# Correlation analysis of modeling approach comparison criteria: methodology proposal

Frédéric Mayart[a] and Jean-Michel Bruel[b] and Brahim Hamid[b]

[a] Psychology M.D. and independant consultant

frederic.mayart@gmail.com

[b] University of Toulouse, France

{bruel|hamid}@irit.fr

July 27, 2013

**Abstract**

The CMA@MODELS workshop aims at defining a set of comparison criteria that provide a basis for understanding, analyzing, and comparing various modeling approaches. While the definition of these criteria has evolved since the Bellairs AOM[1] workshop in April 2011, they have not been studied yet in terms of dependencies and correlation. The aim of this paper is to propose a study that will be conducted by the time of the CMA@MODELS'2013 workshop and that assesses the level of dependence between those factors and hopefully the way they can be improved accordingly.

## 1  Introduction

One possible final goal of defining a set of criteria to define modeling approaches [?] is to help people, especially from industry, picking up the good approaches or artifacts according to their own purpose. For example in the case of the definition of a particular DSL (Domain Specific Language) one will select paradigm A and complement it with approach B using the comparison criteria as a basis to guide his/her selection.

The authors of the comparison criteria have, through different workshops, managed to get several different assessments made by defenders of particular modeling approaches. The assessments, and more precisely the feedback on building these assessments, have been useful to the users to improve and precise the comparison criteria.

From our point of view the experiment is mature enough to support an analysis of the criteria themselves in the sense that we should be able to analyse if, for example, some of the criteria are not redundant with respect to another one, or if some criteria are always impacted by another, etc.

The goal of this paper is to suggest the kind of analysis that could be conducted on the criteria, and to provide arguments that such analysis could be use-

---

[1] http://www.cs.mcgill.ca/joerg/SEL/AOMBellairs2011.html

ful. In section 2 we will briefly provide an overview of the comparison criteria, their organisation and the purchased goal, in order to make the paper readable outside the scope of the CMA specialists. In section 3 we will discuss potential analysis techniques and we will go into details into the ones that we recommend. In section 4 we will provide an overview of how the study should be conducted and the kind of expected results we hope can be obtained. We will additionally provide some advices. We will conclude in section 5.

## 2    Criteria

The comparison criteria, which definition started at the Bellairs AOM workshop in April 2011 [**?**] aim to provide a basis for understanding, analyzing, and comparing various modeling approaches. They are organized in two groups:

- criteria related to general modeling dimensions, and

- those related to key modeling concepts.

The modeling dimensions criteria characterize a modeling approach by (i) the development phases and activities for which it is most applicable and (ii) the languages or notations used. The key modeling concepts criteria are used to classify a modeling approach in terms of the modeling building blocks, attributes, or qualities targeted for optimization or improvement by the approach. These criteria are intended to play a particularly important role in situating existing approaches in the current body of work, and in identifying the considerations that must be made when developing new modeling approaches. In order to fulfil their role, it is important that the criteria be defined in terms that are widely and precisely understood by the modeling community. We

have identified the following as central key modeling concepts: modularity, first class entities including units of encapsulation, and composability including composition rules and composition operators. The key modeling concepts criteria can be applied to either an approach or to the models that are produced using the approach. The focus of this document, however, is the assessment of modeling approaches, and not the models they produce. Furthermore, key modeling concepts can conceivably be interpreted differently depending on which modeling dimension is being considered. While we recognize this fact, we present the modeling dimensions and key modeling concepts in a more orthogonal view.

In Fig. 1 we provide an extract of the structures modeling dimensions and key modeling concepts.

## 3    Potential    analysis techniques

In this kind of exercise where a set of criteria is aimed at describing a particular topic (in our case modeling approaches) for a certain purpose (such as comparison, classification, artefact retrieving, . . . ) it is recommended to conduct a correlation analysis.

If such purpose is to find whether *correlations exist among comparison criteria*, multivariate analysis methods like PCA (Principal Component Analysis) come to mind [**?**]. The idea is to explore how different subsets of methods would correlate while not doing so with other subsets due to specific comparison criteria, which is also a pre-requisite to cluster analysis. Finding underlying factors (vectors) that would best explain which variables aggregate and how subsets separate is just a first step.

