

プログラム・テストの妥当性評価に関する実験報告

— ソフトウェア品質 = テストの質 × テストの量 —

大場 穗 (日本アイ・ビー・エム 製品保証)

〔概要〕

人工的に作ったバグをテスト対象プログラムに移植しておき、テスト完了後移植したバグがテスト中にどの程度除去されたかによってプログラム・テストの妥当性を量的に評価しようとする方法に関する実験結果を報告する。

1. はじめに

プログラム・テスト期間中に発見される機能欠陥(バグ)の数の推移は、巨視的に次の平均値函数 $m(t)$ をもつ非定常ボアソン過程でよく記述されると言われている。^{1) すなわち、テスト開始後 t 時間目における発見欠陥数 $m(t)$ は：}

$$m(t) = \bar{N}(1 - e^{-bt}), \quad (1)$$

で与えられる。ここで、パラメータ \bar{N} はそのテストで発見可能な機能欠陥の総数を、パラメータ b はテストの効率を示す。(1) 式のモデルは図1に示されている曲線の性質をもつ。すなわち、

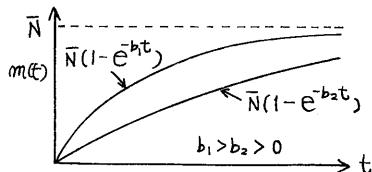


図1. 平均値函数 $m(t)$ の性質

\bar{N} が一定の場合、 b の値が大きいほど短時間で欠陥を数多く発見できる。ところでテストで発見可能な欠陥の総数 \bar{N} とは、テスト開始から完了までに行ったテストと同様のものを無限に実施したと仮定した場合に発見できると推定される欠陥の総数を言う。従って、 \bar{N} はテストを行う作業担当者の欠陥発見能力や、テストに必要なツール類に

依存して定まるもので、対象プログラム中に潜伏するバグそのものの総数を示すとは限らない。テストで発見可能な欠陥の総数 \bar{N} とプログラム中に潜伏するバグの総数が一致しない原因として、次の2点が考えられる：

- ① テストが不十分で発見の難度の高いバグが未発見のままテストが完了した。
- ② プログラム中に潜伏したバグが少量のためパラメータ N の推定誤差が大きい。

以上のような問題のため、平均値函数 $m(t)$ をもつ非定常ボアソン過程モデルでソフトウェアの品質やテストの妥当性を評価することは実用上問題がある。すなわち、前掲の非定常ボアソン過程モデルではテストの量的側面のみしか評価することができない。従って、テストの質が良い場合や、プログラムの機能が単純なオフライン環境で動作するソフトウェアの場合には良く適合するが、プログラムの機能が複雑で高度なテスト技術を必要とするオンライン環境で動作するソフトウェアの場合、非定常ボアソン過程モデルは現実に適合しない例が観測される。

上記のような問題点を改善し、プログラム・テストの妥当性を客観的に評価する方法として、H. ミルズによて提案された Capture-Recapture 法を基礎に^{2,3)}、非定常ボアソン過程モデルによるテストの量的評価に対しこうにテストの質的評価を加味する尺度 SSPQL を提案し、その実験結果について報告する。

2. SSPQL

プログラム・テストの妥当性評価、

すなむちテスト完了時ににおけるソフトウェアの品質評価尺度，SSPQL (Shipping Software Product Quality Level) は，次式のように定義される⁴⁾：

$$SSPQL = \alpha(t_m) \times \gamma(t_m), \quad (2)$$

ここで， $\alpha(t_m)$ はテスト開始から完了までの時間を t_m とするとき，テスト完了時ににおけるテストの精度を示す尺度であり， $\gamma(t_m)$ はテスト完了時ににおけるテストの網羅度 (Test Coverage) を示す尺度である。テストの精度はテストの質的側面を評価する尺度に相等する。すなむちプログラム中に潜伏するバグをテスト中にどの程度の確実さで捕捉したかを示す。またテストの網羅度はテストの量的側面を評価する尺度に相等する。すなむちプログラム中に潜伏しているバグの中で，テスト中に欠陥として発見可能なもののうちどの程度の量をテスト中に発見したかを示す。⁽²⁾ 式で定義される SSPQL は，テストの精度と網羅度の両方の尺度が高い(1に近い) 値をとらない限り高い値とならない。このことは，いかにバグの捕捉率の高いテストを行っても，テストがプログラムの一部の機能のみに対して行われていれば，プログラム中に潜伏するバグは十分に除去されてしまう。ソフトウェアの品質も十分でないという現実に良く合致する。また同様に，いかに網羅的なテストを行ったとしても，バグを見逃していればソフトウェアの品質は十分使用に耐えるものとはなりえない。その意味で，平均値函数 $m_c(t)$ をもつ非定常ホアン過程のモデルでは，テストの精度が十分に高くない場合，発見可能な欠陥数とプログラム内に潜伏する欠陥数との差が大きくなるので，現実と合致しない場合がある。

