
Abstract

Context. Prioritization is essential part of requirements engineering, software release planning and many other software engineering disciplines. Cumulative Voting (CV) is known as relatively simple method for prioritizing requirements on a ratio scale. Historically, CV has been applied in decision making in government elections, corporate governance, and forestry. CV prioritization results are special type of data – compositional data.

Objectives. The purpose of this study is to aid decision making by collecting knowledge on the empirical use of CV and developing a method for detecting prioritization items with equal priority.

Methods. We present a systematic literature review of CV and CV result analysis methods. The review is based on search in electronic databases and snowball sampling of the primary studies. Relevant studies are selected based on titles, abstracts, and full text inspection. Additionally, we propose Equality of Cumulative Votes (ECV) – a CV result analysis method that identifies prioritization items with equal priority.

Results. CV has been used in not only in requirements prioritization and release planning but also in software process improvement, change impact analysis, model driven software development, etc. The review has resulted in a collection of state of the practice studies and CV result analysis methods. ECV has been applied to 27 prioritization cases from 14 studies and has identified nine groups of equal items in three studies.

Conclusions. We believe that collected studies and CV result analysis methods can help the adoption of CV prioritization method. The evaluation of ECV indicates that it is able to detect prioritization items with equal priority.

Keywords:

Cumulative voting, Hundred-dollar test, \$100 test, requirements prioritization, Systematic review

1. Introduction

Software products are becoming larger and more complex. Each product is usually affected by a large number of factors such as product functional requirements, quality attributes, or software process improvement issues. Since time, funding, and resources are limited, it is seldom possible or efficient to fully address all the factors. Therefore, the level of attention to a particular factor must be decided according to its importance (i.e. business value), cost, risk, volatility, dependencies between the factors and other criteria. These type of decisions are made by product stakeholders: users, clients, managers, sponsors, developers, and other persons associated with the product. In order to make decisions regarding a large number of factors it is highly advisable to prioritize the factors in a systematic way [1].

One of the prioritization methods used in software engineering is Cumulative Voting (CV) [2]. The main advantage of CV is that it is relatively simple and fast, yet produces priorities in ratio scale [1, 3]. This allows us not only to determine what prioritization items are more important but also how much more important they are. (Ratio scale prioritization is particularly important in software release planning and cost-value analysis [4, 5].)

Prioritization is usually performed by multiple stakeholders where individual priorities are combined into a single priority list. Each stakeholder's preferences may have different weight in the final priority. Such prioritization provides more information than just the priorities of factors. It may be useful to analyze the results of the prioritization to assess disagreement between stakeholders, measure stakeholder satisfaction with the results or find distinct groups of stakeholders.

The purpose of this study is to help industry practitioners and academia researchers in adopting, using and developing CV, while the importance of prioritization in software engineering and the prospectiveness of CV constitutes a need to do further research in this area.

This study presents a systematic literature review of the empirical use of CV and CV result analysis methods. A new method for CV result analysis, called Equality of Cumulative Votes (ECV), is proposed. The method identifies prioritization items with *equal* priority. ECV is evaluated using a considerable amount of data, which was obtained from the primary studies identified by the systematic review (through the kindness of the authors of said studies).

The remainder of this paper is structured as follows. The background is presented in Section 2. Section 3 describes related studies. In Section 4 research questions and methods are presented. The design of the systematic

40 review is presented in Section 5 and ECV is presented in Section 6. Section 7
41 presents the results of the study and Section 8 is a discussion section.

42 **2. Background**

43 This section presents definitions and places this study in a context. In the
44 coming sections we will cover: a description of software requirements priori-
45 tization methods; examples of CV result analysis methods; and a description
46 of compositional data analysis and CV.

47 *2.1. Prioritization Methods*

48 Some of the most popular prioritization methods are the analytical hi-
49 erarchy process (AHP), cumulative voting (CV), ranking, numerical assign-
50 ment, top-ten, the planning game, minimal spanning tree, bubble sort and
51 binary search tree [1, 6]. Ranking and numerical assignment methods per-
52 form prioritization on an ordinal scale. AHP and CV are, on the one hand,
53 considered to be harder to use and also more time consuming compared to
54 other methods but, on the other hand, produce priorities in ratio scale.

55 Prioritization can be used not just to decide which factors to address, but
56 also to determine the order in which they need to be handled. In market-
57 driven software development a small part of a very large number of require-
58 ments need to be selected and divided into several releases to maximize return
59 on investment. While in bespoke requirements, focusing on early delivery of
60 value can help reduce the risk of project cancellation.

61 Ratio scale priorities have several advantages over ordinal scale priori-
62 ties. Ratio scale shows not just the order of items but also relative distance
63 between them. This enables the priority of a group of items to be calculated
64 by summing up the priorities of individual items [4]. It is possible to say
65 that one item or set of items has higher priority than another set of items.
66 Supposing stakeholders have to choose between several low priority items
67 and one item with higher priority; with ordinal scale, the item with high-
68 est priority will always be selected first. However, if priorities are given on
69 a ratio scale, it is possible that lower priority items will be selected if their
70 cumulative priority is higher. Knowing the relative importance of sets of pri-
71 oritization items helps in software release planning. Ratio scale allows the
72 combining of multiple priority factors by calculating ratios between them.
73 One example of this is the cost-value ratio that shows which requirements
74 give more value for less money [5].

75 2.2. Prioritization Result Analysis

76 Different studies use and analyze CV in different ways. Disagreement
77 between stakeholders happens when two or more stakeholders have assigned
78 a different priority to one prioritization item. If the level of disagreement is
79 high it may indicate potential conflicts between stakeholders. Such conflicts
80 may be of technical character, as well as social or cultural.

81 The satisfaction a stakeholder has with the final prioritization results is
82 determined by the difference between the results and the individual priorities
83 of the stakeholder. A smaller level of difference leads to higher satisfaction.
84 In the end, stakeholder satisfaction is important because it is necessary to
85 achieve stakeholder commitment.