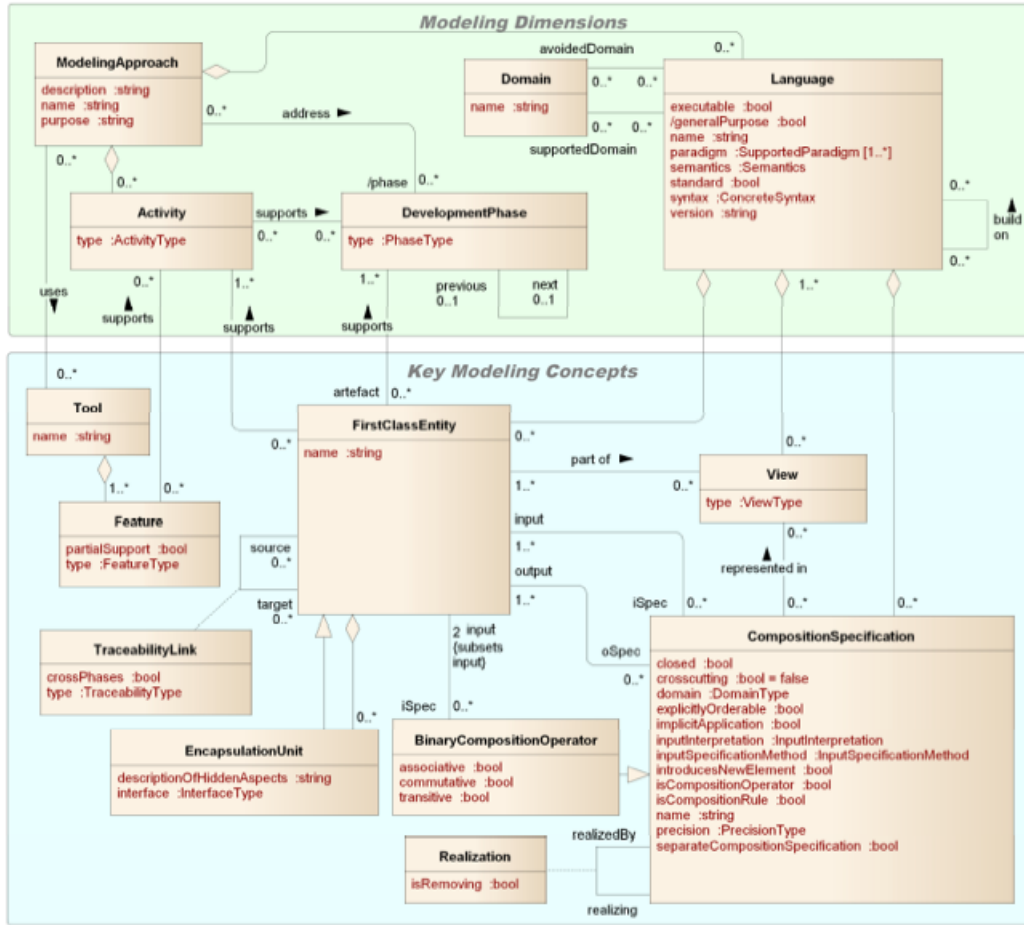An even better purpose is to try to find out whether comparison criteria or

Figure 1: Extract of the metamodel for modeling approaches (from [**?**])

subsets combine with other subsets, at different levels, in order to form "families" or types of methods. Do potential subsets of comparison criteria make some methods more similar to each other inside a given cluster than those determined by other subsets? Cluster analysis might help identify such combinations and thus describe which comparison criteria tend to be tied together and which tend to separate from each other. Clustering, with classification methods, could also be a way to see whether there are *comparison criteria that are essentially repeated or superfluous*, as well as whether there is *a subset of all the comparison criteria that is sufficient to differentiate between the modeling ap-*proaches.

# 4 Proposed study

Multidimensional Statistics use two main methods: Factorial Analysis methods that consist in projecting a cloud of points on a vector space, while loosing as little information as possible; and Classification Methods, that try to cluster those points.

Factorial Analysis methods regroup three main techniques: Principal Component Analysis (PCA, with several quantitative variables), Correspondence Analysis (CA, two quantitative variables, represented by a contingency table) and Multiple Correspondence Anal-

3

ysis (MCA, more than two variables, all qualitative).

In Statistics, variables are nominal (like gender: Male or Female) or numeric, i.e. qualitative or quantitative. Numeric variables can be divided in ordinal variables where the number just represents a position in a scale (rank on a preference scale, level of knowledge), and "true" numeric variables (like Weight, Age, or any measure).

While nominal variables support simple independence tests as we only know their cardinal values (occurrences) to compare them, numeric variables can be computed to assess direct correlations (on ranks or on their values).

In the proposed metamodel, variables describing methods(comparison criteria and keys) are of different types. That's why, even if it follows the same philosophy as the usual PCA when doing a factorial analysis, MCA (Multiple Correspondence Analysis) [?] should be preferred. Quantitative variables can always be chunked down into qualitative ones (classes like low, medium or high value for example), while the reverse is not possible. Once all variables are converted in qualitative ones and put into a disjunctive table (like contingency table for Simple Correspondence Analysis), the analysis can be conducted.