3. テストの精度

テスト完了時，すなむちテスト開始

後 t_m 時間経過後，におけるテストの精度 $\alpha(t_m)$ は次式で定義される：

$$\alpha(t_m) = \frac{\bar{N}_c}{N_c} = \frac{m_c(t_m)}{N_c}, \quad (3)$$

ただし， N_c はプログラムに潜伏する欠陥の総数， $\bar{N}_c = m_c(t_m)$ は潜伏する欠陥のうちテストで発見可能なもう1つの総数とする。ところで，現実にプログラム中に潜伏する欠陥の総数を十分に精度で推定することは困難である。そこで，本小論では人工的に作成した欠陥(制御欠陥と呼ぶ)をテスト開始直前にプログラムへ移植し，その制御欠陥の捕捉率を基礎に， $\alpha(t_m)$ を推定する方法を採用する。すなむち， N_c をプログラムに移植した制御欠陥の総数， \bar{N}_c をテストで発見可能な制御欠陥の総数とするとき， $\alpha(t_m)$ は：

$$\alpha(t_m) = \frac{\bar{N}_c}{N_c} \approx \frac{\bar{N}_c}{N_c}, \quad (4)$$

で推定される。⁽⁴⁾ 式において N_c は既知であるので，非定常ホアン過程モデルの平均値函数 $m_c(t)$ のパラメータが推定できれば， $\bar{N}_c = m_c(t_m)$ なる関係から $\alpha(t_m)$ が推定できる。ここで時刻 t_m でに発見された制御欠陥の総数を \bar{N}_c とすれば，テスト期間中， $(t_0, 0)$ ， (t_1, y_1) ，…， (t_m, y_m) を一連のデータが得られるので，最尤法により $m_c(t)$ の平均値函数のパラメータを推定できる。例として ⁽¹⁾ 式の平均値函数を考えれば， (t_i, y_i) ($i=0, 1, \dots, m$) をデータに対して，Goel & Okumoto の方法により尤度函数を最大にする $m_c(t)$ のパラメータを推定できる。すなむち：

$$m_c(t) = \bar{N}_c (1 - e^{-bc}) , \quad (5)$$

なる \bar{N}_c および bc を推定できる。ところで，テストが十分に行われていれば，推定値 \bar{N}_c は， $N_c \geq \bar{N}_c$ を満足するので，テストの精度 $\alpha(t_m)$ は：

$$0 \leq \alpha(t_m) \leq 1, \quad (6)$$

を満足する。

4. テストの綱ら度

テスト完了時、すなはちテスト開始後 t_m 時間経過後、にみけるテストの綱ら度 $\gamma(t_m)$ は次式で定義される：

$$\gamma(t_m) = \frac{m_i(t_m)}{\bar{N}_i} = \frac{m_i(t_m)}{m_i(t_\infty)}, \quad (7)$$

ここで、 \bar{N}_i はテストで発見可能な真の欠陥の総数、そして $m_i(t_m)$ は時刻 t_m までに発見された真の欠陥の総数とする。 (7) 式において $m_i(t_m)$ はテスト期間中に発見された真の欠陥の数であり、また \bar{N}_i は非定常ホアソン過程モデルの平均値函数 $m_i(t)$ によって与えられる。すなはち、テスト中時刻 t_m までに発見された真の欠陥の総数を Z_i とすれば、テスト期間中、 $(t_0, 0), (t_1, Z_1), \dots, (t_m, Z_m)$ 、なる一連のデータが得られるから、最尤法により $m_i(t)$ なる平均値函数のパラメータを推定することができる。例として (7) 式の平均値函数を参考されば、 (t_i, Z_i) ($i=0, 1, \dots, m$) なるデータに対して、尤度函数を最大にする $m_i(t)$ のパラメータ：

