線下の面積は分類アルゴリズムの精度の高さを表す

予測結果の項目ごとの精度にRatingを付与し、その平均が最大値10に近いアルゴリズムが最適アルゴリズムとなる。

(3)メトリクスデータ管理

再利用元と再利用先のバグ発生原因のプロセス情報メトリクスをデータベースで管理し、Weka の木構造解析を行うことで、再利用元と再利用先について統計的手法を用いて比較、検証する(図 2)。

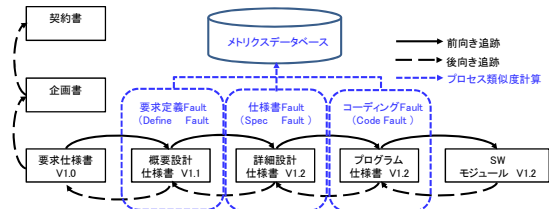


図 2 要求追跡方向と Fault 発生エラーの関係図

6. 実データへの提案方法の適用

6. 1. データセット

(1)プロダクトの選択

ソフトウェアのメトリクスと Fault データとして Promise Software Engineering Repository[4] を利用する対象は、宇宙機器(CM1, KC1), 科学データ処理(KC2), 機載ソフトウェアである。

(2)メトリクスの選択

Promise で公開されているプロダクトのメトリクスは 22 種類用意されているが、プロダクトに共通するメトリクスのみを選択して使用する。本研究で使用したメトリクスは主としてソース容量、ソフトの構成の複雑度、工数と個別にソフトウ

エアの性格を表す個別案件で分類できるようにした(表2).

表2 複雑度マトリクスの種類と内容

マトリクスの種類	内容
ソース容量	McCabe's line count of code
	Halstead "volume"
	Halstead "program length"
複雑度	McCabe "cyclomatic complexity"
	McCabe "essential complexity"
	Halstead "difficulty"
	Halstead "intelligence"
工数	Halstead "effort"
その他	Halstead's count of lines of comments
	IOCodeAndComment: numeric
	total operators
	total operands of the flow graph

6. 2. 仮説1 構造類似性バグの検証

再利用元と再利用先を比較し、予測アルゴリズムによってマトリクスの類似度を評価する。

(1) 検証方法

表1のデータを、Wekaを使用して予測精度が最も近いデータを類似プロダクトデータとし、再利用元ソフトウェアと再利用ソフトウェアと定義する。

表2で示す複雑度のマトリクスによる予測精度が、

再利用元予測精度 < 再利用予測精度

の関係であれば、再利用元ソフトウェアに基づいて仕様変更したことでソフトウェア構造が複雑となったことを表す。

ソフトウェアモジュール中に Fault の発生割合が、

再利用元ソフトウェア < 再利用ソフトウェア

の関係であれば、再利用ソフトウェアで Fault が引き継がれていることを表す。

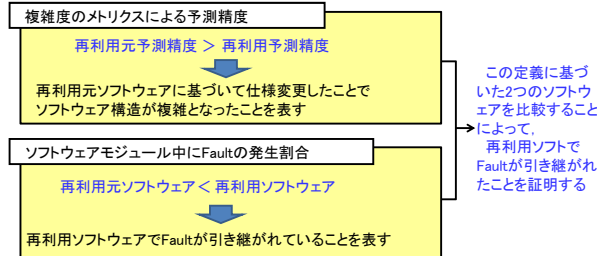


図2 構造類似性バグの検証方法概念

(2) 最適アルゴリズムの選択

最適アルゴリズムの選択方法としては、Wekaの予測結果の項目ごとの精度にRatingを付与し、その平均が最大値10に近いアルゴリズムが最適アルゴリズムとなる。

Wekaで実装されている各カテゴリで代表的なアルゴリズムを抽出し、Precision(適合率)、Recall(再現率)、F-measure(F値)、ROC areaの観点でRating付与した他のアルゴリズムと比較したところ、MultilayerPerceptronが最も予測精度が高いことが分かった(表3)。