86 In some cases a part of stakeholders may form a group of some kind and,
87 therefore, prioritize requirements similarly. It may be useful to detect whether
88 a group of stakeholders has different preferences than all other stakeholders.
89 As an example, in [7] domain experts, technical experts, managers, project
90 managers, testers, and developers use CV to prioritize software process im-
91 provement issues and the CV results are analyzed using disagreement charts
92 and satisfaction charts. Finally, principal component analysis (PCA) is used
93 to identify distinct groups of stakeholders.

94 The same items can be prioritized by the same stakeholders multiple
95 times from different perspectives. In this case it is useful to determine corre-
96 lation between the priorities in different perspectives to assess the differences
97 between the perspectives. As an example, in [8] CV is used by developers,
98 testers, and managers to prioritize quality attributes. The same quality at-
99 tributes are prioritized from two perspectives: the perceived situation today
100 and the perceived ideal situation. Correlation between the two perspectives
101 is evaluated using the Spearman rank correlation matrix. This allows an
102 analysis of how well the company balances the priorities of software quality
103 attributes.

104 In [9] change impact issues are prioritized by developers, testers, man-
105 agers, and system architects. The prioritization is done with respect to three
106 perspectives: strategic, tactical, and operative. In order to determine corre-
107 lation between the perspectives, CV results are analyzed using the Kruskal-
108 Wallis test. In [10] the results of [9] are further analyzed using PCA, bi-plot,
109 and ternary plot. In this case, PCA is used to find correlated issues, bi-
110 plot shows variance, correlation, difference between the priorities of issues,
111 and the viewpoints of stakeholders, while ternary plots are used to show the
112 relative number of issues that received high, medium, and low priority.

113 As can be seen above, from the examples given, prioritization has been
114 performed with various stakeholders, using different perspectives and, in the

115 end, also analyzed using various techniques. We will next describe in more
116 detail one of the more common methods to manage prioritization issues —
117 cumulative voting — which has been used in software engineering for some
118 time, but has its roots in corporate governance and biology.

119 2.3. Cumulative Voting

120 CV is a prioritization method for prioritizing a list of items [2]. CV has
121 many synonyms in literature: hundred dollar method, hundred dollar test,
122 hundred point method, 100\$ dollar method, 100\$ dollar test, 100\$ point
123 method. Before being applied in software engineering CV was used for polit-
124 ical elections [11] and corporate governance [12]. CV has also been applied
125 in e.g. decision making in forestry [13], voting in social networks [14] and in
126 computer algorithms for consensus clustering [15] (as a method for combining
127 the results of different clustering algorithms).

128 In CV a stakeholder is given 100 points, imaginary dollars or units of
129 percentages that can be spent on the prioritization items. In the simplest
130 case, the stakeholder can spend any amount of points on any number of items
131 as long as the total amount adds up to 100. The more points assigned to an
132 item, the higher the priority of the item (and implicitly, the lower priority
133 to the other items). The stakeholder may spend all the points on just one
134 item or distribute them among all or some of the items. Once again, this is
135 the simplest case; other variants exist, which we will see next.

136 Often prioritization is done by more than one stakeholder. The final
137 priority of an item can be calculated by adding up the points each stakeholder
138 has spent on it. Sometimes the vote of some stakeholders may be more
139 important than the votes of others. For example, a manager may be more
140 influential and shareholders may have different amount of shares. In such
141 a case the priorities of each stakeholder may be multiplied by an individual
142 coefficient or a different amount of points for prioritization.

143 Worth mentioning in this context is that it is advisable to randomize the
144 order of items in a prioritization list. This is necessary in order to minimize
145 the effect of order on the prioritization results, which has shown to have an
146 effect [16].

147 2.3.1. Benefits and Drawbacks of Cumulative Voting

148 Compared to analytical hierarchy process (AHP), CV is faster and easier
149 to learn and use [1, 3]. AHP benefits from consistency check, but CV does
150 not require this because all prioritization items are evaluated simultaneously
151 [3].

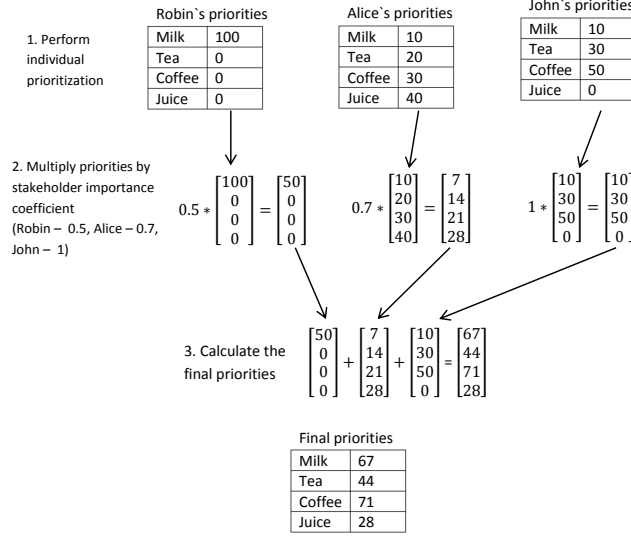


Figure 1: Example of CV with several stakeholders.

152 There are, however, a few problems with CV. First of all, it cannot be
 153 repeated for the same stakeholders and prioritization items due to stake-
 154 holder bias [2] (c.f. Section 2.3.4). Secondly, CV becomes more difficult if
 155 the number of prioritization items increases [17].

156 2.3.2. Example of Cumulative Voting with Several Stakeholders

157 Let us give an example of CV with several stakeholders. Suppose Robin,
 158 Alice, and John are three friends who want to buy some beverages in a store.
 159 They have different preferences but do not want to buy too many drinks.
 160 Therefore, they decide to use CV to decide what to buy. Each of the friends
 161 distributes 100 points between four items: milk, tea, coffee, and juice (Step
 162 1 in Figure 1). Each of them will spend a different amount of money on
 163 the purchase, hence, their priorities are multiplied by different coefficients
 164 (Step 2 and the stakeholder importance coefficient in Figure 1). The final
 165 beverage priorities are calculated by summing up the weighted priorities of
 166 stakeholders (Step 3 in Figure 1).

