

Exploring Design Level Class Cohesion Metrics

Kuljit Kaur, Hardeep Singh

Department of Computer Science and Engineering, Guru Nanak Dev University, Amritsar, India.
Email: kuljitchahal@yahoo.com

Received November 17th, 2009; revised December 15th, 2009; accepted January 25th, 2010.

ABSTRACT

In object oriented paradigm, cohesion of a class refers to the degree to which members of the class are interrelated. Metrics have been defined to measure cohesiveness of a class both at design and source code levels. In comparison to source code level class cohesion metrics, only a few design level class cohesion metrics have been proposed. Design level class cohesion metrics are based on the assumption that if all the methods of a class have access to similar parameter types then they all process closely related information. A class with a large number of parameter types common in its methods is more cohesive than a class with less number of parameter types common in its methods. In this paper, we review the design level class cohesion metrics with a special focus on metrics which use similarity of parameter types of methods of a class as the basis of its cohesiveness. Basically three metrics fall in this category: Cohesion among Methods of a Class (CAMC), Normalized Hamming Distance (NHD), and Scaled NHD (SNHD). Keeping in mind the anomalies in the definitions of the existing metrics, a variant of the existing metrics is introduced. It is named NHD Modified (NNDM). An automated metric collection tool is used to collect the metric data from an open source software program. The metric data is then subjected to statistical analysis.

Keywords: Design Metrics, Class Cohesion Metrics, Cohesion among Methods of a Class, Normalized Hamming Distance, Scaled NHD

1. Introduction

In the object oriented paradigm, cohesion of a class refers to the degree to which members of the class are interrelated. Chidamber and Kemerer defined the first metric to measure cohesiveness of a class [1]. Since then, several class cohesion metrics have been proposed (discussed in the next section). Empirical studies report that class cohesion metrics are useful to assess software design quality [2,3], to predict fault proneness of classes [4-6], and to identify reusable components [7,8]. Existing class cohesion metrics mainly fall into two categories – metrics which can be computed at design level (high level) and metrics which can be computed one step later *i.e.* at source code level (low level). Design level class cohesion metrics use the limited amount of information available about a class at this level *i.e.* only the class attributes, and method signatures. Method implementation is not completely defined at design level. So some assumptions are made. Different class cohesion metrics defined at design level are based on different assumptions.

1) One school of thought assumes that the types of method parameters match the types of the attributes ac-

cessed by the method. It is further assumed that the set of attribute types accessed by a method is the intersection of this method's parameter types and the set of parameter types of all the methods in the class [9,10].

2) Another school of thought assumes that the set of attribute types accessed by a method is the intersection of the set of this method's parameter types and the set of its class attribute types [11].

In this paper we review the design level class cohesion metrics based on the first assumption. Keeping in mind the anomalies in the definitions of the existing metrics, a variant of the metrics is introduced. The paper is organized as follows: Section 2 reviews the related work. Section 3 explains the existing design level class cohesion metrics and introduces a modified version as well. Section 4 presents the statistical analysis of the data collected from an open source project. Section 5 concludes the paper.

2. Related Work

A number of class cohesion metrics are defined in the low level metrics category [1,12-26]. However, there are only a few proposals for design level class cohesion metrics [9-11].

The metric, named Cohesion among Methods of a Class (CAMC) captures the information about parameter types of methods of a class [9]. A class is cohesive if all methods of the class use the same set of parameter types. Methods which use same type of parameter types are assumed to process related kind of information. CAMC metric values lie in the range [0, 1]. Counsell *et al.* point out some anomalies in definition of this metric, and propose a new metric named Normalized Hamming Distance (NHD) [10]. It is a normalized metric which measures average agreement between each pair of methods on their parameter types. A variant of the NHD metric called Scaled NHD (SNHD) is introduced in the same paper. It addresses shortcomings of both CAMC and NHD, as claimed by the authors [10]. This research finds anomalies in the definitions of NHD and SNHD as well, and proposes a modified version of the NHD metric - NHD modified (NHDM). The NHDM metric gives statistically significant results.