$$m_i(t) = \bar{N}_i(1 - e^{-b_i t}), \quad (8)$$

なる \bar{N}_i および b_i を推定できる。

ところで、一般に中小規模ソフトウェアのテストにおける欠陥発見過程を追跡すると、図2の(A)または(B)の2種類の曲線のどちらかに分類できる⁵⁾。

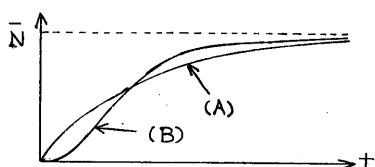


図2. 欠陥発見過程の追跡曲線

のことから、われわれは現実のテスト結果でソフトウェアの品質評価やテストの妥当性評価を行う場合、 (7) 式の平均値函数だけでは不十分なことを知ることができる。現実のテスト結果が図2(B)のような曲線となる原因と

して、以下の2点が考えられる：

① テスト開始直後はテスト担当者がプログラムやテスト環境（テストシステム、テスト・ツール等）に習熟していないため、有効なテストを実施していない空白時間（遅れ）が生じる。

② テスト担当者が欠陥らしい現象を発見してから欠陥と確認するまでに再試行等現象を再現させるための努力をするため、テスト全体にめたり時間遅れを生じる。

上記①の現象は、 (7) 式のモデルを以下のように変形することで記述される：

$$m_i(t) = \bar{N}_i(1 - e^{-b_i(t-d)}). \quad (9)$$

また、上記②の現象は次々平均値函数モデルによつて記述される⁶⁾：

$$m_i(t) = \bar{N}_i(1 - (1+b_i t)e^{-b_i t}). \quad (10)$$

(9) 式のモデルは、 (7) 式の平均値函数を時間軸上でシフトさせたもので、基本的に (7) 式と同じ性質をもち、自動制御で言う time-lag をもつたモデルである。 (10) 式のモデルは、 (7) 式のモデルが自動制御で言う 1 次遅れの系であるのに対し、2 次遅れの系の性質をもつている。著者らの実験結果では、3,000 LOC (Line of Code) から 20,000 LOC 程度の規模のモジュールの欠陥発見過程を分析する場合には、 (9) 式または (10) 式のモデルのどちらかが良く適合した。また、10,000 LOC 以上の規模のプログラム全体の欠陥発見過程を 1 回の平均値函数で表現することは誤差のため実用上問題がある。モデルの適合性については、一般的に機能の単純なモジュール、小規模なモジュールに対しては (9) 式の time-lag をもつモデルが良く適合するのに対し、機能の複雑な制御プログラム等のモジュールなどでは (10) 式の Delay をもつモデルが良く適合した。また、80,000 LOC 以上のプログラムでは、 (10) 式のモデルが巨視的に適合していた。

5. 実験例 I

実際に制御欠陥を 3,000 ル 0 C 程度の小さなプログラムへ移植し、オンライン環境での機能試験を完了した直後 S S P Q L を評価した例について述べる。

実験対象のプログラムは構造化プログラミング用マクロ命令を用いてアセンブリ言語で書かれた、オンライン端末装置を制御するための入出力制御用プログラムである。実験データは通信回線や端末装置を接続しないオンラインの環境で、通信回線および端末装置の状態や応答をシミュレートする「バック・ツール」を用いて、プログラムの機能が正しく動作することを検査する目的の試験にあって得られたものである。

本実験において、制御欠陥はプログラムの基本設計を担当した評価者の立場の人間と、実際にプログラムの詳細設計、コード化、および機能試験を担当した被験者の立場の人間とにより、ほぼ半数ずつ移植された。移植された欠陥は次の 3 種類に分類される：

- ① 欠陥を通過した場合常に異常終了となるもの、
- ② 欠陥を通過しても通常は機能異常を生じないが、特定の条件またはその後の処理で機能異常の原因となるもの、
- ③ 上記①または②に該当し、事前のテストまたは精査で欠陥と判断されたバグを残留させたもの。⁷⁾