167 2.3.3. Stakeholder Bias

168 Prioritization using CV may be biased if a stakeholder knows the pref-
 169 erences of other stakeholders. She may manipulate the results by spending

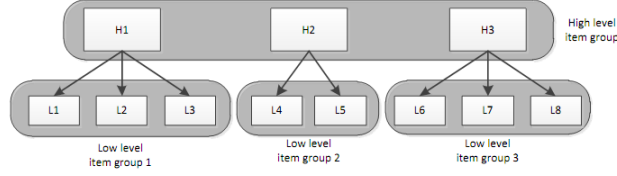


Figure 2: Example of prioritization item hierarchy.

more points on items that are important to her but not the other stakeholders. On the one hand, stakeholder bias makes it unreasonable to repeat CV with the same prioritization items and stakeholders. On the other hand, this property of CV may be useful in giving more power to important minority stakeholders, such as security experts or software testers. Suppose the same software requirements are prioritized for a second time using CV. A developer might know that all vital functionality is selected by other stakeholders, but his toy feature is left out. In effect, the developer could spend all his points on this feature to put it in the next release.

Stakeholder bias may be mitigated by setting a maximum priority that can be assigned to an item. This way each stakeholder is forced to distribute the money between several prioritization items [4].

Another bias is that people in general tend to assign round priority values. This is likely caused by lack of objective judgement criteria. Either way it seems to be a problem not acknowledged by many since all prioritization is largely based on expert opinion.

2.3.4. Scalability of Cumulative Voting, Hierarchical Cumulative Voting

The standard CV approach has a low scalability. If the number of prioritization items is high, stakeholders may lose sight of the bigger picture and assign priorities to a limited number of items. One, unsophisticated, solution to the problem is to provide more points for prioritization (1,000 or 10,000 instead of 100); however, one could take another approach.

When the number of prioritization items is high they can usually be grouped hierarchically by forming a tree structure (Figure 2) and, thus, parent-child dependencies will exist between many items.

In [4] the authors propose a method for prioritizing hierarchically structured items called Hierarchical Cumulative Voting (HCV). It may be seen as combination of the hierarchical part of the Analytical Hierarchy Process (AHP) [1, 18] and the CV prioritization method. Since items are prioritized in smaller sets, stakeholders do not lose sight of the bigger picture during

200 prioritization, and the prioritization of a large number of requirements is
201 considered easier.

202 2.3.5. *Compensation Factors*

203 HCV deals with the problem of prioritization scalability but it comes at
204 a cost. Low level item groups may consist of different numbers of items, but
205 the number of points spent on each group is the same, i.e. in a small-sized
206 group, the same amount of points is distributed among fewer items. Hence,
207 items in smaller groups are statistically more likely to have a higher priority,
208 on average, compared to items in larger groups. To balance this difference
209 each low level prioritization item can be multiplied by a compensation factor
210 [4].

211 As an example, suppose an item (A) in a group of 10 items is assigned
212 60 points. Hence, A will receive 600 compensated points. In this case it is
213 impossible for any item in a group smaller than 6 items to compete with A .
214 Even if item (B) in a group of 5 is assigned the maximum number of points
215 (100), the maximum compensated priority value B can receive is 500.

216 In [17] the authors suggest that compensated prioritization is more fa-
217 vorable compared to uncompensated. But neither compensated nor uncom-
218 pensated prioritization is perfect and, as a general rule, it is better to keep
219 the size of prioritization item groups similar.

220 2.3.6. *HCV Execution*

221 According to [4], HCV is conducted with the following steps (Steps 4–5
222 are optional):

- 223 1. Construct hierarchy. Prioritization items need to be divided into one
224 high and several low level item groups. Each low level item group is
225 child to exactly one high level item. And each high level item has
226 one low level item group. One low level item may belong to several
227 item groups. Even if part of the items are not logically connected they
228 can be grouped separately and assigned a fake parent item, e.g. ‘misc.
229 items’. HCV does not, as far as we know, provide any directions on
230 creating a requirements hierarchy.
- 231 2. Each high and low level item group is prioritized separately using CV.
232 The stakeholder may prioritize all item groups at once or one by one.
233 But it should be possible to prioritize groups in any order and repeat-
234 edly, because the stakeholder might learn more about the items while
235 performing the prioritization.
236 In particular the stakeholder is likely to learn more about a high level
237 item when prioritizing its low level item group [19]. Some stakeholders

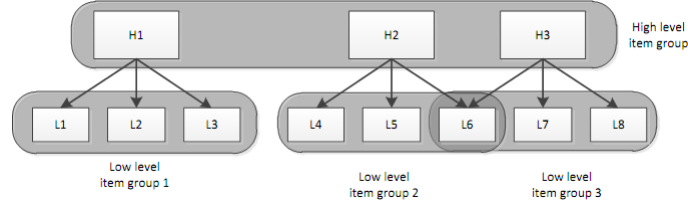


Figure 3: Overlapping prioritization item hierarchy example.

may prioritize only part of the groups and each group may be prioritized by different stakeholders.

3. The priority of each low level item is normalized by dividing it with the sum of all low level priorities of each item in all groups.
4. The final priority of each low level item is calculated by multiplying it with the priority of its parent high level item.
5. Then apply the compensation factor to all low level requirements as described in Section 2.3.5.
6. Finally, when multiple stakeholders have performed the prioritization, priorities of low level items are combined as in standard CV.

It is possible that one low level item is child of more than one high level requirement and, thus, belongs to two or more low level requirement groups (see Figure 3). Such requirements participate in the standard HCV prioritization process and are prioritized two or more times with each group they belong to. At the end of the prioritization they receive several priority values. These values must be summed together to form the final priority of the item. (This is done because the item adds value to both parts of hierarchy.)