Dallal proposes another metric for measuring cohesion of a class at design level [11]. Similarity based Class Cohesion (SCC) metric is based on the second assumption discussed above. This metric is not analyzed in this paper as the automated tool developed for this research does not support collection of this metric.

3. Design Metrics

This section describes the class cohesion metrics computable with information available at design level. At design level, information regarding name of the class, its attributes (names, and data types), and method signatures is available. Method signature includes name of the method and its parameter list which describes names of the parameters and their data types. A Class does not have a detailed or algorithmic description of its methods available at this level.

3.1 CAMC

The CAMC metric measures the extent of intersection of individual method parameter type lists with the parameter type list of all methods in the class [9]. This metric computes the relatedness among methods of a class based upon the parameter list of the methods. It is assumed that methods of a class, having access to similar parameter types, process closely related information.

The CAMC metric uses a parameter-occurrence matrix (PO matrix) that has a row for each method and a column for each data type that appears at least once as the type of a parameter in at least one method in the class. The value in row i and column j in the matrix is 1 when the i th method has a parameter of the j th data type and is 0 otherwise. In the original version of the metric [9], the PO matrix has an additional column of all 1s. This column represents the 'self' parameter that corresponds to the type of the class itself which is by default one of the pa-

rameters of every method. In this discussion, the original version of the metric is referred to as CAMC_s (Cohesion among methods of a class with 'self' parameter) and metric definition without the 'self' parameter is named as CAMC [10].

The CAMC metric is defined as the ratio of the total number of 1s in the PO matrix to the total size of the matrix.

$$\text{CAMC}(C) = \frac{\sigma}{kl} \quad \text{where } \sigma = \sum_{i=1}^k \sum_{j=1}^l PO[i][j]$$

CAMC suffers from the following anomalies:

1) CAMC gives false positives – the metric gives a non-zero value for a class with no parameter sharing in its methods.

2) CAMC can not differentiate between two classes having same number of 1s but with different patterns of 1s in their PO matrices.

3) Smaller classes take high values for the cohesion metric than the larger classes with same properties.

3.2 NHD

Counsell *et al.* [10] suggested an alternative of CAMC. It is based on the definition of hamming distance. NHD measures agreement between rows in the PO matrix. NHD metric for a class with k methods and l unique parameter types (set obtained from union of parameter types received by all its methods) is defined as:

$$\text{NHD} = \frac{2}{lk(k-1)} \sum_{i=1}^{k-1} \sum_{j+1}^k a(i,j)$$

where $a(i,j)$ is value of the cell at (i,j) th location in the PO matrix. Another easy way to compute NHD is to first find the sum of disagreements between methods for all the parameter types and then subtract it from 1.

$$\text{NHD} = 1 - \frac{2}{lk(k-1)} \sum_{j=1}^l c_j(k-c_j)$$

where c_j is the number of 1s in the j th column of the PO matrix.

A variant of NHD (with self parameter), NHD_s can be defined for a PO matrix with an additional column of all 1s.

NHD suffers from the following anomalies:

1) NHD metric also gives false positives. The metric removes the first anomaly of the CAMC for a class with $k = l = 2$. The metrics fails to give correct answer for higher values of k and l (e.g. when $k = l = 3$, and there is no parameter sharing among methods, NHD metric gives a non-zero value).

2) NHD does not give different answers for classes with different properties – metric fails to distinguish a class with no parameter sharing in its methods from a class with substantial amount of parameter sharing in its methods.

3) Class size influences metric value. As size of the class increases, value of the NHD metric also increases (even if the PO matrix gets sparser).

