

1 INDEPENDENT VARIABLES

Previous work indicates that structural, evolution-based, and code smells are the state-of-the-art features to represent change-prone classes. Following these findings, these code features will be used as the independent variables in this work. The metrics to be used as independent variables are depicted in Tables 1 to 3. Table 1 contains the structural metrics to be used. These metrics are grouped in cohesion, coupling, complexity, inheritance and size.

Table 1: Structural metrics.

Metric	Definition	Source
Cohesion metrics		
LCOM2	Lack of cohesion in methods. The number of pairs of methods in the class using no common attributes minus the number of pairs of methods that do. If this difference is negative, however, LCOM2 is set to zero.	(Chidamber and Kemerer, 1994) [2]
LCOM3	Lack of cohesion in methods. Treats each method pair as an individual entity, and determines the difference between the amount of similar and different pairs.	(Li and Henry, 1993) [7]
TCC	Tight class cohesion. The percentage of pairs of public methods of the class which are directly connected methods. Two methods are called connected, if they use common attributes, directly or indirectly.	(Bieman and Kang, 1995) [1]
LCC	Loose class cohesion. The percentage of pairs of public methods of the class which are directly or indirectly connected. If there are methods m_1, \dots, m_n , such that m_i and m_{i+1} are connected for $i = 1, \dots, n - 1$, then m_1 and m_n are indirectly connected.	(Bieman and Kang, 1995) [1]
Coupling metrics		
CBO	Coupling Between Object Classes. A count of the number of other classes to which it is coupled.	(Chidamber and Kemerer, 1991) [2]
RFC	Response For a Class. A set of methods that can potentially be executed in response to a message received by an object of that class.	(Chidamber and Kemerer, 1991) [2]
FANIN	Fan-in. The number of external classes that invoke methods from the analyzed class.	(Henry and Kafura, 1981) [5]
FANOUT	Fan-out. The number of external method invocations made by the analyzed class	(Henry and Kafura, 1981) [5]
Complexity metrics		
WMC	Weighted Methods Per Class. The sum of the cyclomatic complexity of the methods of a class.	(Chidamber and Kemerer, 1991) [2]
AvgCyclomatic	Average Cyclomatic Complexity. Average cyclomatic complexity for all nested functions or methods.	(McCabe, 1976) [9]
SumCyclomatic	Sum Cyclomatic Complexity. Sum of cyclomatic complexity of all nested functions or methods.	(McCabe, 1976) [9]
MaxCyclomatic	Max. Cyclomatic Complexity. Maximum cyclomatic complexity of all nested functions or methods.	(McCabe, 1976) [9]
Inheritance metrics		
IFANIN	Base Classes. The number of immediate base classes and interfaces.	(Destefanis et al., 2014) [3]
DIT	Depth of Inheritance Tree. The number of nodes between the root of the inheritance tree and the analyzed class.	(Chidamber and Kemerer, 1991) [2]
NOC	Number of Children. Number of immediate sub classes subordinated to a class in the class hierarchy.	(Chidamber and Kemerer, 1991) [2]
OR	Override Ratio. Ratio of methods in a class that are overrides from a superclass.	(Lanza et al., 2006) [6].
Size metrics		
NIM	Instance Methods. The number of instance methods in a class.	(Lorenz and Kidd, 1994) [8]
NIV	Instance Variables. The number of instance variables in a class.	(Lorenz and Kidd, 1994) [8]
LOC	Total lines of code. Count all lines of executable code within the system, class, or method.	(Lorenz and Kidd, 1994) [8]
CLOC	Lines with Comments. The number of lines containing comments in a class.	(Lorenz and Kidd, 1994) [8]
NOPA	Public Fields. Number of Public Fields.	(Lanza et al., 2006) [6]
STMTC	Statements. The number of statements in a class.	(Lorenz and Kidd, 1994) [8]
WOC	Weight of a Class. The number of methods in the interface of the class, divided by the total number of interface members.	(Lanza et al., 2006) [6]
CountSemicolon	Semicolons. The number of semicolons in a class.	(SciTools, 2021) [11]
CountStmtDecl	Declarative Statements. The number of declarative statements in a class.	(SciTools, 2021) [11]
RatioCommentToCode	Comment to Code Ratio. The ratio of comment lines to code lines in a class.	(SciTools, 2021) [11]
CountDeclMethodDefault	Local Default Visibility Methods. The number of local methods with default visibility in a class.	(SciTools, 2021) [11]

Continued on next page

Table 1 – continued from previous page

Metric	Definition	Source
AvgLine	Average Number of Lines. The average number of physical lines between methods of a class.	(SciTools, 2021) [11]
AvgLineComment	Average Number of Lines with Comments. The average number of lines containing comments between methods of a class.	(SciTools, 2021) [11]
CountDeclClassMethod	Class Methods. The number of class methods in a class.	(SciTools, 2021) [11]
CountDeclClassVariable	Class Variables. The number of class variables in a class.	(SciTools, 2021) [11]
NPM	Private Methods. The number of local private methods in a class.	(SciTools, 2021) [11]
CountDeclMethodProtected	Protected Methods. The number of local protected methods in a class.	(SciTools, 2021) [11]
NPRM	Public Methods. The number of local public methods in a class.	(SciTools, 2021) [11]

Table 2 contains the process evolution-based metrics to be used. These metrics quantify the change history of the software between different releases considering the class scope. It considers newly created classes, changes in existing classes, frequency of changes, different weights to changes introduced to a class if it was changed in a far or in a recent past, the number of changes and the percentage of lines of code changed.

Table 2: Software evolution-based metrics (Source: Elish and Al-Khiaty, 2013 [4]).

