

University of Groningen

Proposing and empirically validating change impact analysis metrics

Arvanitou, Elvira Maria

IMPORTANT NOTE: You are advised to consult the publisher's version (publisher's PDF) if you wish to cite from it. Please check the document version below.

Document Version

Publisher's PDF, also known as Version of record

Publication date:
2018

[Link to publication in University of Groningen/UMCG research database](#)

Citation for published version (APA):

Arvanitou, E. M. (2018). *Proposing and empirically validating change impact analysis metrics*. [Thesis fully internal (DIV), University of Groningen]. University of Groningen.

Copyright

Other than for strictly personal use, it is not permitted to download or to forward/distribute the text or part of it without the consent of the author(s) and/or copyright holder(s), unless the work is under an open content license (like Creative Commons).

The publication may also be distributed here under the terms of Article 25fa of the Dutch Copyright Act, indicated by the "Taverne" license. More information can be found on the University of Groningen website: <https://www.rug.nl/library/open-access/self-archiving-pure/taverne-amendment>.

Take-down policy

If you believe that this document breaches copyright please contact us providing details, and we will remove access to the work immediately and investigate your claim.

Downloaded from the University of Groningen/UMCG research database (Pure): <http://www.rug.nl/research/portal>. For technical reasons the number of authors shown on this cover page is limited to 10 maximum.



university of
 groningen

Proposing and Empirically Validating Change Impact Analysis Metrics

PhD Thesis

to obtain the degree of PhD at the
 University of Groningen
 on the authority of the
 Rector Magnificus Prof. E. Sterken
 and in accordance with
 the decision by the College of Deans.

This thesis will be defended in public on

Friday 13 July 2018 at 09.00 hours

by

Elvira Maria Arvanitou

born on 16 September 1988
 in Thessaloniki, Griekenland

Supervisors

Prof. P. Avgeriou

Prof. A.N. Chatzigeorgiou

Co-supervisor

Prof. A. Ampatzoglou PhD

Assessment Committee

Prof. A.C. Telea

Prof. F. Arcelli Fontana

Prof. T. Mens

ISBN: 978-94-034-0752-4

to my baby....

Samenvatting

Bij het onderhoud van de software is de software onderhevig aan veranderingen als gevolg van het verhelpen van fouten, veranderende eisen, aanvullende eisen, enz. Het belang van het laag houden van de onderhoudskosten is benadrukt in de literatuur met empirisch bewijs, waaruit blijkt dat de kosten van deze fase ongeveer 50%-75% uitmaken van de totale kosten van softwareontwikkeling. Deze kosten kunnen nog verder oplopen als de software wordt gekenmerkt door (a) *veranderingsgevoeligheid*, d.w.z. de waarschijnlijkheid dat een softwareartefact om interne redenen verandert (d.w.z. het oplossen van bugs in dat artefact of het wijzigen van eisen met betrekking tot het artefact); en (b) *instabiliteit*, d.w.z. de waarschijnlijkheid dat een softwareartefact verandert als gevolg van veranderingen in andere artefacten van het systeem.

Het proces dat veranderingsgevoeligheid en instabiliteit onderzoekt wordt *change impact analysis* genoemd; dit proces is niet alleen belangrijk tijdens het onderhoud, maar ook tijdens de andere ontwikkelingsfasen, zoals requirements engineering, ontwerp, implementatie en testen. Change impact analysis is gebaseerd op de kwantificering van veranderingsgevoeligheid en instabiliteit, om beslissingen te kunnen nemen over welke veranderingen moeten worden uitgevoerd, en hoe en wanneer. In de literatuur zijn echter de volgende beperkingen vastgesteld met betrekking tot de praktijken van effectbeoordelingen van veranderingen:

- In de fasen van het ontwerp van het architecturale ontwerp en de behoeftetechniek ontbreken meetwaarden die de instabiliteit van het artefact en de gevoeligheid van de verandering bepalen. Daarom kan de impactanalyse in die fasen niet kwantitatief zijn.
- De implementatie- en detailontwerpfasen worden ondersteund door meetwaarden voor instabiliteit en veranderingsgevoeligheid, maar dergelijke meetwaarden zijn niet nauwkeurig omdat ze de bovengenoemde parameters niet combineren.