3.3 SNHD

SNHD is the Scaled NHD metric proposed to interpret values of the NHD metric in a more varied range. Proponents of the NHD metric are of the opinion that NHD metric can take values at two extremes: the minimum or the maximum. But they admit that it is not clear as to which of these extremes represents a cohesive class. However without giving any clear explanation they state that classes at both the extremes may be cohesive. They define these extreme values as NHD_{min} and NHD_{max} respectively [10]. SNHD metric value helps to know how close the NHD metric is to the maximum value of the NHD value in comparison to the minimum value. SNHD is defined as follows:

$$SNHD = \begin{cases} 0 & \text{if } NHD_{min} = NHD_{max} \text{ and } \sigma < kl, \\ 1 & \text{if } \sigma = kl, \\ 2\left(\frac{NHD - NHD_{min}}{NHD_{max} - NHD_{min}}\right) - 1, & \text{otherwise} \end{cases}$$

The SNHD metric values lies in the range [-1,1]. SNHD = -1 implies that $NHD = NHD_{min}$, and SNHD = 1 implies that $NHD = NHD_{max}$. NHD is closer to its minimum or maximum value depending upon whether SNHD is getting values close to -1 or +1 respectively. A class is considered non-cohesive if SNHD metric value for the class is 0.

SNHD_s is defined by considering the ‘self’ parameter. SNHD suffers from these Anomalies:

- 1) Difficult to calculate and interpret.
- 2) False negatives – SNHD metric gives 0 value for a class with good degree of cohesion.

3.4 NHDM

Keeping in view the anomalies of the cohesion metrics discussed above, this research proposes a variation of the NHD metric. This variat is named as Normalized Hamming Distance Modified (NHDM) metric. The NHD metric ignores the method pairs with zero values in a column of the PO matrix. It counts only those methods pairs which do not agree, and ignores all other method pairs irrespective of whether they agree on a 0 or a 1. NHDM counts the method pairs which agree on a 0, as a disagreement. NHDM for a class with k methods and l unique parameter types, of all its methods, is defined as:

$$NHDM = 1 - \frac{2}{lk(k-1)} \sum_i^l (c_j(k - c_j) + \frac{1}{2}z_j(z_j - 1))$$

where c_j is the number of ones and z_j is the number of zeroes in the j th column of the PO matrix for the class.

Similarly NHDMs is defined by including the ‘self’ parameter in the PO matrix.

This metric removes the anomalies present in the definition of CAMC, NHD, and SNHD metrics. NHDM

gives correct results. It gives different results for classes with different properties. NHDM metric values are independent of the class size.

4. Data Analysis

Cohesion metrics discussed above are collected from an open source software system available at www.sourceforge.net. The software is a JAVA based charting library, and it consists of 884 classes. For automated collection of metrics, a tool CohMetric is developed.

4.1 Descriptive Analysis

Histograms in **Figures 1 to 4** show metrics distributions. **Table 1** presents the descriptive statistics. It can be observed that majority of the CAMC metric values lie close to 0 (see **Figure 1**). On average a class’s cohesion value is 0.21. NHD metric takes values in a higher range (**Figure 2**). Average NHD metric value is 0.66. SNHD is 0 for maximum of the classes. Its values lie more on the

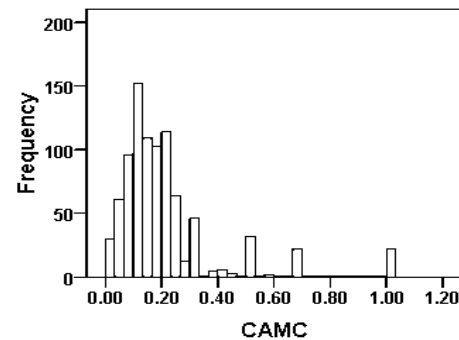


Figure 1. Distribution of CAMC metric

Table 1. Descriptive statistics for cohesion metrics

Metric	Average	Std Dev	Metric	Average	Std Dev
CAMC	0.21	0.18	CAMC _s	0.48	0.21
NHD	0.66	0.21	NHD _s	0.81	0.12
SNHD	-0.43	0.51	SNHD _s	0.63	0.42
NHDM	0.05	0.16	NHDM _s	0.38	0.22