表3 アルゴリズムの予測精度評価結果

Category	Name	Precision	Rating	Recall	Rating	F-Measure	Rating	ROC Area	Rating	Rating Avg.
trees	J48	0.970	8	0.952	8	0.957	8	0.972	4	7.0
	J48graft	0.950	2	0.937	5	0.941	5	0.874	5	4.3
	RandomForest	0.928	3	0.937	5	0.929	2	0.966	3	3.3
	LADTree	0.900	1	0.921	1	0.905	1	0.989	1	1.0
function	Logistic	0.937	4	0.937	5	0.937	3	0.983	7	4.8
	MultilayerPerceptron	0.987	10	0.984	10	0.985	10	1.000	10	10.0
bayes	BayesNet	0.977	9	0.968	9	0.971	9	0.997	9	9.0
	BayesSimple	0.950	5	0.937	5	0.941	5	0.986	8	5.8
meta	AttributeSelectedClassifier	0.957	7	0.952	8	0.954	7	0.964	2	6.0
	PART	0.957	7	0.952	8	0.954	7	0.983	7	7.3

• Precision (適合率) : Trueと分類された中で、実際にTrueであるモジュールの割合
 • Recall (再現率) : 全体のTrueモジュールの中で正しくTrueと分類された割合
 • F-measure (F値) : 適合率と再現率の調和平均
 • ROC area : ROC曲線下の面積は分類アルゴリズムの精度の高さを表す

(3) 類似プロダクトデータの選択

表3のソフトウェアに最適アルゴリズムMultilayerPerceptronの予測精度を算出した(図3, 4)

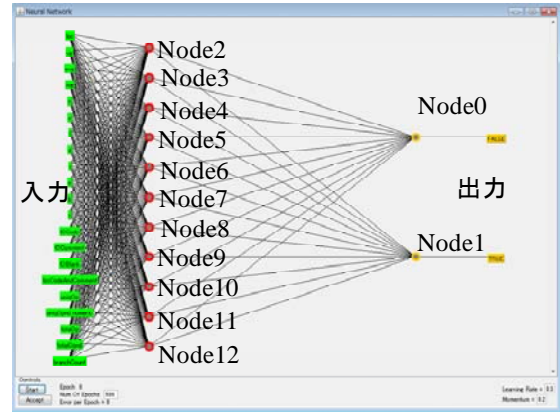


図3 WekaによるMultilayerPerceptronの入出力関係図

```

    == Summary ==
    Correctly Classified Instances 436      87.5502 %
    Incorrectly Classified Instances 62      12.4498 %
    Kappa statistic -0.0156
    Mean absolute error 0.1568
    Root mean squared error 0.3121
    Relative absolute error 87.6482 %
    Root relative squared error 104.785 %
    Total Number of Instances 498

    == Detailed Accuracy By Class ==

    TP Rate  FP Rate  Precision  Recall  F-Measure  ROC Area  Class
    0.969    0.98    0.901    0.969    0.933    0.734    FALSE
    0.02     0.031    0.067    0.02    0.031    0.734    TRUE

    Weighted Avg. 0.876  0.886  0.819  0.876  0.845  0.734

    == Confusion Matrix ==
    a b <- classified as
    435 14 | a = FALSE
    48  1 | b = TRUE
  
```

図4 WekaによるMultilayerPerceptronの計算結果

得られた予測精度をもとに表4に示す通りRatingを付与し、同レベルの類似性をもったソフトウェアをCM1, KC1, KC2とした。JM1とPC1はRating結果が他の3つのソフトウェアと比較して大きく異なっているため、類似プロダクトデータとして選択ができなかったことが分かった(表4)。

表4 類似プロダクトデータの選択結果

Category	Precision	Rating	Recall	Rating	F-Measure	Rating	ROC Area	Rating	Rating Avg.
CM1	0.819	2	0.876	4	0.845	4	0.734	3	3.3
JM1	0.769	1	0.810	1	0.741	1	0.690	1	1.0
KC1	0.835	4	0.859	3	0.828	2	0.771	4	3.3
KC2	0.834	3	0.847	2	0.835	3	0.828	5	3.3
PC1	0.921	5	0.936	5	0.917	5	0.723	2	4.3