上記①の欠陥は、発見の比較的容易なバグである。例としては、GETMAIN マクロで獲得したバッファを FREEMAIN マクロで開放する際に、開放する記憶域の長さ指定を獲得時の指定と矛盾させておくといふような當時異常終了となるバグである。上記②のバグは、発見の困難度の高い制御欠陥である。例として、制御ブロックの参照・更新時に

参照・更新フィールドの相対番地をスライせていくなどのバグである。制御ブロックのポインタを更新するルーチンにこの種の制御欠陥を移植した場合、この欠陥の通過時には何の機能異常も発生しないが、その後の処理で機能異常を発生する。上記③の欠陥は、プログラムに本来潜在していたもので、そのプログラムに潜伏しているバグがどの程度除去されたかを知る目安となる。例としては、制御ブロック内に定義されているフラグ(1 ビットのスイッチ)をリセットするために排他的論理和を実行する命令を使用していたものが事前のコード精査で確認され、それを残留させたものなどがある。この例の場合、この命令実行時に当該フラグのオン状態が保証される場合は問題を生じないが、当該フラグがオフ状態である場合には、遂にオン状態に反転しその後の誤動作の原因となる。以上のことにも加えて、制御欠陥の移植に際しては、コード精査で誤りの多く発見されたモジュールや、論理の複雑化モジュールなどには多くの制御欠陥を移植するよう配慮した。

図 3 に本実験例での真の欠陥および制御欠陥の発見過程を示す。実験期間は 10 労働日(約 200 人時)であった。

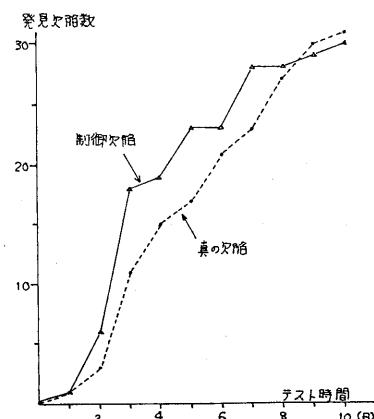


図 3. 実験結果

図3から明らかなように、本実験例の場合、(9)式の平均値函数モデルは適合していなかった。本例の真の欠陥の場合について(9)式および(10)式のモデルを仮定した平均値函数のパラメータの推定結果を図4および図5に示す。ちなみに、実測値とモデルによる推定値との2乗誤差 E^2 :

$$E^2 = \sum_{i=1}^m [Z_i - m_i(t_i)]^2, \quad (11)$$

の比較では、(9)式の(Tim-lag)モデルが(10)式のモデルに勝っていた。しかし、推定されたパラメータ \bar{N}_i を比較すると、(9)式のモデルの場合が52、(10)式のモデルの場合が38となっており、大きな差があった。ここで、制御欠陥の漏ら度が90%に達していたため、真の欠陥の漏ら度が60%以下

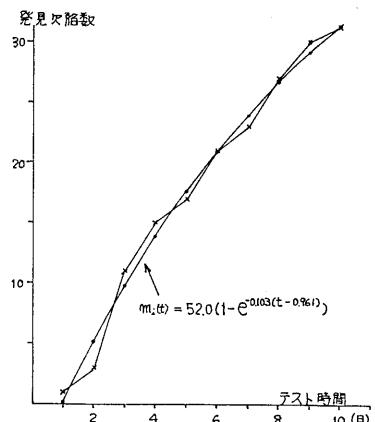


図4.

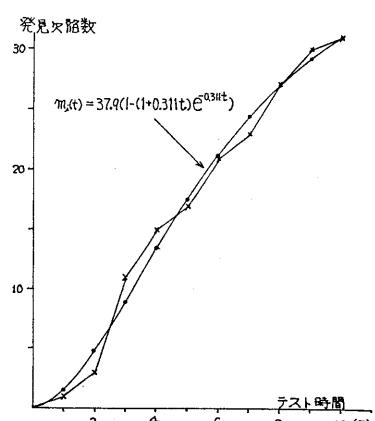


図5.