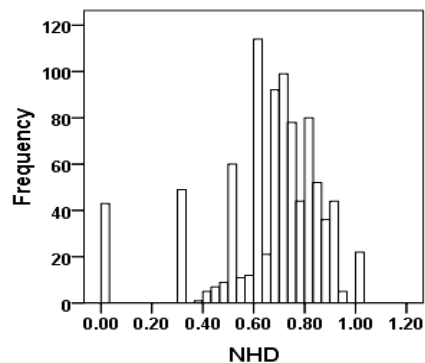


Figure 2. Distribution of NHD metric

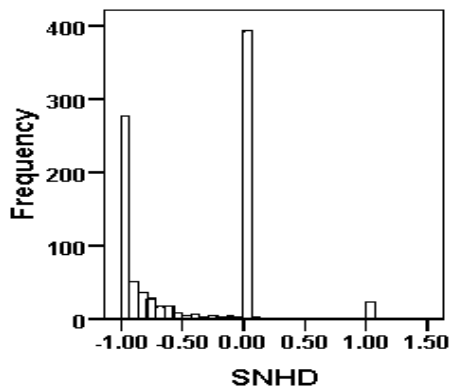


Figure 3. Distribution of SNHD metric

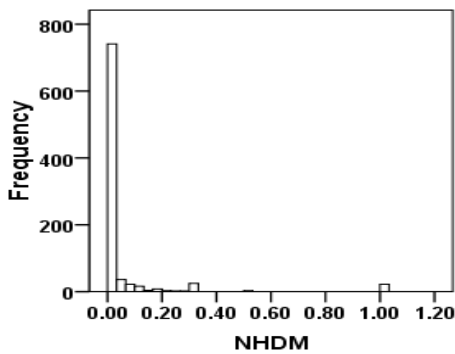


Figure 4. Distribution of NHDM metric

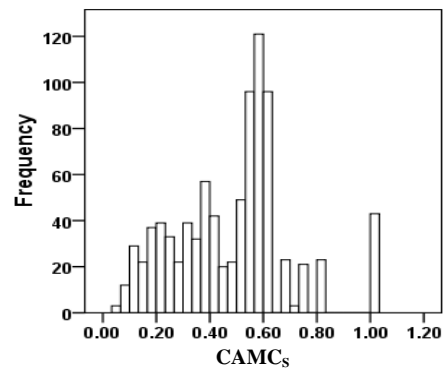


Figure 5. Distribution of CAMCs metric

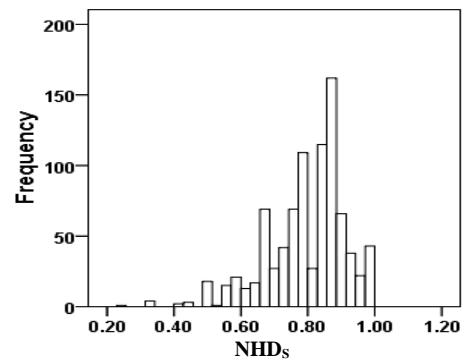


Figure 6. Distribution of NHDs metric

left side of 0 which implies that majority of the classes has NHD more close to NHD_{min} than NHD_{max} . Average SNHD for a class is -0.43 and standard deviation is also very high (Figure 3). NHDM takes very low values (Figure 4). For majority of the classes it is 0. Its average value is just 0.05. As earlier stated, it may be due to the reason that it does not give false positives.

4.2 Metric Variants

Variants of these cohesion metrics are defined on the basis of the assumption that all the methods of a class by default receive the class type itself (self) as one of the parameter types. $CAMCs_s$, NHD_s , $SNHD_s$ and $NHDM_s$ are defined as variants of CAMC, NHD, SNHD, and NHDM respectively. Cohesion metrics which consider the 'self' parameter are expected to give higher values as the class methods agree on at least one parameter type. Table 1 gives a comparison of averages of cohesion metrics and their variants. All the metrics in this category (which consider self parameter type) have higher averages than their counterparts. The observation is that metric variants, which consider 'self' as one of the parameter types, take values in higher range. It is also confirmed by the descriptive analysis of these metrics as shown in Figures 5 to 8. It is worth noting that $SNHD_s$ takes values in the range from 0 to 1 more frequently, in contrast to SNHD which takes values in the range from 0 to -1 . It