Metric	Definition
BOC	Birth of a Class. The first time the class appears.
TACH	Total Amount of Changes. It is the sum of added lines, deleted lines, and twice changed lines between release $n - 1$ and release n .
FCH	First Change. The first time the class has been exposed to changes.
LCH	Last Change. The last time the class has been exposed to changes.
CHO	Change Occurred. It is a binary metric that indicates whether or not the class has been exposed to changes from release $n - 1$ to n .
FRCH	Frequency of Changes. The number of times (in term of releases) the class has been changed.
CHD	Change Density. Change density of a class C is its change size ($TACH(C)$) normalized by the size of the class (its total lines of code (LOC)).
WCD	Weighted Change Density. It is a cumulative frequency of change density (CHD) that favor the latest occurrence of changes over the old ones.
WFR	Weighted Frequency of Changes. Is a cumulative frequency of changes that favor the latest occurrence of changes over the old ones.
ATAF	Aggregated Change Size Normalized by Frequency of Change. This is obtained from accumulating size of changes of the class in the past and normalizing by frequency of changes.
LCA	Last Change Amount. It is defined as the last change size of the class when moving from release $i - 1$ to release i .
LCD	Last Change Density. This metric is defined as its last change size (LCA) normalized by the size of the class.
CSB	Changes since the Birth. It is computed by comparing the size of the first version of a class with its current version.
CSBS	Changes since the Birth Normalized by Size. It is the CSB normalized by the size of the first version of the class
ACDF	Aggregated Change Density Frequency. It is obtained from cumulating density of changes introduced to the class in the past, and then this accumulated amount is normalized by the frequency of changes.

We considered sixteen code smell types that are described in Table 3 to be used as change-prone class predictors. The code smells are associated with symptoms of software maintainability problems. These types were chosen because they capture varied characteristics of code involving size, inheritance, complexity, coupling and cohesion at the class and method level.

Table 3: Types of code smells (Source: Rêgo, 2018 [10]).

Metric	Definition
Brain Class	Complex classes that centralize functionalities
Brain Method	Methods that centralize the intelligence of a class
Complex Class	Classes presenting a overly high cyclomatic complexity
Data Class	A class exposes its attributes, thus violating the information hiding principle
Dispersed Coupling	Occurs when a method calls methods from a large number of provider classes
Feature Envy	A method accesses the data of another object more than its own data
God Class	One class monopolizes the processing, and other classes primarily encapsulate data
Intensive Coupling	When a method calls many other methods from a few classes
Lazy Class	Understanding and maintaining classes always costs time and money. So if a class doesn't do enough to earn your attention, it should be deleted.

Continued on next page

Table 3 – continued from previous page

Metric	Definition
Long Method	A method contains too many lines of code. Generally, any method longer than ten lines should make you start asking questions
Long Parameter List	More than three or four parameters for a method
Message Chains	Message chains occur when a client requests another object, that object requests yet another one, and so on. These chains mean that the client is dependent on navigation along the class structure. Any changes in these relationships require modifying the client.
Refused Bequest	Subclass uses only some of the methods and properties inherited from its parents, the hierarchy is off-kilter. The unneeded methods may simply go unused or be redefined and give off exceptions.
Shotgun Surgery	Making any modifications requires that you make many small changes to many different classes.
Spaghetti Code	Declare a number of long methods without parameters
Speculative Generality	There's an unused class, method, field or parameter.

REFERENCES

- [1] James M Bieman and Byung-Kyoo Kang. 1995. Cohesion and reuse in an object-oriented system. *ACM SIGSOFT Software Engineering Notes* 20, SI (1995), 259–262.
- [2] S.R. Chidamber and C.F. Kemerer. 1994. A metrics suite for object oriented design. *IEEE Transactions on Software Engineering* 20, 6 (1994), 476–493. <https://doi.org/10.1007/s11219-020-09525-y>
- [3] Giuseppe Destefanis, Steve Counsell, Giulio Concas, and Roberto Tonelli. 2014. Software metrics in agile software: An empirical study. In *International Conference on Agile Software Development*. Springer, 157–170.
- [4] Mahmoud O Elish and Mojeeb Al-Rahman Al-Khiaty. 2013. A suite of metrics for quantifying historical changes to predict future change-prone classes in object-oriented software. *Journal of Software: Evolution and Process* 25, 5 (2013), 407–437.
- [5] Sallie Henry and Dennis Kafura. 1981. Software structure metrics based on information flow. *IEEE transactions on Software Engineering* 1, 5 (1981), 510–518.
- [6] Michele Lanza and Radu Marinescu. 2010. *Object-Oriented Metrics in Practice: Using Software Metrics to Characterize, Evaluate, and Improve the Design of Object-Oriented Systems*. Springer.
- [7] Wei Li and Sallie Henry. 1993. Object-oriented metrics that predict maintainability. *Journal of Systems and Software* 23, 2 (1993), 111–122. [https://doi.org/10.1016/0164-1212\(93\)90077-B](https://doi.org/10.1016/0164-1212(93)90077-B) Object-Oriented Software.
- [8] M. Lorenz and J. Kidd. 1994. *Object-oriented software metrics - a practical guide*. Prentice-Hall, Inc.
- [9] T.J. McCabe. 1976. A Complexity Measure. *IEEE Transactions on Software Engineering* SE-2, 4 (1976), 308–320. <https://doi.org/10.1109/TSE.1976.233837>
- [10] Diego Cedrim Gomes Rêgo. 2018. *Understanding and Improving Batch Refactoring in Software Systems*. Ph. D. Dissertation. Pontifícia Universidade Católica do Rio de Janeiro, Rio de Janeiro.
- [11] Scitools. 2021. What Metrics Does Undertand Have? <https://support.scitools.com/t/what-metrics-does-undertand-have/66>.