- Er is een gebrek aan hulpmiddelen die het proces van het berekenen van veranderingsgevoeligheid en instabiliteitsmetriek voor implementatie en gedetailleerde ontwerpfasen kunnen automatiseren.

Om de bovenvermelde beperkingen aan te pakken, hebben *we* als *onderdeel van dit doctoraat een reeks methoden en instrumenten voorgesteld die zowel instabiliteit als veranderingsgevoeligheid met verhoogde nauwkeurigheid kunnen kwantificeren in drie belangrijke ontwikkelingsfasen* (d.w.z. requirements engineering, ontwerp en implementatie).

Als eerste stap om dit doel te bereiken, hebben we eerst de literatuur bestudeerd om ontwerp-tijd kwaliteitsattributen en meetwaarden te onderzoeken, die kunnen worden gebruikt om ze te kwantificeren. Specifiek hebben we een corpus van meer dan 150 primaire studies beoordeeld. De belangrijkste bevindingen benadrukken het belang van veranderingsgevoeligheid en instabiliteit in de huidige onderzoeksliteratuur. Met name stabiliteit en veranderingsgevoeligheid zijn de meest bestudeerde kwaliteitsattributen, na onderhoudbaarheid (respectievelijk 15 en 18 studies). De studie bevestigde echter het gebrek aan meetwaarden in verschillende ontwikkelingsfasen (vooral in eisen en architectuur); de meeste beschikbare meetwaarden bevinden zich op het broncodeniveau.

Gebaseerd op de belangrijkste bevinding van de vorige studie, hebben we als volgende stap de mogelijkheid van codemeetwaarden onderzocht om toepasbaar te zijn op artefacten in een andere ontwikkelingsfase, namelijk architectuur. Daarom onderzochten we meerdere codemeetwaarden, met betrekking tot hun vermogen om fijn- en grofkorrelige veranderingen in softwareonderhoud vast te leggen. Intuïtief passen meetwaarden die in staat zijn om grofkorrelige veranderingen vast te leggen beter bij het architectuurniveau; integendeel, fijnkorrelige veranderingen moeten worden verwaarloosd op het architectuurniveau. De empirische evaluatie suggereerde dat sommige meetwaarden gevoeliger zijn voor veranderingen, en dus meer geschikt zijn voor evaluaties op methodiek- en klasseniveau, terwijl andere stabielere zijn, d.w.z. dat ze grootschaliger veranderingen vereisen om hun

waarden te beïnvloeden, en dus meer geschikt zijn voor het architectuurniveau (bv. pakketten, componenten, enz.). Op basis van de resultaten is de enige meetwaarde die gerelateerd is aan instabiliteit en in staat is om grofkorrelige veranderingen vast te leggen de Response voor a Class (RFC)-meetwaarde, door klasse-niveaumeetwaarden samen te voegen tot het architectuurniveau met behulp van de gemiddelde (AVG)-functie. RFC combineert koppeling (d.w.z. gebruik van de openbare interface van andere klassen) en grootte (d.w.z., aantal lokale methoden) eigenschappen en daarom kan het niet worden beschouwd als een zuivere veranderbaarheid of instabiliteitsmeetwaarde. Daarom wordt het voorstel van nieuwe gevoeligheids- en instabiliteit meetwaarden voor veranderingen in de architectuur noodzakelijk geacht.