2.3.7. Example of Hierarchical Cumulative Voting

In this section we will give a short example of HCV. Suppose six requirements for a mobile phone operating system need to be prioritized: ‘reminder alarm’, ‘specify repeated event’, ‘hide contact’, ‘add picture to phonebook’, ‘search contact’, ‘make video call’. Three high level requirements can be identified: ‘Calendar’, ‘Phonebook’, ‘Call’. The low level requirements are then grouped as sub-requirements of high level requirements as shown in Figure 4. The ‘Search contact’ requirement is a sub-requirement and has two parent requirements: ‘Phonebook’ and ‘Call’. The computation of the final priorities of requirements is shown in Table 1.

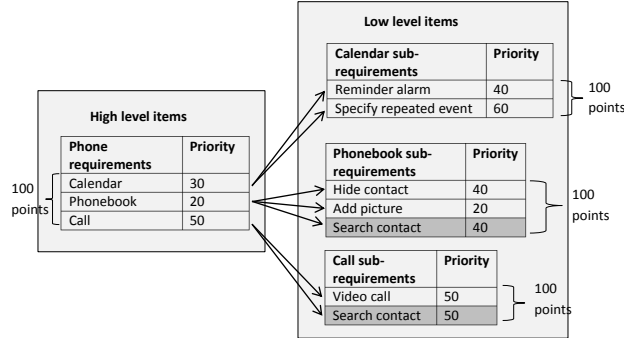


Figure 4: Example of hierarchical cumulative voting, requirement hierarchy

Table 1: Example of hierarchical cumulative voting.

Phone requirements	Compensation factor	Sub-requirements	Priority calculation	Final priority
Calendar	2	Reminder alarm	$40*30*2$	2400
Calendar	2	Specify repeated event	$60*30*2$	3600
Phonebook	3	Hide contact	$40*20*3$	1600
Phonebook	3	Add picture	$20*20*3$	800
Phonebook & Call	3 & 2	Search contact	$40*20*3 + 50*50*2$	7400
Call	2	Video call	$50 * 50 * 2$	2500

After requirements are grouped, and a hierarchy is defined, each group of requirements are then prioritized using CV. The final priority of a low level requirement is computed by multiplying the priority of the requirement with the priority of its parent high level requirement and the compensation factor. The compensation factor in this particular case is the number of elements in a group, two for the ‘calendar’ and ‘call’ sub-requirements and three for the ‘phonebook’ sub-requirement.

2.4. Compositional Data Analysis

CV results can be seen as a special type of data, i.e. compositional data. Compositional data does not contain absolute values. It shows only the relative weight of a component in a whole. In [10] the authors propose the use of compositional data analysis for the statistical analysis of CV.

A compositional data item is a vector (x) of positive components with a constant sum k :

$$x = (X_1; X_2; \dots; X_n) \text{ where } x_i \geq 0 \text{ and } \sum_{j=1}^n x_j = k \quad (1)$$

280 The property of the sum of the items being restricted is called the con-
 281 stant sum constraint. In CV, priorities assigned by a stakeholder to the
 282 items of a prioritization set is a compositional data vector with a constant
 283 sum of 100. The value of k (i.e. 100 in this case) is arbitrary and does not
 284 affect the analysis of the data because the information is contained in the
 285 ratios between the components of the vector. The vector can sum up to any
 286 number but still hold the same data, i.e. vectors (1, 2, 7) and (10, 20, 70)
 287 are in this case considered equivalent.

288 The priority of an item is relative to the priority of the other items in
 289 the set. Hence, the priority of an individual item is meaningless without
 290 context, i.e. the complete set of items. The same item may receive different
 291 priority when put in two different prioritization sets. If the item is put in a
 292 set of items with high priority it will receive a lower relative priority. This
 293 also holds true the other way around i.e. if the item is put in a set with low
 294 priority items its priority will be higher.

295 Compositional data analysis has, however, serious limitations. Ordinary
 296 unconstrained variables are free to take any positive or negative values,
 297 whereas, compositional data values can only be positive and have a con-
 298 strained maximum value. Moreover, components of compositional data vec-
 299 tors are not independent from each other. The fact that an item is assigned
 300 70 priority points means that the next item can take only values between 0
 301 and 30. Hence, there is a negative correlation between the items.

302 Standard parametric statistical tests require that data vectors have mul-
 303 tivariate normal distribution. Vector $X = (X_1, X_2, \dots, X_n)$ is considered to
 304 have multivariate normal distribution if any linear combination of its parts
 305 is normally distributed, and linear combination is defined by:

$$Y = a_1 X_1 + a_2 X_2 + \dots + a_n X_n \quad (2)$$

306 where Y is the product of lineal combination and a_i is any real number.
 307 Now, since the sum of priorities assigned in CV must add up to 100 (or any
 308 other constant number) at least one linear combination of X is not normally
 309 distributed because it must always add up to 100:

$$Y = 1 \cdot X_1 + 1 \cdot X_2 + \dots + 1 \cdot X_n = 100 \quad (3)$$

310 In our opinion, the above indicates, quite strongly, that CV results do
 311 not follow a multivariate normal distribution and, hence, it follows that they
 312 should be analyzed using non-parametric statistical tests [20].

313 2.4.1. Problem of Zeroes

314 Compositional data analysis requires that ratios between any components
 315 in a vector can be computed. But computing a ratio with a zero value is,
 316 in this case, meaningless. This is a problem since CV allows stakeholders to
 317 assign zero priorities to some prioritization items (we would even strongly
 318 argue that this is very common).

319 In compositional data there are two types of zeroes: essential and rounded.
 320 Essential zeroes mean that a data component is not present. Rounded zeroes
 321 mean that the component is present but its value is very low. We, as others
 322 have before us, conjecture that zeroes in CV results are rounded because the
 323 priority of an item is a completely abstract notion and the instrument for
 324 measuring priority is human judgement [10].

325 Before compositional data analysis can be applied to CV results, we must
 326 first remove zeroes in the data. One approach can be to forbid stakeholders to
 327 assign zero priorities. This approach is used in e.g. [7]. But this can add some
 328 unnecessary complexity to the prioritization process and, explicitly, delimits
 329 an expert's freedom. In [10] the authors propose the use of a multiplicative
 330 replacement strategy (as defined in [21]) for CV result analysis.