So, it is possible to do a factorial analysis with a mix of qualitative and quantitative data, though sometimes it can lead to weird results because there is more or less some kind of *loss of information*.

## 4.1 Objectives

We propose to study instances (methods), variables (comparison criteria) and their modalities [?].

### 4.1.1 Methods

Two methods are close if questionnaires have been answered the same way. The focus will not be on instances (methods) *per se* but more on sets: *are there groups of methods?* This is a common approach in social studies where we are not much interested in individuals themselves but groups of individuals inside the population. And here we have individuals responding and giving information about the criteria they follow when using their method. So, doing a correlation analysis on methods is "more or less" the same as analyzing individuals and observe how they regroup and under which factors.

### 4.1.2 Comparison criteria and their modalities

First of all, we want to study the relationships between variables and the associations between modalities. A modality is the "value" a variable can take. Qualitative variables and quantitative ordinal variables have discrete values, and usually a finite set of modalities (ex. Gender, two modalities: M or F). Two modalities are close if they have been taken together often. That means two comparison criteria characteristics (from different criteria) that are often cited together by a group of individuals (and thus from different methods) will geometrically appear close on the plot graphs generated by the factor analysis. We are looking for such plots clouds.

Second of all, one or more synthetic continuous variable(s) are to be looked for to resume the qualitative variables. Here is when a statistician would use the principal components to interpret the relations between variables. Using a representation by modalities is easier to show how vector dimensions separate or aggregate the different criteria and gives more precision.

Third of all, the goal is to characterize instance subsets (groups of indi-

Figure 2: Actual assessment form (spreadsheet)

viduals/methods) by modalities of comparison criteria. A Hierarchical Cluster Analysis is the logic and common follow-up to a MCA. We want to regroup the methods in a few number of classes corresponding to "profiles" of Comparison Criteria. The result is a hierarchical tree easy to interpret. Methods would appear as leaves, clustering into small branches, then bigger ones, etc., up to a a trunk. Classes can then be described by the criteria variables and/or their modalities, by the factorial dimensions, or the individuals/methods.

Free solutions exist to make such complex analyses, like the now famous statistical software $R$, with a large support community and plenty of specific packages being developped, a bit *à la* LaTeX.

## 4.2 Possible improvements

We can imagine that modernizing the way people fill the form can be useful. In Figure 2 we illustrate the actual form (Excel spreadsheet) in comparison with a potential Google Form (see Figure 3).

The advantage is the default results available through the form (charts, etc.) and that is always possible to generate an Excel spreadsheet. For the data mining involved in such studies it greatly reduces the time necessary for encoding the variables and automatically filling the result tables to be used by a Statistical Software, while avoiding potential human errors.

Another improvement can be done on the questionnaire itself. The possible loss of information due to the mix of variables and therefore the need to convert quantitative variables into qualitative ones, for the appropriate Factor Analysis, can be minimized. Some variables appear to be disjunctive in their modalities, with only two mutually exclusive "values" (either/or). Others may have a middle-ground modality. Even more, and whenever it is possible, preference scales (with 5 to 7 ranks) should be used as long as they pragmatically reflect a position by the user between two extreme points on a given Comparison Criteria (*formal* vs *informal* language). Or an intensity rank for single variable modalities with no opposite, meaning "more or less", scaling from "never" to "always" or "none" to "a lot", etc. (like *Aspect Oriented*).

Open questions imply the questionnaire should be processed twice. A simple content analysis on already gathered data (like counting the most common occurrences of keywords, the ones that cover around 80% of the total as in a Pareto diagram for example) would lead to new multiple choice questions with collected keywords as new variable modalities.

Figure 3: Proposed assessment form (Google Form)

# 5    Conclusions

In this paper we have presented some insights about how the comparison criteria could benefit of some rigorous studies. Correlation analyses and integration improvements can be explored through advanced statistical methods such as Multiple Component Analysis and Hierarchical Clustering. We believe such tools can help give insights about many questions relating to similarities and differences between modeling approaches, and/or comparison criteria characteristics complex relationships. Data collection can somehow be enhanced too with a little revamping of the questionnaire so it better feeds the statistical tables and minimizes loss of information.

We do not have ourselves the required resources to conduct those studies. Moreover those studies have to be discussed during the workshop. But we hope the proposed approaches will be useful for the workshop organizers. Our purpose is to sketch out a general guideline of what can be done and how, with resource efficiency and overall pragmatism in mind.