Voortbouwend op de resultaten van de hierboven genoemde studies, zijn we overgegaan tot de belangrijkste bijdragen van deze thesis, namelijk *het voorstel van methoden (voor berekening van de meetwaarden) en de ontwikkeling van overeenkomstige tools voor het kwantificeren van veranderingsgevoeligheid en instabiliteit in de behoefte-, architectuur- en implementatiefase*. In het bijzonder hebben we methoden gedefinieerd die rekening houden met zowel instabiliteit als veranderings vatbaarheidsstatistieken, het verstrekken van een uniforme meetwaarde die kwantitatief de veranderingsimpact analyse kan leiden. Aanvankelijk richtten we ons op het broncodeniveau, waarvoor we twee maatstaven hebben voorgesteld: de **Ripple Effect Measure** (REM) voor het vastleggen van de instabiliteit van klassenafhankelijkheden, en de **Change Proneness Measure** (CPM) voor het kwantificeren van de vatbaarheid voor klassenwisselingen. Beide meeteenheden zijn empirisch gevalideerd op open-sourcesoftwareprojecten (OSS) door hun beoordelingsvermogen te vergelijken met bestaande meetwaarden. De validatie is uitgevoerd op basis van de 1061-1998 IEEE Standard for Software Quality criteria. De resultaten van de studie suggereerden dat de voorgestelde meetwaarden betere voorspellers van veranderingsgevoeligheid en instabiliteit zijn dan de bestaande. In het bijzonder heeft REM een 38% sterkere correlatie vergeleken met de beste voorspeller van instabiliteit in de literatuur (dat wil zeggen, koppeling tussen objecten), terwijl CPM een 48% sterkere correlatie vertoont in vergelijking met de beste voorspeller van veranderbaarheid in de literatuur (dat wil zeggen, berichtdoorgangskoppeling). Daarnaast bieden we bewijs dat een meetwaarde die zowel veranderbaarheid en instabiliteit combineert een hogere

nauwkeurigheid biedt in vergelijking met het gebruik van de twee factoren afzonderlijk.

Ten slotte hebben we de bovengenoemde meetwaarden als uitgangspunt genomen en deze aangepast aan het architecturale ontwerp- en behoeftepeil. De voorgestelde meetwaarde: *Module Change Proneness Measure* (MCPM) en *Requirement Ripple Effect Metric* (R2EM) zijn empirisch gevalideerd in respectievelijk een OSS- en een industriële context. De evaluatie van beide meetwaarden is net als voorheen zeer positief. Enerzijds is MCPM gemiddeld 23% een nauwkeurigere voorspeller in vergelijking met bestaande pakketstatistieken (bijv. Efferente en Afferente Koppeling). Aan de andere kant heeft R2EM bewezen sterk gecorreleerd te zijn (ca. 60%) met het deskundig oordeel van software engineers.

Voor alle vier bovengenoemde meetwaarden hebben we tools ontwikkeld die hun berekening kunnen automatiseren op basis van bestaande software artefacten, waardoor we de schaal van onze empirische evaluaties kunnen vergroten (waardoor we meer vertrouwen krijgen in de resultaten) en de mogelijke toepasbaarheid van de voorgestelde methoden in de praktijk kunnen vergroten.

Abstract

Along software maintenance the software is subject to changes due to bug fixing, changing requirements, additional requirements, etc. The importance of keeping the cost of maintenance low has been highlighted in the literature with empirical evidence, suggesting that the cost of this phase is approximately 50%-75% of the total cost of software development. This cost can be further increased if the software is characterized by (a) **change proneness**, i.e., the probability of a software artifact to change due to the internal reasons (e.g., fixing bugs in that artifact or changing requirements related to the artifact); and (b) **instability**, i.e., the probability of a software artifact to change due to changes in other artifacts of the system.

The process that investigates change-proneness and instability is called **change impact analysis**; this process is important not just during maintenance but also during the other development phases, e.g. requirements engineering, design, implementation, and testing. Change impact analysis is based on the quantification of change proneness and instability, in order to make decisions on which changes to perform, how and when. However, the literature has identified the following limitations on change impact analysis practices:

- The architectural design and requirements engineering phases are completely lacking metrics that capture artifact instability and change proneness. Therefore, change impact analysis at those phases cannot be quantitative.
- The implementation and detailed-design phases are supported by metrics for instability and change proneness, but such metrics lack accuracy, since they do not combine the aforementioned parameters.
- There is a lack of tools that can automate the process of calculating change proneness and instability metrics for implementation and detailed-design phases.

To tackle the aforementioned limitations, as part of this PhD, *we* have **proposed a set of methods and tools that can quantify both instability and change proneness with increased accuracy at three main development phases** (i.e. requirements engineering, design, and implementation).