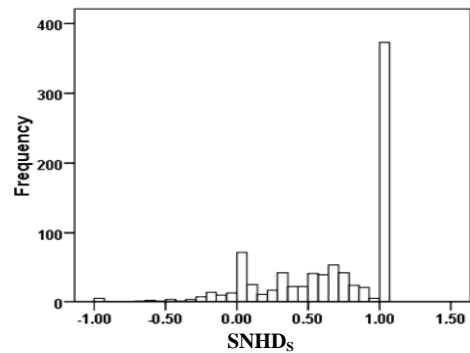


Figure 7. Distribution of SNHDs metric

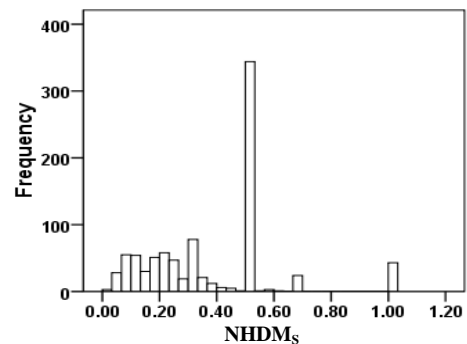


Figure 8. Distribution of NHDMs metric

implies that a class whose NHD value is more close to NHD_{max} is more cohesive.

4.3 Size Independence

Figures 9 to 12 present the relation between cohesion metrics and class size (measured in terms of number of methods). CAMC metric value is higher for small classes and is lower for large classes. NHD takes large values for classes with larger number of methods. This is in line with the earlier findings about these two metrics [10]. As shown in Figure 11, SNHD is close to 1 for