(4) 仕様変更前後のソフトウェアの選択

表5に示す通り、類似プロダクトデータの中で仕様変更後のソフトウェアを選択した。

再利用元ソフトウェアから再利用ソフトウェアを作成すると、ソフトウェア構造が複雑化して各データ間の類似性が低くなる。したがって、複雑度を表すマトリクスとFault発生の有無に限定したうえで、Wekaの予測結果を基にRatingを付与した。

全体のソフトウェア数が3であるため、複雑度が低いソフトウェアをRating 3.0とし、仕様変更前でソフトウェアの複雑度が高くなる前の状態を定義する。逆に複雑度が高いソフトウェアをRating 1.0とし、仕様変更後でソフトウェアの複雑度が高くなった状態を定義する。

再利用元類似性 Rating > 再利用ソフト類似性 Rating の関係が成り立てば、類似性の高低によって仕様変更前後の関係となる(表5)。

表5 仕様変更前後のソフトウェアの選択結果

Category	Precision	Rating	Recall	Rating	F-Measure	Rating	ROC Area	Rating	Rating Avg.
CM1	0.843	3.0	0.892	3.0	0.859	3.0	0.691	1.0	2.5
KC1	0.830	2.0	0.856	2.0	0.830	2.0	0.788	2.0	2.0
KC2	0.814	1.0	0.830	1.0	0.816	1.0	0.829	3.0	1.5

ただし、この時点では、CM1, KC1, KC2 のいずれのソフトウェアが再利用元あるいは、再利用先ソフトウェアであるのか不明であるため、CM1, KC1, KC2 のマトリクス比較によって ①式を満たす組み合わせを○で表記し、それ以外を×で表記する(表 6)。

表 6 再利用元と再利用先複雑度 Rating マトリクス比較

	CM1 平均Rating 2.5	KC1 平均Rating 2.0	KC2 平均Rating 1.5
CM1 平均Rating 2.5		○	○
KC1 平均Rating 2.0	×		○
KC2 平均Rating 1.5	×	×	

仕様変更前の再利用元ソフトウェアと、仕様変更後の再利用先ソフトウェアを①式を満たす組み合わせをまとめると以下の 3 パターンの組み合わせとなる(表 7)。

表 7 アルゴリズムの予測精度評価結果

仕様変更前		仕様変更後	
再利用元ソフトウェア	Rating Avg.	再利用先ソフトウェア	Rating Avg.
CM1	2.5	KC1	2.0
CM1	2.5	KC2	1.5
KC1	2.0	KC2	1.5

これによって、再利用元ソフトウェアと再利用先ソフトウェアの組み合わせの定義を決めることができた。

(5) 再利用ソフトウェアの Fault の引き継ぎ確認結果

複雑度のマトリクスの算出結果を基に定義した再利用元ソフトウェアと、再利用先ソフトウェアの Fault の発生割合が再利用元の Fault 発生割合 < 再利用先 Fault 発生割合であれば再利用元ソフトウェアの Fault が引き継がれ、再利用先のソフトウェアの Fault の発生が増大したといえる。そこで、各ソフトウェアのモジュール中の Fault 発生割合を求めた(表 8)。

表 8 各ソフトウェアモジュール中の Fault 発生割合

Category	①Faultなし モジュール数	②Faultあり モジュール数	②/①*100
CM1	449	49	10.9 %
KC1	1,783	326	18.3 %
KC2	415	107	25.8 %

次に表 6 に示した通り、再利用元ソフトウェアと再利用先ソフトウェアの組み合わせの定義結果に対して、表 7 の複雑度のマトリクスの算出結果を基に定義した再利用元ソフトウェアと、再利用先ソフトウェアの Fault の発生割合が同じであれば、以下の関係が成り立つ。

再利用元複雑度 Rating > 再利用先複雑度 Rating

複雑度の高低によって再利用元(仕様変更前)と再利用先(仕様変更後)の関係が成立する(複雑度小=Rating 大, 複雑度大=Rating 小)。

再利用元の Fault 発生割合 < 再利用先 Fault 発生割合
再利用元ソフトウェアの Fault が引き継がれ、再利用先のソフトウェアの Fault の発生が増大したことを明らかにすることができる。