As a first step to achieve this goal, we have first reviewed the literature to investigate design-time quality attributes and metrics that can be used to quantify them. Specifically we have reviewed a corpus of more than 150 primary studies. The main findings highlight the importance of change proneness and instability in the current research literature. In particular, stability and change proneness are the most frequently studied quality attributes, after maintainability (15 and 18 studies, respectively). However, the study confirmed the lack of metrics in different development phases (especially in requirements and architecture); the majority of available metrics are at the source-code level.

Based on the main finding of the previous study, as a next step we explored the ability of code metrics to be applicable at artifacts in a different development phase, namely architecture. Thus, we investigated multiple code metrics, with respect to their ability to capture fine- and coarse-grained changes along software maintenance. Intuitively, metrics that are able to capture coarse-grained changes are more fitting to the architecture level; on the contrary fine-grained changes should be neglected at the architecture level. The empirical evaluation suggested that some metrics are more sensitive to changes, and are thus more fitting for method and class level assessments, whereas others are more stable, i.e., they require larger-scale changes for their values to be affected, and are thus more fitting for the architecture level (e.g., packages, components, etc.). Based on the results, the only metric that is related to instability and is able to capture coarse-grained changes is the Response for a Class (RFC) metric, by aggregating class-level metrics to the architecture-level with the use of the average (AVG) function. However, RFC combines coupling (i.e., use of other classes' public interface) and size (i.e., number of local methods) properties, and therefore it cannot be considered as a pure change proneness or instability metric. Thus, the proposal of novel change proneness and instability metrics for architecture is considered necessary.

Building on the results of the aforementioned studies, we proceeded to the main contributions of this thesis, i.e., the *proposal of methods (for metrics calculation) and the development of corresponding tools for quantifying change proneness and instability at the requirements, architecture, and implementation phase*. In particular, we defined methods that consider both instability and

change proneness metrics, providing a unified metric that can quantitatively guide change impact analysis. Initially, we focused on the source-code level, for which we proposed two metrics: the ***Ripple Effect Measure*** (REM) for capturing class dependencies instability, and the ***Change Proneness Measure*** (CPM) for quantifying class change proneness. Both metrics have been empirically validated on open source software (OSS) projects, by comparing their assessment power to existing metrics. The validation has been performed based on the 1061-1998 IEEE Standard for Software Quality Metrics. The results of the study suggested that the proposed metrics are better predictors of change proneness and instability compared to existing ones. In particular REM has 38% stronger correlation compared to the best predictor of instability in the literature (i.e., Coupling Between Objects), whereas CPM shows 48% stronger correlation compared to the best predictor of change proneness in the literature (i.e., Message Passing Coupling). In addition to that, we provide evidence that a metric combining both change proneness and instability offer higher accuracy, compared to using the two factors in isolation.

Finally, using the aforementioned metrics as a starting point, we tailored them to fit the architectural design and requirements level. The proposed metrics: ***Module Change Proneness Measure*** (MCPM) and ***Requirement Ripple Effect Metric*** (R2EM) have been empirically validated in an OSS and an industrial context, respectively. Similarly as before, the evaluation of both metrics has been very positive. On the one hand, MCPM is on average 23% a more accurate predictor compared to existing package metrics (e.g., Efferent and Afferent Coupling). On the other hand, R2EM has proven to be strongly correlated (approx. 60%) to the expert opinion of software engineers.

For all four aforementioned metrics we have developed tools that can automate their calculation from existing software artifacts, enabling us to expand the scale of our empirical evaluations (increasing our confidence of the results), and increasing the possible applicability of the proposed methods in practice.