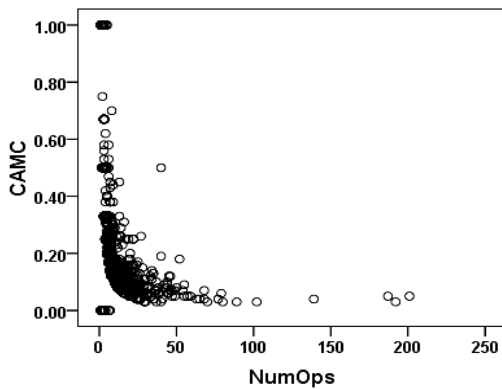


Figure 9. Scatter diagram of CAMC and class size

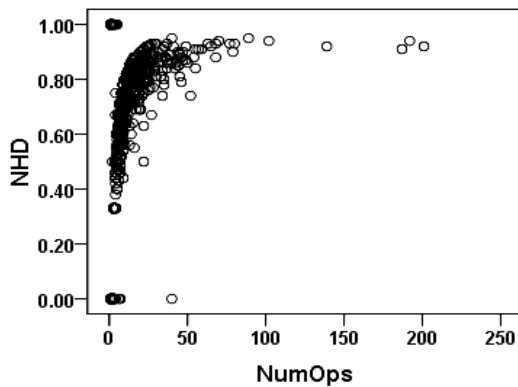


Figure 10. Scatter diagram of NHD and class size

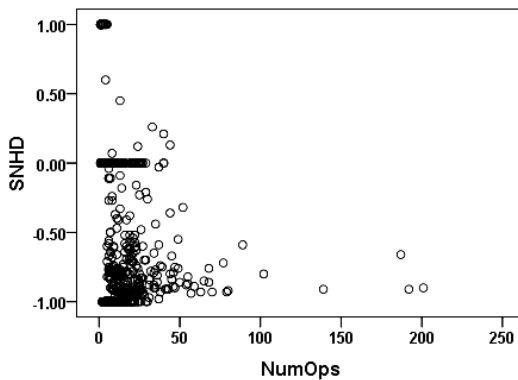


Figure 11. Scatter diagram of SNHD and class size

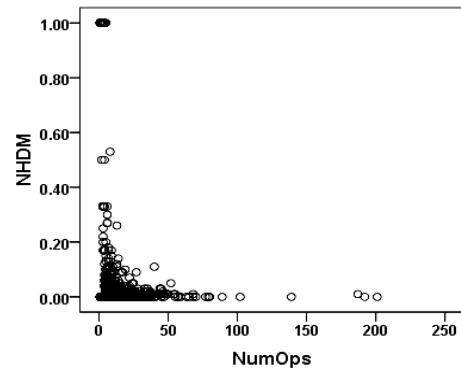


Figure 12. Scatter diagram of NHDM and class size

some comparatively small classes. For larger classes, SNHD lies in the range $[-1, 0]$. NHDM takes values near 0 for most of the classes. However small classes have metric value in the higher range. However if size of the parameter occurrence (PO) matrix is taken into consideration then it is found that it does not have significant correlation with any of the metrics (see Table 2). Here l represents the number of parameter types, k is the number of methods of the class, and lk is the size of the parameter occurrence matrix. This result is unlike the previous studies on these metrics [10,27].

4.4 Metrics Inter-Dependencies

The parametric Pearson’s correlation coefficient between each pair of cohesion metrics is given in Table 3. All the correlation figures are significant at $p = 0.01$ level. Metric variants such as $CAMC_s$, NHD_s , $SNHD_s$, and $NHDM_s$ are moderately correlated with their counterparts. NHD and NHD_s have the highest correlation coefficient in this category. $NHDM$ and $CAMC$ are strongly correlated. Similar is the case for their variants $NHDM_s$ and $CAMC_s$. $SNHD$ is moderately correlated with $NHDM_s$ and $CAMC_s$. Unlike the previous studies, the correlation analysis for this data set does not show any significant correlation in NHD and CAMC [10,27]. However the scatter plot of values for these two metrics shows a negative trend. CAMC and NHD show a negative relationship in the scatter diagram given in Figure 13. CAMC is very low for the classes for which NHD is very high. On average the NHD metric takes values in higher range. This implies that this metric pair does not have a linear covariation.

Principal Component Analysis (PCA) is used to identify the metrics measuring orthogonal dimensions. Rotated principal components are obtained using the varimax rotation technique. Three principal components are extracted which capture 93.28% of the data set variance (shown in Table 4). Metrics with significant loading coefficients in a particular dimension are highlighted in bold. An analysis of the table shows that $NHDM_s$ and $CAMC_s$ and $SNHD$ contribute significantly to the first

Table 2. Correlation in cohesion metrics and size

	CAMC	CAMC _s	NHD	NHD _s	SNHD	SNHD _s	NHDM	NHDM _s
l	-.222	-.696	.337	.003	-.356	-.539	-.128	-.645
k	-.307	-.520	.350	.219	-.071	-.179	-.177	-.429
lk	-.158	-.357	.210	.138	.067	-.223	-.075	-.298

Table 3. Correlation analysis among metrics

	CAMC	CAMC _s	NHD	NHD _s	SNHD	SNHD _s	NHDM	NHDM _s
CAMC _s	0.575							
NHD	-0.267	-0.542						
NHD _s	-0.372	-0.024	0.623					
SNHD	0.341	0.654	-0.043	0.334				
SNHD _s	-0.253	0.347	0.271	0.678	0.478			
NHDM	0.854	0.466	0.122	0.107	0.403	-0.043		
NHDM _s	0.456	0.962	-0.356	0.249	0.726	0.520	0.480	

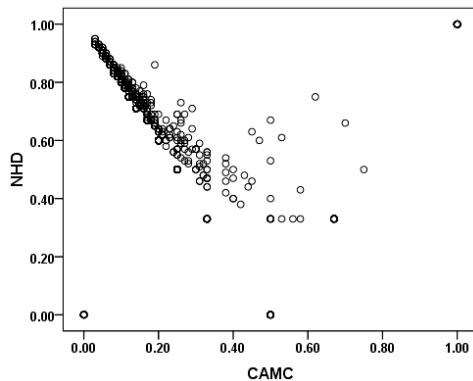


Figure 13. Scatter diagram shows correlation in CAMC and NHD metrics

Table 4. Principal Components Matrix

	PC1	PC2	PC3
Eigen Value	3.72	2.39	1.35
Percent	46.45	29.93	16.90
Comm. percent	46.45	76.38	93.28
CAMC	0.25	-0.32	0.90
CAMC _s	0.91	-0.27	0.30
NHD	-0.39	0.89	0.11
NHD _s	0.27	0.90	-0.13
SNHD	0.84	0.21	0.27
SNHD _s	0.66	0.59	-0.31
NHDM	0.24	0.15	0.95
NHDM _s	0.95	-0.01	0.25