This method replaces rounded zeroes with small values using the expres-
 sion

$$r_j = \begin{cases} \delta_j, & \text{if } x_j = 0, \\ (1 - \frac{\sum_{k|x_k=0} \delta_k}{c})x_j, & \text{if } x_j > 0, \end{cases} \quad (4)$$

331 where δ_j is the imputed value and c is the constant sum constraint.
 332 In order for the total sum of components to stay constant, the equation
 333 subtracts some value from the items with a priority higher than zero. More
 334 is subtracted from components with higher values than from components
 335 with lower values (and the value of the imputed δ_j is arbitrary).

336 2.4.2. Isometric log-ratio transformation

337 In order to apply standard statistical methods to compositional data it
 338 must be transformed to remove the inherent correlation of the values. Com-
 339 positional data analysis proposes special transformations that change the
 340 compositional data values to unconstrained real values. One such transfor-
 341 mation is isometric log-ratio (*ilr*) transformation (as proposed by [20, 22]):

$$\begin{aligned}
z &= (z_1, \dots, z_{D-1}), \\
z_i &= \sqrt{\frac{i}{i+1}} \log \frac{\sqrt[i]{\prod_{j=1}^i x_j}}{x_{i+1}} \text{ for } i = 1, \dots, D-1
\end{aligned} \tag{5}$$

where x is the vector that is being transformed and z is the vector that is created. It should be noted that z is shorter than x by one element.

After compositional data vectors are transformed using zero replacement and *ilr*, any standard statistical tests can be applied.

3. Related Work

A systematic review of requirements prioritization methods is presented in [23]. The study focuses on prioritization method comparison and selects eight relevant studies. Two of the studies use CV. These studies are also revealed by the systematic literature review conducted as part of this study. Khan [23] concludes that there is little research on requirements prioritization and studies usually deal with a small number of requirements.

The systematic literature review presented in this paper does not reveal any CV result analysis methods that allows to identify prioritization items with equal priority. Thus, this problem is not addressed in any way.

4. Methodology

This section covers the research questions of this study and the methods used to answer them.

4.1. Selection of Research Methods

The main purpose of this study is to collect knowledge on the use of CV in order to help software engineers and researchers in adopting it.

One way of collecting this knowledge is to conduct an empirical study. A survey in a large number of software companies can be used to quantify the level of adoption of CV in industry (similarly to the study by [24]), while a case study can be used to receive qualitative feedback on the use of CV [25].

Knowledge on the empirical use of CV can also be obtained from existing studies. This may be done by means of a systematic literature review. Several studies have used CV in industry as well as in academic settings. Nevertheless, there are no studies that provide an overview of the current state of the practice in this field (as reported by research studies). Therefore,

371 before continuing with the refinement of CV and conducting new empirical
372 studies (i.e. case study or experiment), a systematic literature review would
373 be required.

374 This paper proposes a new method for CV result analysis, called Equality
375 of Cumulative Votes (ECV). (ECV groups prioritization items into groups
376 of items with similar priority.) As will be presented later, the systematic
377 review did not reveal any methods that solve this problem; however, ECV
378 needs to be evaluated and, hence, applied to CV results.

379 There are two options to obtain CV results in order to test ECV. One
380 is to conduct a new empirical study. The second option is to collect CV
381 results from existing studies. The latter approach also has the added ben-
382 efit of trying to replicate the results from previous studies and, if the CV
383 results from other studies are used, a larger amount of data can be obtained.
384 Moreover, the generalizability of the evaluation increases when prioritization
385 results from different sources and domains are used. On the other hand, the
386 main benefit of conducting a separate empirical study is the possibility to
387 control the conditions of CV.

388 In our study we evaluated ECV by obtaining data from previously con-
389 ducted studies as found by the systematic literature review. In order to
390 obtain the data, authors of relevant primary studies were contacted.

391 In short, this study consists of two parts: a systematic literature review
392 (SLR) of CV and an evaluation of ECV based on the data from the primary
393 studies found in the SLR.

394 *4.2. Research Questions*

395 The systematic review should focus on catching studies that empirically
396 use CV. Information about place, time, scale, and domain of the studies
397 should be collected and the results of the review will hopefully aid academic
398 researchers by identifying paths for further investigation of CV. Hence, the
399 first research question is:

400 **RQ 1.** What is the state of practice in empirical studies that use CV?

401 The level of trust in research results considering CV is determined by the
402 quality of the studies that use CV, hence this study includes an evaluation
403 of the quality of primary studies identified by the systematic review.

404 Next, a valuable aspect of decision making is the analysis of prioritization
405 results. Thus, the second research question is:

406 **RQ 2.** What CV result analysis methods have been presented in papers as
407 identified by RQ 1?

408 Finally, the evaluation of ECV answers the third research question:

409 **RQ 3.** Is ECV capable of identifying prioritization items with equal prior-
410 ity?

411 5. Systematic Literature Review

412 This section presents the design of the systematic literature review. For
413 the results of the execution please see Section 7.1 and 7.2.

414 Table 2 presents an overview of activities performed during the system-
415 atic literature review. The review protocol was developed by one researcher
416 and evaluated by another researcher. Studies were searched for in two itera-
417 tions. The first search was performed by using databases. The second search
418 was performed using snowball sampling [26] (snowball sampling examines the
419 references of primary studies revealed by the first search). References that
420 are relevant to the review, i.e. they pass the selection criteria, are then added
421 to the set of primary studies.

422 The search for papers was performed by a single researcher. Study se-
423 lection, on the other hand, was performed by two researchers. First, one
424 researcher examined all found studies. Next, another researcher re-examined
425 all studies classified as primary studies in addition to 20 randomly selected
426 excluded studies to ensure the quality of the selection.

427 To ensure the quality of the review, the quality evaluation and data ex-
428 traction was performed independently by two researchers. Inter-rater anal-
429 ysis was performed using Krippendorff’s Alpha statistics [27, 28].