であるとは考えにくないので、本例については(10)式の(Delay)モデルが適合していると判断した。

次に、制御欠陥のデータに関する解析では、(9)式および(10)式のモデルで試算した結果、2乗誤差 E^2 およびパラメータ \bar{N}_i の推定値の両方ともほとんど差がなかった。図6に(10)式の(Delay)モデルによる推定結果を示す。

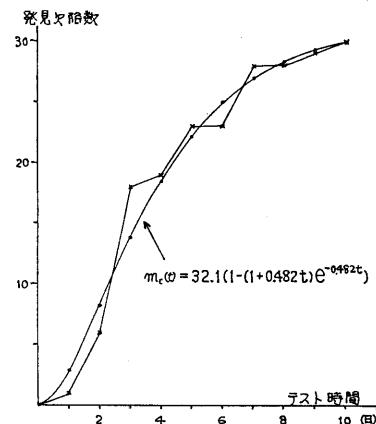


図6. 制御欠陥

以上、データ解析結果より、本例におけるテストの漏ら度とは:

$$\begin{aligned} r(t_{10}) &= \frac{31}{m_i(t_{10})} \\ &= \frac{31}{37.9(1-(1+0.32t_{10})e^{-0.312t_{10}})} \\ &= 0.818, \end{aligned} \quad (12)$$

となる。また、テストの精度 α は:

$$\begin{aligned} \alpha(t_{10}) &= \frac{m_i(t_{10})}{37} \\ &= \frac{31.3(1-(1+0.498t_{10})e^{-0.498t_{10}})}{37} \\ &= 0.846, \end{aligned} \quad (13)$$

となる。従って、SSPQLは:

$$\begin{aligned} SSPQL &= 0.818 \times 0.846 \\ &= 0.692, \end{aligned} \quad (14)$$

となる。

本例においてテスト完了時における

SSPQL 約 0.7 と低いのは、テスト担当者が制御プログラムのテストに未経験であるためであり、現象として発生している欠陥を異常として認識できていない点や、テスト・ケースが正常処理の流れのテストに偏っている点から考えると、妥当な値であろう。ちなみに、本実験の完了後、現実のオンライン環境におけるシステム統合テスト(5 労働日、延 30 人時)において同一テスト担当者により、5 種類の新しい真の欠陥が発見された。また、それに続く約 3 週間の実使用期間中に確認された真の欠陥が 2 種類(機能影響の小さなもの)、他に確認されていない異常現象が数種類報告された。以上のことがら、テストの網羅度、精度ともほぼ妥当値であったと言えよう。

6. 実験例Ⅱ

本実験例は、全体で約 10,000 LOC の 3 フックモジュールから構成されるプログラムについて SSPQL を実験的に適用した例の一端である。ここでは、約 4,000 LOC の入出力実行制御モジュールの評価例を紹介する。プログラムはアセンブリ言語で書かれており、簡単な構造の IF 文、CASE 文、GO TO 文、BEGIN 文、END 文、および DEBUG 文等のマクロ命令が準備されている以外は全てアセンブリ命令でコード化されている。

本例の入出力実行制御モジュールは、プログラムの基本設計を担当した当該アセンブリ言語でのプログラミング経験者が、詳細設計、コード化、およびデバッグを担当した。制御欠陥の移植は、他の関連モジュールの詳細設計、コード化の担当者、事前の詳細設計精査やコード精査のデータを基に行なった。特に本例の制御欠陥は、品質の要求水準が高かったため、発見の困難度の高いバグを多く移植した。

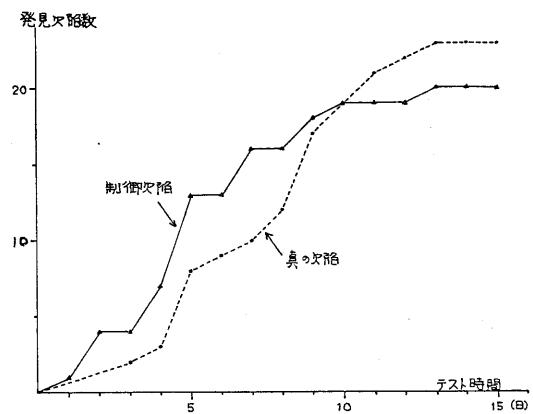


図 7. 実験結果

図 7 に本例での真の欠陥および制御欠陥の発見過程を示す。ラスト期間は 14 日間、発見した真の欠陥の総数は 23、発見した制御欠陥の総数は 20 である。図 7 から明らかのように、本例の場合 (1) 式の平均値回数モデルは適合していない。真の欠陥について (9) 式および (10) 式のモデルを仮定し、それらの平均値回数のパラメータを推定した結果、実測値とモデルによる推定値との二乗誤差の比較、およびパラメータ N; の推定値の妥当性の比較とともに (10) 式の (Delay) モデルがよくデータに適合していた。図 8 に (10) 式のモデルによる推定曲線と実測値を示す。