dimension: PC1. SNHD_s is moderately significant in two dimensions: PC1 and PC2. NHD and NHD_s both load significantly on PC2. NHDM and CAMC both load significantly on PC3.

It is worth mentioning here that NHDM and NHDM_s have the maximum variance among all the metrics in this analysis. So metrics measuring different dimensions are:

- PC1: NHDM_s, CAMC_s
- PC2: NHD, NHD_s
- PC3: NHDM, CAMC

5. Conclusions

Cohesion is one of the important design properties to realize a quality software product. Many empirical studies exist which relate the cohesion design property with other properties of interest such as maintainability, reusability, and reliability. Several metrics have been proposed to compute cohesion at class level in object oriented systems. In this paper design level cohesion metrics such as CAMC, NHD, SNHD have been investigated using empirical data. In view of the anomalies present in the existing metrics' definitions, a modified version of the NHD metric is proposed and is named as NHDM (NHD Modified) ss. Statistical analysis of the metrics data shows that CAMC and NHD are influenced by the size of class (measured in terms of number of methods). None of the studied metrics correlates with the size of the Parameter Occurrence matrix (PO matrix) of the class. Principal Component Analysis of the data shows that NHDM and CAMC both give similar results but NHDM has more variation in its values. Similar is the case for NHDM_s and CAMC_s. SNHD or SNHD_s does not contribute significantly to any dimension. NHD and NHD_s are not significantly related to any of the other metrics.

REFERENCES

- [1] P. Chidamber and C. Kemerer, "Towards a Metrics Suite for Object Oriented Design," *Proceedings of 6th ACM Conference on Object Oriented Programming, Systems, Languages and Applications*, Phoenix, Arizona, 1991, pp. 197-211.
- [2] L. Briand, J. Wust, J. Daly and D. Porter, "Exploring the Relationships between Design Measures and Software Quality in Object Oriented Systems," *Journal of Systems and Software*, Vol. 51, No. 3, 2000, pp. 245-273.
- [3] J. Bansiya and C. Davis, "A Hierarchical Model for Object Oriented Quality Assessment," *IEEE Transactions on Software Engineering*, Vol. 28, No. 1, 2002, pp. 4-17.
- [4] T. Gyimothy, R. Ferenc and I. Siket, "Empirical Validation of Object-Oriented Metrics on Open Source Software for Fault Prediction," *IEEE Transactions on Software Engineering*, Vol. 31, No. 10, 2005, pp. 897-910.
- [5] Z. Zhou and H. Leung, "Empirical Analysis of Object-Oriented Design Metrics for Predicting High and Low