430 5.1. Data Sources and Search Strategy

431 This SLR was designed based on the guidelines by Kitchenham [29]. First
432 a trial search in electronic databases was conducted. In order to scale the
433 review to a manageable, yet sufficient size, databases were searched with
434 different search strings. Relevant papers that were found during the trial
435 search were used to extract additional search strings. The trial search re-
436 vealed that the number of studies that use CV is not very large. Therefore,
437 we decided to include not only software engineering studies but also studies
438 in other research areas, such as forestry or corporate governance, since one
439 key aspect we intended to investigate was analysis methods for CV.

440 Since CV is frequently used in studies without mentioning this in the
441 abstract, full text search in databases is preferable. Unfortunately not all
442 databases support full text search. Full text search was performed in the

Table 2: Review activities.

Review phase		Researchers involved
Trial search in databases		A
Develop review protocol		A
Evaluate review protocol		B
Paper search and selection from databases	Search in databases	A
	Search string validation	A
	Selection based on metadata	A and B
	Selection based on full text	A and B
Pilot data extraction (3 papers)		A
Paper selection from the reference lists	Selection based on metadata	A and B
	Selection based on full text	A and B
Data extraction		A and B
Data synthesis		A

A – Cumulative voting	E – hundred dollar method
B – 100 dollar method	F – hundred dollar test
C – 100 dollar test	G – hundred point method
D – 100 point method	

443 IEEE Xplore and Springer Link databases. In ACM Digital Library, In-
444 spec/Compendex, ISI Web of Knowledge, and SCOPUS only metadata was
445 searched. Search strings consisting of a Boolean expression (A or B or C or
446 D or E or F or G), where:

447 Search strings contained only synonyms of CV and they did not limit the
448 research area to software engineering. The search was performed indepen-
449 dently using each of the search strings in each database. All search results
450 were combined and documented using reference management software. The
451 quality of the search strings and the selection of electronic databases were
452 validated against a previously known core set of papers—[3, 30, 10, 31]—
453 checking that all papers from the core set were found by the search.

454 5.2. Study Selection

455 To select relevant papers a set of criteria were designed. The criteria for
456 paper selection are presented in Tables 3 and 4.

457 Papers were selected in two phases: based on metadata and based on full
458 text.

459 Obviously, the main criterion for inclusion of a paper is that it must
460 present empirical use of CV or present an analysis of the results of using
461 CV. However, there are papers that pass this criterion but are not relevant
462 for this review. CV is frequently used in computer algorithms. There is
463 a significant difference between the way that humans and computers make
464 decisions. Since this review is concerned with human decisions we excluded
465 papers that present CV that is not performed by humans. In addition, only
466 papers that were written in English were selected and duplicate studies were
467 automatically excluded by the citation management software used in this
468 review.

469 5.3. Quality Evaluation

470 The goal of quality evaluation is to determine the best primary studies
471 according to some measure of quality. Since the number of studies that use
472 CV is not large, quality evaluation was not used as an exclusion criterion.

473 Study quality obviously depends on the correctness of the study process
474 including planning, operation, analysis and interpretation of the results (is
475 the study right?) The correctness of the process can be measured by eval-
476 uating the description of the study or replicating the study. Thus, to gain

Table 3: Paper search and selection in the databases.

Selection phase	Inclusion criteria	Number of papers selected
Search in databases	published from 2001 until 2011 (databases last accessed Feb. 20, 2011)	256
	contains search strings	
Selection based on metadata	exclude duplicates and tables of contents	177
	written in English	
Selection based on full text	full text is available	127
	study involves empirical use of CV or presents analysis of empirical use of CV	58
	CV is done by humans and not software	25

Table 4: Paper selection from the reference lists of the selected papers.

Selection phase	Inclusion criteria	Number of papers selected
Selection from references	papers included in the reference lists of relevant papers found in databases	467
Selection based on metadata	written in English	462
	reference is already revealed by search in databases	450
Selection based on full text	full text is available	329
	study involves empirical use of CV or presents analysis of empirical use of CV	15
	CV is done by humans and not software	

477 the trust of industry practitioners and other researchers, the process of the
478 study must be rigorously described. In short, the description must facilitate
479 replication of the study as well as the presentation of limitations and validity
480 threats.

481 Even the most correct and rigorously described study is useless if it does
482 not contribute to the industry or research community (is it the right study?)
483 The topic of the research ought to address important goals and issues. The
484 findings of the study should also be significant, i.e. there must be a high
485 probability of the results of the study being true. The significance of the
486 findings depends on how realistic the study is, the correctness of the process
487 and the results of the study, as well as the statistical significance of the
488 findings.

489 **Realism** of a study depends on the context, scale, and subjects of the
490 study. The study should be conducted in a **setting** that is similar or equal
491 to the setting in which the findings of the study are intended to be used.
492 Hence, studies that are conducted in an industrial setting are in many cases
493 valuable. The **subjects** of a study should be similar to the people who are
494 supposed to use the findings of the study. The subjects ought to have appro-
495 priate work experience, role in the organization, skills, cultural background,
496 motivation, and so forth. The **scale** of a study refers to the size of the study
497 objects. In the case of this systematic review the scale of a study is mea-
498 sured as the number of prioritization items. Study in academia may have a
499 large number of prioritization items. At the same time, an industrial study,
500 with professionals as subjects, may involve a smaller number of prioritization
501 items.

502 Each study may have a different level of realism. Some studies involve
503 industry practitioners in an academic setting to simulate real word practice in
504 a laboratory environment. Other studies may involve academic researchers
505 that execute a project. For example, researchers may be developing open
506 source software. On the reality scale these studies are somewhere in between
507 the purely academic and industrial studies.

508 The **type** of the research study can be considered as a criterion for the
509 evaluation of study realism. [32] suggest that study designs that are more
510 rigorous (e.g. experiments) are more realistic than observational studies (e.g.
511 case study) due to a higher level of control. On the other hand [33] rate study
512 designs based on other criteria, i.e. how frequently each type of study de-
513 sign is used in an industrial or academic setting. If a study design is used
514 more in an industrial setting, then it is considered more realistic. For in-
515 stance, in software engineering, case studies are frequently used in industrial
516 settings, whereas, experiments are usually performed in academia using stu-

dents as subjects. Therefore, [33] argue that case studies are more realistic than formal experiments. Obviously the effect of study design on the study realism may be interpreted in different ways. Therefore, we will not use this parameter in our quality evaluation.