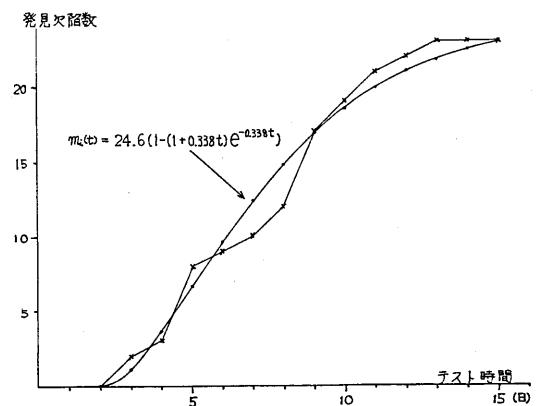


図 8. 真の欠陥

次に、制御欠陥のデータについては全部で26個の制御欠陥が移植されたが、テスト中に確認されたものは全部で20個であった。本例のテストはモジュールの特定機能を検定する目的で行われるモジュールの他の機能には注目していないため、テスト対象外の下位モジュールがあった。テスト対象となる下位モジュールに移植された制御欠陥の数は、26個中21個である。従って、テスト対象である下位モジュール中に未発見の制御欠陥が1個残存している。制御欠陥の発見データについて、(9)式および(10)式のモデルを仮定し、平均値函数のパラメータを推定した結果、実測値とモデルによる推定値との2乗誤差の比較では(10)式のモデルが約1.5倍勝っていたが、推定したパラメータ \bar{N}_c の値に差はなかった。図9に(10)式の(Delay)モデルによる推定値と実測値との関係を示す。

以上のデータ解析の結果よりテストの網羅度は：

$$\begin{aligned} r(t_{15}) &= \frac{23}{m_c(t_\infty)} \\ &= \frac{23}{24.6(1-(1+0.338t_\infty)e^{-0.338t_\infty})} \\ &= 0.935, \end{aligned} \quad (15)$$

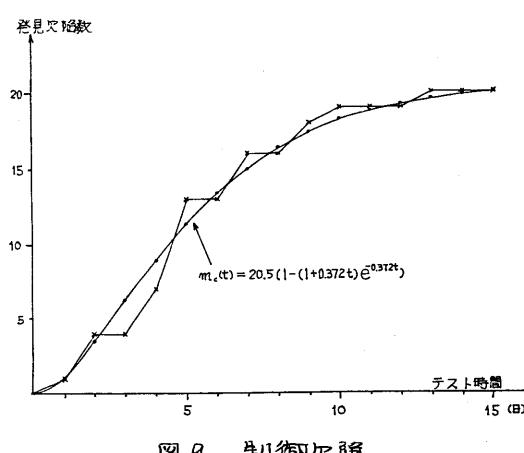


図9. 制御欠陥

となる。また、テストの精度 α は、テスト対象となる下位モジュールに移植された制御欠陥が21個であるところから：

$$\begin{aligned} \alpha(t_{15}) &= \frac{m_c(t_\infty)}{21} \\ &= \frac{20.5(1-(1+0.372t_\infty)e^{-0.372t_\infty})}{21} \\ &= 0.976, \end{aligned} \quad (16)$$

となる。従って、本例の場合SSPQLは：

$$\begin{aligned} SSPQL &= 0.935 \times 0.976 \\ &= 0.913, \end{aligned} \quad (17)$$

となる。

本例の場合、SSPQLが約0.9とかなり高い値である原因として、プログラムの設計・コード化および本例のテスト担当者が、テスト・ツールおよびプログラミング言語の経験がある点が挙げできる。ちなみに、プログラミング言語とテスト・ツールに未経験のプログラマが担当したモジュールについては、SSPQLは約0.6であった。本例のテスト完了後、新たに発見されたテスト対象下位モジュールの真の欠陥は、1個のみであった。