図 5 に示した通り、【A】再利用元ソフトウェアの Fault の発生割合、と【B】再利用先ソフトウェアの Fault の発生割合を比較すると全ての組み合わせで【A】<【B】の関係が成り立ち、実データ上で Fault が引き継がれていることを明らかにすることができた。

これにより、再利用先のソフトウェアの Fault 発生を予測する手段として、予測値の最も高いアルゴリズムを用いて、モジュールの複雑度を算出し、Rating によって Fault 発生の高さを数値化することで Fault モジュールを特定することができる。

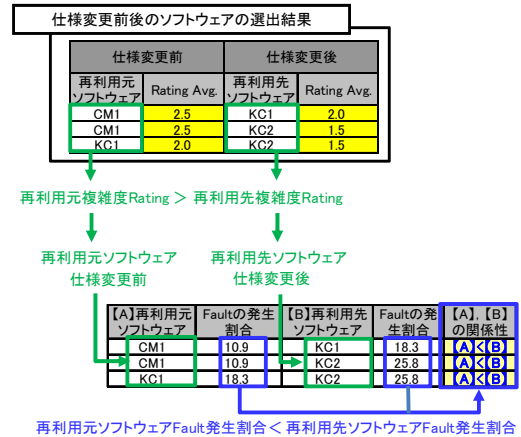


図 5 再利用ソフトウェアの Fault の引き継ぎ確認結果

6. 3. 仮説 2 プロセス類似性バグの検証

要求仕様書に記述されている要求項目が正当なものであるかどうかを確認するためには、仕様書の前提となった文書に当たって確認を取らなければならない。こうした作業を支援する手段として要求追跡がある。

要求を追跡するには、文書の作成時に、追跡リンクを設定しておく、要求追跡を実現するためには、それぞれの要求と、要求変更が公的な文書として文書化され、構成管理の下に管理されている必要がある。

再利用ソフトウェアにおける要求仕様書のプロセスへの期待は、熟練分析者の経験と知識を基にした属人的情報を記録し、その情報から、犯さなくとも済むはずの失敗を再び繰り返さなくすむ可能性を高めることである。

要求仕様書には、ユーザーの要求が、出来るだけ正確かつ完全な形で表現されていなければならない、しかし、限られた時間内に、ユーザーが対象領域の知識と、システムに対する個別の要求のすべてを分析者に伝えることは不可能に近い。またそれらの要求を分析者が正しく理解し、正確に記述することも困難である。そのような場合、類似した問題領域での分析者の過去の経験が、ユーザーの知識の欠如と、分析者の理解不足を補うことができる。

そのため要求追跡のデータ管理方法として ①要求定義のエラー、②コードエラー、③仕様書のエラーの分類を定義し、Fault 発生原因を特定することで得られる検証結果をマトリクス化することで、プロセス類似度に Fault 発生の予測値を求めることを可能とする。

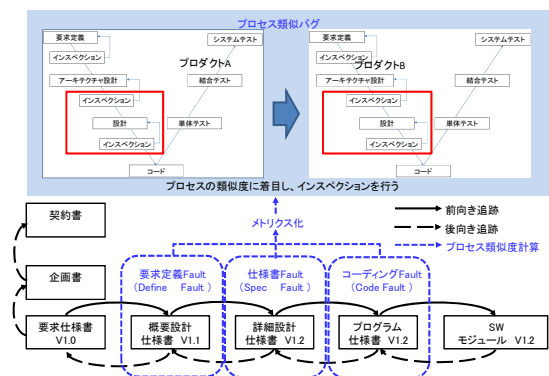


図 6 要求追跡の方向とプロセス類似バグ計算図

再利用データとそれ以外のデータを分類し、Fault の原因を①要求定義のエラー、②コードエラー、③仕様書のエラーのいずれかのエラーであるか Weka の木構造によって

検証する。

(1) 検証方法

表9に示す通り、再利用ソフトウェアでFaultが発生したデータを①～③に分類した。Wekaで実装されているTreesカテゴリの中から代表的なアルゴリズム4個のうち最も精度が高いものを選択する。