The statistical significance of the results of a study can be used to evaluate the significance of the study findings. This measure will not be used, because the studies that are evaluated belong to very different research areas, i.e. the significance levels of the findings of the studies are not directly comparable for meta-analysis. Additionally, sometimes no result is more interesting than a significant result. If study results do not conform to the expectations of researchers, this may reveal important gaps in existing knowledge.

The ultimate goal of research, at least in software engineering, is in many cases industry impact. However, most of the time ideas need to be developed and validated in academia before industry professionals will risk to adopt them. Therefore, academic impact is important as well. Academic impact is usually measured by the number of citations. Academic impact is also measured for particular researchers, using the number of papers she has published and the number of citations of her papers. This measure will not be used in our quality evaluation because it is somewhat biased. The number of citations is likely to be lower for newer papers and the number of papers that a researcher has published gives little information about the actual quality or impact of her research.

5.3.1. Rating of the Studies

The quality evaluation in our review is based on the evaluation of: (i) Study realism. (ii) Study scale. (iii) Availability of raw results of CV. (iv) Quality of the research methodology.

Realism of the studies is rated in three aspects: subjects, setting, and scale. The subjects and setting is rated according to Table 5. The total rating of study realism is determined by summing up the ratings of the two aspects. For instance, if a study is conducted with industry professionals as subjects in an academic context the study will receive rating 1.

In order to rate the scale of a study the number of prioritization items was counted. If a paper presents several prioritization cases only the prioritization with the largest number of the prioritization items is considered. If HCV is used all of the prioritization items on different levels are counted together. However, if an item is present in several groups in the hierarchy it is counted only once.

The availability of raw results of CV is rated separately because it is especially important for our purposes (and for most other researchers in

Table 5: Rating of study reality level

Aspect	Contribute to relevance (rating 1)	Do not contribute to relevance (rating 0)
Subjects	Industry professionals	Academia students or teachers, or other
Context	Industrial	Academia

Table 6: Research data availability rating

Rating	Study rating criteria
0	CV results was not provided in the paper and we was unable to obtain the results from the authors.
1	CV results are not provided in the paper but the data was obtained from the authors. Part of the data is lost or corrupted.
2	CV results are not provided in the paper but all the data was obtained from the authors.
3	All CV results are included in the paper or reference is given to online source where all the data can be accessed.

order to replicate a study). The data availability rating criteria is given in Table 6. If the data of a study are not available it is not possible to validate the results of the study and, hence, the credibility of the findings is lower. Ideally the data collected in the study should be presented directly in the paper. An alternative may be to make the data freely available online and reference the online source.

The quality of the research methodology of a paper is rated according to a checklist presented in Appendix C. The checklist is based on guidelines for presenting research studies as presented in [34, 35] and the guidelines for quality evaluation of research studies presented in [33, 29]. Evaluation is done with regard to the rigor of the description and correctness of the research process and reasoning. Checklist items represent issues that research studies should implement and present in research paper. The checklist also contains item descriptions or questions that are used to evaluate the quality. Each item in the checklist is rated according to criteria presented in Table 7. The final rating of correctness of the research process of a study is computed by summing up the ratings assigned to all items in the checklist.

Study rating criteria was validated during a trial data extraction. Two researchers each rated three randomly selected papers. Afterwards, differences in ratings were discussed and study rating criteria were updated to avoid differences in interpretation.

As a result of the rating each study was assigned four rating values on an ordinal scale. In order for us to perform a more advanced analysis of the quality evaluation results these ratings were then converted into ratio

Table 7: Rating of correctness of research process

Rating	Study rating criteria
0	No description provided.
1	Only basic information is provided about the checklist item. Or significant validity threats exist with regard to this item.
2	Description is sufficient. Some minor questions are left unanswered. Validity threats may exist but they are not likely to affect the results of the study.
3	Description is rigorous and clear. Questions presented in quality evaluation checklist in Appendix C are answered. Decisions of the study are well justified, alternatives are discussed. No unhandled validity threats can be identified.

Table 8: Example of rating values

Study	Realism	Research data availability	Correctness of research process	Number of prioritization items
ST1	2	0	15	6
ST2	1	3	20	69
ST3	0	3	10	6

scale ranks. For each study, the number of studies that have received lower ratings is counted. The resulting number is the rank of the study; thereby, the quality of a study is expressed as four rank values.

An example of rating values is shown in Table 8. Table 9 shows ranking values computed for the studies in Table 8. We can observe that study realism level rating for ST3 is 0. There are no studies that have a lower study realism. Therefore, realism ranking for ST3 is 0. ST1 on the other hand has the highest realism rating. Since ST1 has higher reality level than both ST2 and ST3 it is assigned reality level rank 2.

5.4. Data Extraction

The goal of data extraction is to understand how and why CV is used and how CV results are analyzed in research studies. Ultimately, this will allow us to answer the first and second research questions in our study.

Data extraction was documented with the help of spreadsheet software. Extracted data items are available from [36].

Table 9: Example of ranking values

Study	Reality level	Research data availability	Correctness of research process	Number of prioritization items
ST1	2	0	1	0
ST2	1	1	2	2
ST3	0	1	0	0

595 6. Equality of Cumulative Votes

596 In the previous section we described the execution of the systematic lit-
597 erature review. In order to perform a more thorough analysis later we here
598 present the design of ECV before presenting the results of the systematic
599 literature review. For the results of the evaluation of ECV please see Sec-
600 tion 7.3 (ECV is implemented in the *R* programming language [37] and the
601 code can be found at [38].)

602 In CV stakeholders may assign similar or equal values to several prior-
603 itization items. As a result the difference between the items is small. The
604 variation in priorities is caused not only by the difference between priorit-
605 ization items but also by human error and lack of information for decision
606 making. For instance, people tend to simplify the task of prioritization by
607 assigning rounded values to items or giving equal values to several items [39].