7. まとめ

以上、ソフトウェアの品質評価尺度およびテストの妥当性評価尺度としてのSSPQLについて述べた。従来の信頼度成長曲線によるソフトウェアの品質評価は前述のように、テストの質に評価が依存するという問題点があった。また、Capture-Recapture法による評価では、制御欠陥の移植の方法（分布のさせ方）によって評価結果が左右されるという難点があった。SSPQLでは、これらの問題を改善するため評価をテストの質と量の独立した2つの評価に分け、テストの量的評価を従来の信頼度成長曲線のモデルによって行い、ま

たテストの質的評価を Capture-Recapture 法によって行うこととした。このことにより、SSPQLでの評価値は信頼度成長曲線のモデルに比較してテストの担当者の能力やテスト・ツールの能力等の影響を小さくするとともに、制御欠陥の移植法の影響をも少さくした。

SSPQL の推定には信頼度成長曲線モデルの現実のデータへの適合性が重要となる。このため、本小論では信頼度成長曲線モデルとして従来から知られている(1)式のモデルとそれを拡張した(9)式の(Time-lag)モデルに加えて、新たに(10)式の(Delay)モデルを採用した。実験例からも明らかのように、小規模プログラムやそのモジュールに対する(10)式の(Delay)モデルが一般的に適合性が良く、実用的である。たゞ、(1)式や(9)式のモデルでは、発見可能な欠陥数を示すパラメータ N の値が非現実的な(大きな)値となる場合が多く、特に絶対数の未知である真の欠陥については有効でない場合が多い。ただし、制御欠陥についてはその分布が他の欠陥に比較して各モジュールへ一様に分散しているため、時間軸上のズレを考慮した(9)式の(Time-lag)モデルも(10)式の(Delay)モデルとほぼ同程度に適合していた。

SSPQL の評価値の妥当性については、さらに中規模ソフトウェアでの実験が必要であろう。ただ現在までの実験結果から、定性的には SSPQL の評価値の高いプログラムのテスト後の信頼度が高いことが確認されている。従来の信頼度成長曲線モデルのみの評価に比較して、信頼度の高いプログラムの品質評価の確度が向上している。しかししながら、逆に信頼度の低いプログラムの品質評価が、従来の評価法に比較して厳しくなる傾向がある。すなまち、SSPQL の評価値が 0.6 から 0.7 程度のプログラムでも基本的な機能に

ついては不安定ながら使用可能である場合がある。これらは、SSPQL の評価が稠度と精度とのカケ算で定義され、現実には精度を評価する時に稠度の影響があることなどが原因と考えられる。今後、この点についての検討が重要であろう。

謝 辞

本研究を行うに際し、ソフトウェアの欠陥に関するデータ収集について御助力を頂きました日本アイ・ビー・エム(株)製品保証担当佐藤武義氏ほか、同所属藤岡義彦氏、および同武田健二氏に感謝します。

[参考文献]

- 1) Goel A., Okumoto K.: "Time-dependent error-detection rate model for software reliability & other performance measures," IEEE Trans. Reliability, Vol.R-28, 1979.
- 2) Duran J, et al.: "Capture-Recapture Sampling for Estimating Software Error Content," IEEE Trans. Software Eng., Vol. SE-7, 1981.
- 3) 伊士誠他: "Capture & Recapture 法による潜在バグの推定法とその応用", 情報学会 ソフトウェア工学研究会資料 19, 1981年7月.
- 4) 篠田東吾他: "制御欠陥と機能試験の妥当性評価", 情報学会 22回全国大会 5C-7, 1981年3月.
- 5) 酒巻恒一他: "ソフトウェアの信頼性:品質管理のためのソフトウェアの信頼性予測," 信学会 信学校報, (信頼性) Vol. 81, 1981年5月.
- 6) 山田茂他: "ソフトウェアエラー発見過程の確率モデルとその応用", 情報学会 ソフトウェア工学研究会資料, 1981年11月.
- 7) 坂本忠彦他: "バグ残留法によるプログラムの品質測定," 情報学会 21回全国大会 7C-8, 1980年5月.