

対象別評価の導入によるソフトウェアプロセス最適化の検討

3W-2

銭谷由彦 山口智治 白石智

NTT ネットワークサービスシステム研究所

1.はじめに

ソフトウェアの開発において、大規模になれば扱うデータも増大し、工程管理・品質管理の評価はよりシビアなものになる。本稿では、大規模ソフトウェア開発における膨大な開発情報を分析し（データマイニング）、会社、個人のレベルを対象としたプロセス評価を可能とし、その値の組み合わせによりプロセスの最適化をはかる方法を提案する。

今回はデータマイニングの手法として相関・判別分析のアルゴリズムの適用を検討しており、報告する。

2.現在の開発支援環境

図1にあるように開発作業のデータを蓄積し、作業標準書に定められた指標値を計算し、進捗遅れや品質に問題があれば開発管理者に警告するシステムを構築している。

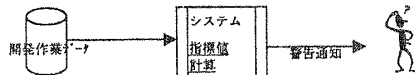


図1 自動警告システム

3.データマイニング

データマイニングとは大規模ソフトウェア等の開発作業において日々蓄積される膨大な生データから、従来のデータベース検索では発見することの出来なかった有用な情報、例えば、データ間や事象間の関連、法則、傾向などを発見する事である。これを実行することにより、従来の作業標準書には載っていない情報を獲得でき、新たなノウハウとして活用できる。このノウハウは開発が続く限り増えつづけ、しかも陳腐化しない。

具体的には、例えば「設計作業でのドキュメント品質の良し悪しが試験作業の進捗に影響をもたらす」といった「法則」や、「このベンダでは製造進捗はいつもほぼこのぐらいである」といった「傾向」を導く事である。

技術的にはソフトウェアプログラムの項目単位（ex.クラス群）毎に多次元の要素で表現し、その多次元要素を既存データベースからデータ取得したものをプロセス特性基本データとする。要素の取り方には留意する。

<ex. n個の要素を持つプロセス特性基本データ>
 クラス群 Aaa = (x1, x2, x3, x4, x5, x6, ..., xn)
 = (設計進捗分析値, ドキュメント品質分析値, 製造進捗分析値, 規模変動分析値, 試験進捗分析値, 試験品質分析値, ..., 問題処理分析値)

そして全てのプロセス特性基本データを多次元に解析（ex. 相関係数計算、判別分析、回帰分析等）し特性を抽出し、「法則」や「傾向」とする。

(方式1)

相関分析・・・要素間の相関を数値で表す

相関分析とは

・・・一般に相関係数（ピアソンの積率相関係数）Rxy は次の様に定義される。

ある2つの変数 x、y について n 個のデータが得られているとき、

$$R_{xy} = \frac{x \text{ と } y \text{ の共分散}}{x \text{ の標準偏差} * y \text{ の標準偏差}}$$

$$R_{xy} = \frac{1/n \sum (x_i - \bar{x})(y_i - \bar{y})}{\sqrt{1/n \sum (x_i - \bar{x})^2} * \sqrt{1/n \sum (y_i - \bar{y})^2}}$$

(ここで、 $\bar{x} = 1/n \sum x_i$, $\bar{y} = 1/n \sum y_i$)

を計算し相関の仮説を立てる分析法である。（図2）
 Rxy は-1 から 1 の値をとる（-1 ≤ Rxy ≤ 1）

例えば x2 と x6 の相関係数の計算結果により仮説を立てる
 x2 と x6 の相関係数が高い x2 と x6 には強い関連がある
 x2 と x6 の相関係数が低い x2 と x6 は関係がない

（相関係数計算により x2（ドキュメント品質分析値）と x6（試験進捗分析値）の相関係数が 0.8（相関高）という値が出たら「ドキュメント品質が比較的悪い状態のまま製造作業に移ればいずれ試験作業時に進捗遅れとして影響が出る。」という仮説が成り立つ。）