Table of Contents

Samenvatting	I
Abstract.....	VI
Table of Contents	IX
Acknowledgements.....	XV
Chapter 1 - Introduction	1
1.1 Software Quality Models, Attributes and Metrics	2
1.2 Software Maintainability, Instability, Change Proneness	4
1.3 Research Design	7
1.3.1 Problem Statement	7
1.3.2 Design Science as Research Methodology	8
1.3.3 Practical Problems and Knowledge Questions	9
1.3.4 Using Empiricism to Answer Knowledge Questions	14
1.3.5 Overview of the Dissertation	18
Chapter 2 – Design-Time Quality Attributes and Metrics	21
2.1 Motivation.....	21
2.2 Related Work	26
2.2.1 Domain- or Technology-Agnostic Studies	27
2.2.2 Domain- or Technology-Specific Studies	29
2.2.3 Overview	31
2.3 Study design	33
2.3.1 Objectives and Research Questions	33
2.3.2 Search Process	34
2.3.3 Article Filtering Phases	39
2.3.4 Keywording of Abstracts (Classification Scheme).....	40
2.3.5 Data Collection	40

2.3.6 Data Analysis	42
2.4 Results	43
2.4.1 Design-time Quality Attributes.....	44
2.4.2 Quantification of Quality Attributes through Software Metrics.....	49
2.5. Discussion.....	62
2.5.1 Interpretation of the Results	62
2.5.2 Synthesis and Applicability of the Results.....	66
2.5.3 Implications for Researchers and Practitioners.....	69
2.6 Threats to Validity	71
2.7 Conclusions.....	72
Chapter 3 – Applicability of Metrics on Different Development Phases	73
3.1 Motivation.....	74
3.2 Related Work	77
3.3 Quality Attributes and Object-Oriented Metrics	79
3.4 Software Metrics Fluctuation.....	81
3.5 Case Study on Assessing the Fluctuation of Metrics.....	85
3.5.1 Study Design.....	85
3.5.2 Results	91
3.5.3 Interpretation of Results	99
3.6 Case Study on the Usefulness of SMF in Metrics Selection.....	101
3.6.1 Study Design.....	101
3.6.2 Results	107
3.7 Implications for Researchers and Practitioners.....	109
3.8 Threats to Validity	111
3.9 Conclusions.....	113
Chapter 4 – A Metric for Class Instability	114
4.1 Motivation.....	114
4.2 Related Work	116

4.3 Ripple Effect Measure.....	117
4.4 Validation Process.....	121
4.5 Theoretical Validation	123
4.5.1 Normalization and Non-Negativity.....	123
4.5.2 Null Value and Maximum Value	123
4.5.3 Monotonicity	123
4.5.4 Merging of Classes	124
4.6 Empirical Validation.....	125
4.6.1 Case study Design.....	126
4.6.2 Results	130
4.7 Discussion	134
4.8 Threats to Validity	136
4.9 Conclusions.....	137
Chapter 5 – A Metric for Class Change Proneness.....	139
5.1 Motivation.....	139
5.2 Related Work.....	142
5.3 Proposed Method.....	143
5.4 Case Study Design	146
5.4.1 Metric Validation Criteria	147
5.4.2 Research Objectives and Research Questions.....	147
5.4.3 Case and Units of Analysis.....	148
5.4.4 Data Collection	149
5.4.5 Data Analysis	150
5.5 Results	152
5.5.1 Correlation, Consistency, Tracking, Predictability and Discriminative Power (RQ ₁).....	152
5.5.2 Reliability (RQ ₂).....	155
5.6 Discussion	157

5.6.1 Interpretation of the Results	157
5.6.2 Implications to Researchers and Practitioners	159
5.7 Threads to Validity	159
5.8 Conclusions	160
Chapter 6 – A Metric for Architectural Change Proneness	162
6.1 Motivation.....	162
6.2 Background Information.....	165
6.2.1 Related Work	165
6.2.2 Metric Validation Criteria	166
6.3 Proposed Method.....	166
6.4 Case Study Design	169
6.4.1 Objectives and Research Questions.	170
6.4.2 Case Selection Units of Analysis and Selection	170
6.4.3 Data Collection& Analysis.....	171
6.5 Results	172
6.5.1 Correlation, Consistency, Tracking, Predictability and Discriminative Power (RQ ₁).....	173
6.5.2 Reliability (RQ ₂)	175
6.6 Discussion	177
6.6.1 Interpretation of the Results	177
6.6.2 Implications to Researchers and Practitioners	177
6.7 Threats to Validity	178
6.8 Conclusions.....	180
Chapter 7 – A Metric for Requirements Change Proneness	181
7.1 Motivation.....	181
7.2 Related Work	183
7.2.1 Design and Source Code Change Proneness	184
7.2.2 Our Earlier Work on Change Proneness	185