- Severity Faults,” *IEEE Transactions on Software Engineering*, Vol. 32, No. 10, 2006, pp. 771-789.
- [6] M. Marcus and D. Poshyvanyk, “Using the Conceptual Cohesion of Classes for Fault Prediction in Object-Oriented System,” *IEEE Transactions on Software Engineering*, Vol. 34, No. 2, 2008.
- [7] J. Lee, S. Jung, S. Kim, W. Jang and D. Ham, “Component Identification Method with Coupling and Cohesion,” *Proceedings of the Eighth Asia-Pacific Software Engineering Conference*, December 2001, pp. 79-86.
- [8] G. Gui and D. Scott, “Measuring Software Component Reusability by Coupling and Cohesion Metrics,” *Journal of Computers*, Vol. 4, No 9, Academy Publishers, 2009, pp. 797-805.
- [9] J. Bansiya, L. Etzkorn, C. Davis and W. Li, “A Class Cohesion Metric for Object Oriented Designs,” *Journal of Object Oriented Programming*, Vol. 11, No. 8, 1999, pp. 47-52.
- [10] S. Counsell, S. Swift and J. Crampton, “The Interpretation and Utility of Three Cohesion Metrics for Object-Oriented Design,” *ACM Transactions on Software Engineering and Methodology*, Vol. 15, No. 2, 2006, pp. 123-149.
- [11] J. Dallal, “A Design-Based Cohesion Metric for Object-Oriented Classes,” *Proceedings of the International Conference on Computer and Information Science and Engineering*, 2007, pp. 301-306.
- [12] W. Li and S. Henry, “Object-Oriented Metrics that Predict Maintainability,” *Journal of Systems and Software*, Vol. 23, No. 2, 1993, pp. 111-122.
- [13] S. Chidamber and C. Kemerer, “A Metrics Suite for Object Oriented Design,” *IEEE Transactions on Software Engineering*, Vol. 20, 1994, pp. 476-493.
- [14] M. Hitz and B. Montazeri, “Measuring Coupling and Cohesion in Object-Oriented Systems,” *Proceedings of International Symposium on Applied Corporate Computing*, 1995.
- [15] J. Bieman and B. Kang, “Cohesion and Reuse in an Object-Oriented System,” *Proceedings of the 1995 Symposium on Software Reusability*, ACM Press, 1995, pp. 259-262.
- [16] B. Henderson-Sellers, L. Constantine and I. Graham, “Coupling and Cohesion (towards a Valid Metrics Suite for Object-Oriented Analysis and Design),” *Object Oriented Systems*, Vol. 3, 1996, pp. 143-158.
- [17] L. Briand, J. Daly and J. Wust, “A Unified Framework for Cohesion Measurement in Object-Oriented Systems,” *Empirical Software Engineering*, Vol. 3, No. 1, 1998, pp. 65-117.
- [18] H. Chae, Y. Kwon and D. Bae, “A Cohesion Measure for Object-Oriented Classes,” *Software Practice and Experience*, Vol. 30, No. 12, 2000, pp. 1405-1431.
- [19] Z. Chen, Y. Zhou and B. Xu, “A Novel Approach to Measuring Class Cohesion Based on Dependence Analysis,” *Proceedings of the International Conference on Software Maintenance*, 2002, pp. 377-384.
- [20] L. Badri and M. Badri, “A Proposal of a New Class Cohesion Criterion: An Empirical Study,” *Journal of Object Technology*, Vol. 3, No. 4, 2004.
- [21] J. Wang, Y. Zhou, L. Wen, Y. Chen, H. Lu and B. Xu, “DMC: A More Precise Cohesion Measure for Classes,” *Information and Software Technology*, Vol. 47, No. 3, pp. 176-180, 2005.
- [22] C. Bonja and E. Kidanmariam, “Metrics for Class Cohesion and Similarity between Methods,” *Proceedings of the 44th Annual Southeast Regional Conference*, ACM Press, New York, 2006, pp. 91-95.
- [23] G. Cox, , L. Etzkorn and W. Hughes, “Cohesion Metric for Object-Oriented Systems Based on Semantic Closeness from Disambiguity,” *Applied Artificial Intelligence*, Vol 20, No. 5, 2006, pp. 419-436.
- [24] L. Fernández and R. Peña, “A Sensitive Metric of Class Cohesion,” *International Journal of Information Theories and Applications*, Vol. 13, No. 1, 2006, pp. 82-91.
- [25] S. Makela and V. Leppanen, “Client Based Object Oriented Cohesion Metrics,” *31st Annual International Computer Software and Applications Conference*, Vol. 2, 2007, pp. 743-748.
- [26] A. Marcus and D. Poshyvanyk, “The Conceptual Cohesion of Classes,” *Proceedings of 21st IEEE International Conference on Software Maintenance*, 2005, pp. 133-142.
- [27] J. Dallal and L. Briand, “An Object-Oriented High-Level Design-Based Class Cohesion Metric,” Simula Technical Report (2009-1), Version 2, Simula Research Laboratory, 2009.