(方式2)

判別分析・・・ある要素の判別を他の要素を分析することで表す判別分析とは

・・・判別したい2群の平均値とサンプルのばらつきも考慮した距離（マハラビスの距離）を求め、より近い方のグループに属していると判定する分析方法。

(1) 既存のデータ(変数 k 個)を集め、第 i 群の重心を (x1i, x2i, ..., xki), 第 i 群に属するデータを用いて計算した分散共分散行列を Σi とする

(2) 判別したいデータとそれぞれの群との重心との偏差ベクトルを

$$U_i = \begin{pmatrix} x1 - x1i \\ x2 - x2i \\ \dots \\ xk - xki \end{pmatrix} \text{ とする}$$

(3) マハラビスの距離 $D_i = U_i^T \Sigma_i^{-1} U_i$ を計算する

(4) 距離が最小の群に属すると判定する。

xn がどちらの群（満足行く結果 or 問題のある結果）になるか事前に判別したい場合、サンプルデータから判別用の値(x1, x2, x3, x4, ..., xk)を取り出し値を計算し判断する。

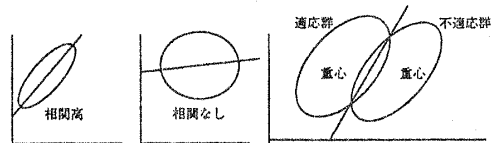


図2 方式1 相関分析

図3 方式2 判別分析

具体例)

設計ドキュメントのレビュー品質（コメント密度、エラー密度、レビュー密度）について作業標準書で定めた閾値に達しないドキュメントがあった場合を例にとる。

通常であればこのようなドキュメントについては再レビュー等の判定をしているが、本データマイニングにおいては必要に応じてプロセス特性基本データに設計作業に携わる開発担当者の個人データを一部追加して分析を行う。通常では見落としがちな問題点や、その問題の解決となる仮説を導くところに重点を置いて膨大な生データから関連を導き出せそうなデータを絞って分析を行うこととする。

Examination of software process optimization by introduction of evaluation according to object

Yoshihiko ZENIYA Tomoharu YAMAGUCHI Satoshi SHIRAIISHI

NTT Network Service Systems Laboratories

9-11 Midori-Cho 3-Chome Musashino-Shi, Tokyo 180-8585 Japan

すなわち

<分析に必要な追加データ>

個人 Bbb= (個人 ID、ベンダ ID、設計スキルレベル、開発経験年数、
エラー発見率、エラー原因率、再レビュー時エラー発見率、
再レビュー時エラー原因率…)

などのデータを抽出し、上記レビュー品質データと相関分析を行う。

それにより、相関のあるもの、ないものが相当数得られるので一部挙げて

- (相関高1) エラー密度が規定値に達しないドキュメントの開発担当者
と開発経験3年未満の開発担当者の相関は高い
- (相関高2) レビューに参加したメンバーの平均開発経験年数とエラー密度は正比例する
- (相関高3) レビュー実施時間(12~24)とエラー発見率は反比例する
- (相関無1) レビュー参加者の設計スキルレベルとエラー発見率は特に相関が無い。

また過去の開発データを利用することにより、後工程のデータも考慮した判別分析も行う。

それにより以下の判別情報が得られるので挙げて

- (判別1) 担当ベンダは、設計工程で再レビューしたドキュメントにより製造したプログラムに対する試験作業において、作業標準値で定めた閾値の4/5以上のバグは試験期間内に検出できる
- (判別2) 担当ベンダが再レビューしたドキュメントの製造作業において製造進捗遅れが発生する可能性は低い(確率20%)
- (判別3) 再レビューしたドキュメントにより製造されたプログラムに対する問題処理票の対応日数は再レビューとならなかったものと判別はできない。