7.2.3 Change Impact Analysis on Requirements	186
7.2.4 Contributions of this Study	187
7.3 Requirements Change Proneness Metric	187
7.3.1 Proposed Method	187
7.3.2 Illustrative Example	191
7.3.3 Proposed Tool-Chain	195
7.4 Case Study Design	196
7.4.1 Research Questions.....	197
7.4.2 Case Selection.....	197
7.4.3 Data Collection	198
7.4.4 Data Analysis	202
7.5 Results	203
7.5.1 Ripple Effect Factors (RQ ₁).....	205
7.5.2 R2EM Efficiency for Testing Prioritization (RQ ₂).....	210
7.6 Discussion	214
7.6.1 Interpretation of Results	214
7.6.2 Implications for Researchers & Practitioners	215
7.7 Threats to Validity	216
7.7.1 Construct Validity	216
7.7.2 Reliability	217
7.7.3 External Validity.....	217
7.8 Conclusion.....	217
Chapter 8 – Conclusions and Future Work	219
8.1 Answers to Research Questions and Contributions.....	219
8.2 Ongoing and Future Work.....	223
8.2.1 Improvement of Metrics Accuracy	223
8.2.2 Industrial Applicability.....	224
8.2.3 Extending the Scope of the Proposed Metrics	224

8.2.4 Propose Metrics for Other Quality Attributes.....	224
Appendix A	226
Appendix to Chapter 2.....	226
Appendix B	241
Appendix to Chapter 7.....	241
References.....	250
Index	265

Acknowledgements

Conducting a PhD research is a challenging, but at the same time a highly constructive process. Now that this process is almost completed, I need to acknowledge that this thesis is a collective outcome of the forces and support of many people that were involved in this research per se, but also in the underlying process. To these people I would like to express my sincere gratitude.

First of all, I would like to thank my supervisors' team for the research and ethical support that they have provided me throughout this endeavor. In particular, I would like to thank my supervisor, Prof. Paris Avgeriou for giving me the chance to start this project and investigate this very interesting topic. His wide knowledge was precious while guiding my research and finalizing the thesis itself. Additionally, I would like to thank Prof. Alexander Chatzigeorgiou for his trust throughout our time-lasting collaboration. His active participation in the course of this project and his research guidance were really helpful, especially in the first years when he helped me understand aspects of the research and the domain that at that point were unclear to me.

Last, but not least, I could not omit from this section Apostolis Ampatzoglou, who is the person that, among others, 'persuaded' me to start this PhD project. Apostolis guided my research all these years, devoting to me large portion of his valuable time. A big thanks for the endless hours that he 'spent' with me either through Skype calls or face-to-face meetings. Along with Apostoli I have started getting to know what research is, in my BSc thesis, continued to research with him in my MSc thesis, and of course during my PhD. All these years I have learned how to perform research properly, but most importantly through our collaboration, I was continuously motivated, inspired and eventually ended-up to love researching. I would also like to thank him for the help that he has provided me all these years in all aspects of the PhD process: when he was always willing and available to discuss with me anything that might be a problem for me, for his understanding all the difficult periods of the project when I was underperforming, but most importantly for the times when he was believing in me even more than myself.

At the personal level, I big thanks goes to my husband Traianos Plougarli, who gave me the option to pursue my PhD degree and encouraged me to continue my studies all these years. Through his priceless support and selfless attitude,

especially in the difficult times of this PhD journey, he helped me to successfully complete on more step in my academic route. In the course of my PhD I became pregnant and gave birth to a beautiful boy. I could not omit thanking him for the tranquility, strength, and courage that he brought to my life. Every day, he makes me want to become a better person and of course a better researcher. Finally, I would like to thank my family, my father Christo, my mother Dimitra, and my brother Niko for their love and continuous support in all the years of my studies. Their support was of the ultimate importance in both good and bad times, but especially in cases when frustration was my main 'PhD feeling'.

This PhD thesis is devoted to all the people that believed in me all these years and helped my bringing this project to a successful ending.

Thank you all!