また、木構造によって定量的評価を行うため、それ以外のカテゴリのアルゴリズムと比較して、Treesのアルゴリズムで得られる予測精度に妥当性があるか検証する。

表9 再利用バグの原因の一覧

Base No.	Development phase	Foundation of Demand	Details	BUG by "YOKOTEN"	Mar day	Fault type	Existence of a requirement specification definition based on "YOKOTEN"
2	A	new		FALSE	10	Define Fault	FALSE
34	B	new old	Improvement in an actuator operation	FALSE	15	Define Fault	FALSE
111	B	new old	Improvement in an actuator operation	FALSE	15	Define Fault	FALSE
4	A	new old		FALSE	10	Spec Fault	FALSE
10	B	new old		FALSE	15	Spec Fault	FALSE
33	B	new old		FALSE	15	Spec Fault	FALSE
15	B	soft program		FALSE	15	Spec Fault	FALSE
16	B	soft program		FALSE	15	Spec Fault	FALSE
20	B	soft program		FALSE	15	Spec Fault	FALSE
29	C	spec system	Improvement in an actuator operation	FALSE	10	Code Fault	FALSE
43	C	spec system	Improvement in an actuator operation	FALSE	10	Code Fault	FALSE
44	C	spec system	Improvement in an actuator operation	FALSE	10	Code Fault	FALSE
47	C	spec system	Improvement in an actuator operation	FALSE	10	Code Fault	FALSE
50	C	spec system	Improvement in an actuator operation	FALSE	10	Code Fault	FALSE
61	E	spec system	Improvement in an actuator operation	TRUE	10	Code Fault	FALSE
81	E	spec system	Change state correction	TRUE	10	Code Fault	FALSE
82	E	spec system	Change state correction	TRUE	10	Code Fault	FALSE
37	C	spec system	Change id	FALSE	10	Code Fault	FALSE
42	C	spec system	Change id	FALSE	10	Code Fault	FALSE
46	C	spec system	Change id	FALSE	10	Code Fault	FALSE
49	C	spec system	Change id	FALSE	10	Code Fault	FALSE
14	C	spec system	Omission of specification	TRUE	10	Code Fault	TRUE
15	C	spec system	Timing control correction	FALSE	10	Code Fault	FALSE
17	C	spec system	Judgment control correction	FALSE	10	Code Fault	FALSE

(2) 最適アルゴリズムの選択

表10に示すとおりTreesのカテゴリを代表するアルゴリズムで分類精度の確認を行ったところ、J48-C 0.25-M2が最も高いことが分かった。他のアルゴリズムには、J48-C 0.25-M2よりも高い分類精度を持つものがあるが、10個のアルゴリズム中7位の分類精度であり、絶対値比較としてもAve. 0.963は95%以上の信頼区間を得ているためJ48-C 0.25-M2のアルゴリズムは妥当性であるといえる。

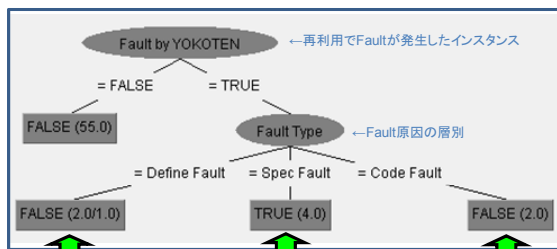
表10 最適アルゴリズムの選択結果

Category	Name	Precision	Recall	F-Measure	ROC Area	Rating
trees	J48	0.970	0.952	0.957	0.972	7.0
	J48graft	0.950	0.937	0.941	0.874	4.3
	RandomForest	0.928	0.937	0.929	0.966	3.3
	J48Tree	0.900	0.921	0.905	0.959	1.0
functions	Logistic	0.937	0.937	0.937	0.983	4.8
	MultilayerPerceptron	0.987	0.984	0.985	1.000	10.0
bayes	BayesNet	0.977	0.968	0.971	0.997	9.0
	BayesSimple	0.950	0.937	0.941	0.986	5.8
meta rules	AttributeSelectedClassifier	0.957	0.952	0.954	0.954	6.0
	PART	0.957	0.952	0.954	0.983	7.3