これらの相関関係、及び、判別情報から次の法則、傾向を導くことが出来る。

- (法則1) 開発経験3年未満の設計ドキュメント数の何割かは再レビューになる
- (法則2) レビューは平均開発年数を考慮してトータルの人選を行う。
- (法則3) レビューは夜遅くなるようなら無理をせず翌日行う
- (法則4) レビューの個人人選は設計スキルレベルは考慮しなくてよい

- (傾向1) 担当ベンダの再レビューしたドキュメントはほぼ目標バグ密度値を満足しそうである。
- (傾向2) 担当ベンダの再レビューしたドキュメントの製造進捗遅れはなさそうである。
- (傾向3) 担当ベンダの再レビューしたドキュメントに対する問題処理対応能力は通常の場合と大差なさそうである。

4. プロセス評価

データマイニングにより抽出された情報を元に、開発作業を会社、個人レベルで細かく評価していき、具体的な数字を提示することにより、各々最善をはかるように取り組ませることがより現実味を帯びたものになる。どこまで値をあげなくてはいけないという明確な目標が与えられるからである。その単位としては、作業標準書で定めた工程単位である場合もあるし、工程をまたがった作業でもそれを一つのプロセスとしてまとめるのが良い場合もある。したがって、

従来の作業標準書にはない柔軟な作業単位の設定が必要である。

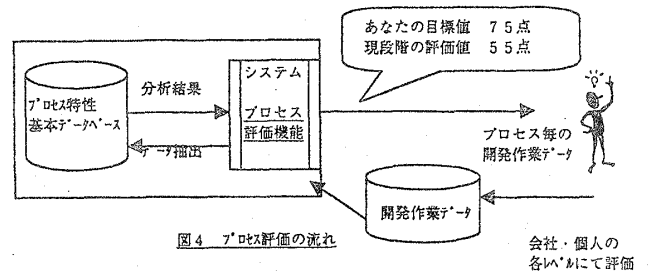


図4 プロセス評価の流れ

会社・個人の各レベルにて評価

5. プロセス最適化

データマイニングにより抽出された情報により開発者の各単位でプロセス評価が実現できれば、それを最適化することが一番求められていることである。つまり一つのプロセスの開始から終了までとり着く最良の道を求めることである。具体的には図5のように、開発作業をプロセス単位で区切って、各プロセスの優良な評価値(>平均値)を設定し、それにとり着くための「法則」や「傾向」をタイムリーに開発当事者に知らせ、その実行状態を把握し、道から外れないようにフォローしていく流れをイメージしている。

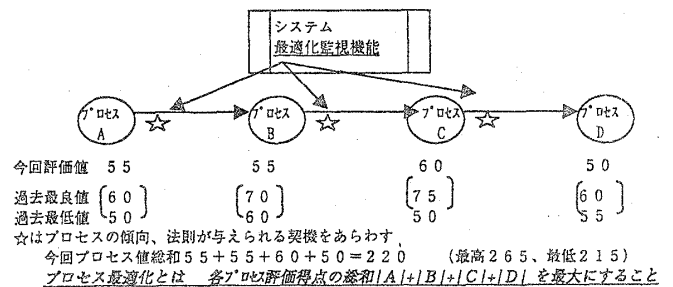


図5 プロセス最適化の流れ

6. 本方式による効果

本方式により、開発作業のプロセス評価対象が明確になり、その評価も客観的に行え、その総和を最大にする努力をフォローすることによりプロセス評価値の最適化、つまりはプロセス最適化につながる効果が見込まれる。

7. まとめ

本稿では、ソフトウェア開発における膨大な情報をデータマイニングによりプロセス評価に役立て、プロセス最適化に利用する構想を示した。今後詳細な検討を進めて実現に必要な仕様を確認する。

参考文献

- [1]現代人の統計2 新版 多変量解析法 柳井・高根共著(朝倉書店)
- [2]ソフトウェアプロセス成熟度の改善 Watts S.Humphrey 著 藤野訳(日科技連)