7. 評価結果

(1) Fault原因の層別化

図7に示す通り、Wekaに実装されている木構造の視覚化により、Fault発生の要因を解析する。再利用ソフトウェアで発生したFaultの中から発生原因を①要求定義Fault、②コーディングFault、③仕様書Faultに層別し再利用ソフトのプロセスの問題点を明らかにした。



YOKOTENに基づく要求仕様定義の存在しないクラス
YOKOTENに基づく要求仕様定義の存在がクラスに属するモノが2つあるが、間違った分類は1である
YOKOTENに基づく要求仕様定義の存在しないクラスに属するモノが4つあるが、間違った分類は0である
YOKOTENに基づく要求仕様定義の存在しないクラスに属するモノが2つあるが、間違った分類は0である

図7 木構造のビジュアル結果

(2) Fault発生箇所の読み取り

図8に示す通り、再利用先のソフトウェアに発生したFaultの原因の中に、ソフトウェア要求仕様書に基づいたプロセスが存在している場合の分類精度を表している。

①要求定義Faultは、クラスに属する総数に対して、50%正しい結果が得られていないので除外した。

②コーディングFaultは、再利用元に要求仕様の定義が無くても、再利用されたプログラムのFaultを引き継いだ。

③仕様書Faultは、再利用元に要求定義があり、要求定義と仕様書に乖離がある

仮説2のプロセスの類似性については①、②、③のいずれかが再利用元のプロセスと同じであれば、Faultが発生し得ると考えられる。この結果から、要求仕様の定義に基づき実装しても、ソフトウェアの複雑性に起因するFaultが発生している可能性が明らかとなった。

また、図8に示した通り、②仕様書Faultは概要設計から詳細設計にプロセスが移行するときに、要求された振る舞いを実現できていないためにFaultが発生した可能性があることを明らかにした。

これによって、6.2(5)の再利用ソフトウェアのFaultの引き継ぎ確認結果から、Faultが発生した再利用元と再利用先との類似度が高ければFault発生の可能性が高まり、また、プロセスに課題があるときは、再利用元のプロセスFaultが再利用先に引き継がれることが分かった。

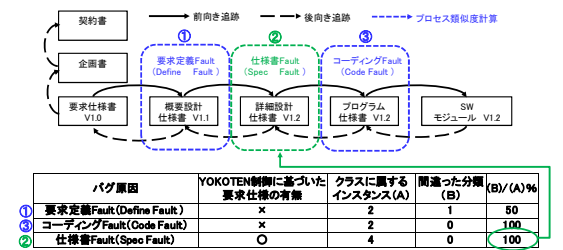


図8 評価結果概要図

8. 今後の課題

本稿の予測手段を適用し、工数低減量を評価する。また開発プロセスはプロダクトの成熟とともに改良されていくため、類似度比較するデータ対象が常に同じであるとは限らない。比較できない追加プロセスは欠損データとなってしまうため、類似度への影響を確認する必要がある。

9. まとめ

統計的観点からデータマイニングツールを活用し、簡単な分析によってソフトウェア開発プロセスとソフトウェアコードの類似度を比較し、Faultが発生する要因を特定、予測を行うことでソフトウェアインスペクション精度を向上する手法を提案した。提案方法の有効性を示すために、実データを基にしたマトリクスを用いて類似度によるソフトウェアFaultの発生を特定できることを確認した。

参考文献

[1] R. Subramanyan, et al., Empirical Analysis of CK Metrics for Object-Oriented Design Complexity, IEEE TOSE., Vol. 29, No. 4, 2003, pp. 297-310.
 [2] 阿萬 祐久, 山下 裕也, 整数計画法を用いた重点レビュー対象モジュールの選択, コンピュータソフトウェア, Vol. 27, No. 4, 2010, pp. 240-245.
 [3] Weka: Univ. of Waikato, <http://www.cs.waikato.ac.nz/ml/weka/>.
 [4] J. S. Shirabad, et al., The PROMISE Repository of Software Engineering Databases, Univ. of Ottawa, <http://promise.site.uottawa.ca/SERepository/>.