608 During prioritization it may be beneficial to know which items are equal.
609 A common example is software release planning where requirements are dis-
610 tributed among several product releases. If two or more requirements are
611 considered equal they can be freely interchanged between the releases, and
612 other criteria, such as cost or effort, may be used as sole indicators for plan-
613 ning that particular release.

614 6.1. Testing Equality of Two Items

615 There are two ways to determine which prioritization items have similar
616 priority. One approach is to find items that are different and consider other
617 items as equal. Another approach is to find items that are equal.

618 The first approach uses statistical tests to evaluate differences between
619 e.g. two population means, in order to determine that two items are different.
620 Populations in this case consist of priorities assigned by all stakeholders to a
621 particular prioritization item. The number of stakeholders that perform the
622 prioritization is frequently small. Hence, the size of the sample is very often
623 too small for statistical tests to detect a significant difference and the tests,
624 thus, identify too many equal items to make any useful conclusions.

625 ECV, in contrast, uses the second approach. It finds items that are
626 similar and the rest of the items are considered different. This method tests
627 the probability of the difference between the means of two items being smaller
628 than the given value. In short, ECV tests the probability of the means of two
629 prioritization items differing by less than 25%. If the probability is higher
630 than 70% the items are considered equal.

631 The input to ECV is an $n \times p$ matrix A that contains the raw results of
632 the prioritization. The columns of the matrix represent prioritization items

while rows represent stakeholders. ECV performs the following operations for the priorities of each of the two prioritization items:

1. Replace zeroes in CV results.
2. Transform the data using *ilr* transformation.
3. Determine distribution function using kernel density estimation.
4. Use the distribution function to find the probability that the difference between two prioritization items is smaller than 25%.
5. Form groups of equal prioritization items.

Since CV results are compositional data, zeroes in A must be replaced with other values. This is done using the multiplicative replacement strategy which is described in Section 2.4.1. Next, two columns are extracted from matrix A to create the new matrix B :

$$B = [a_{*,k} a_{*,l}] \quad (6)$$

where a is an element of matrix A , and k and l are the columns that represent items that are tested for equality.

The *ilr* transformation is then applied to each row of the matrix B and the new vector C is obtained. The equation for calculating elements of C using *ilr* transformation is:

$$c_i = ilr(b_{i1}, b_{i2}) = \sqrt{0.5} \log(b_{i1}/b_{i2}) \quad (7)$$

where c_i is the i^{th} element of C and b_{i1} and b_{i2} are the first and second elements in the i^{th} row of B . Each value c_i represents a ratio between k and l . The mean of the values of C can be interpreted as an average ratio between the items that expresses the difference between the items.

After the data is transformed into log-ratios statistical test can be applied. The purpose of the test is to determine what the probability is of the relative difference between two prioritization items k and l being less than 25%. This means determining the probability of the ratio k/l between the items k and l as being in the range of $\frac{3}{4}$ to $\frac{4}{3}$. Or in terms of log-ratios it means determining the probability of $ilr(k, l)$ being between $ilr(3, 4)$ and $ilr(4, 3)$. Hence, the objective of the test is to determine the probability of the sample mean (i.e. mean value of C) laying between the two values.

The probability that the mean takes a particular value can be expressed in the form of a cumulative distribution function. The probability of the mean being between two values a and b (where a is smaller than b) can be

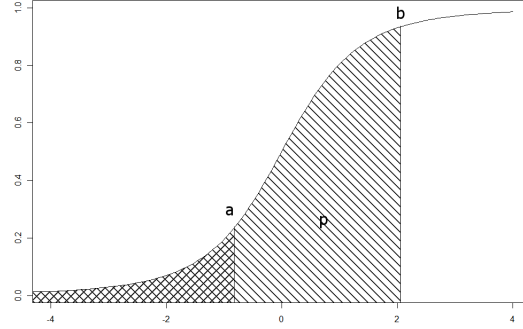


Figure 5: Cumulative distribution function of the ratio k/l between the items k and l (area p denotes probability that k/l is between $\frac{3}{4}$ and $\frac{4}{3}$.)

determined by subtracting the probability of the mean being smaller than a from probability of the mean being smaller than b .

However, CV result data may or may not be normally distributed. If the data is normally distributed a Student's t distribution function can be used.

Otherwise a non-parametric estimation of the distribution function is needed. In our case, the CV result data obtained from the primary studies identified by the systematic review, were tested for normality using the Anderson-Darling test. The tests we performed indicated, quite strongly, that in most of the prioritization cases the data is not normally distributed. Hence, our recommendation is that, in general, a non-parametric approach should be used to determine the probability density function, and one such, common, approach would be to use the kernel density estimation. (In our implementation of ECV in the R programming language, kernel density estimation is performed using the package *ks*.)

To determine the probability of \bar{x} being between a and b the following equation is used:

$$p = P(b) - P(a) \quad (8)$$

where P is the cumulative distribution function obtained by applying kernel density estimation on ilr -transformed priority values denoted by vector C . Variable a is equal to $ilr(3, 4)$ and b is equal to $ilr(4, 3)$. (A graphical interpretation of Equation (8) is presented in Figure 5.) The area that is denoted by letter p represents the probability computed by the equation.

After both prioritization items are tested for equality it may be convenient to display the equality of different items in the form of a table. Please see Table 10 for an example.

Table 10: Example of equality table

prioritization items	i1	i2	i3	i4
i1	equal	equal	-	equal
i2	equal	equal	-	-
i3	-	-	equal	-
i4	equal	-	-	equal

6.2. Grouping Prioritization Items

When equal items are determined they must be divided into groups of equal items. Division must be performed in such a way that each two items in a group are equal. The test for equality of the items described in Section 6.1 is not transitive. Hence, if prioritization item A is equal to B and B is equal to C then it does not automatically imply that A is equal to C . Therefore, there may be several ways to group the equal items. The two possible division criteria that we have considered in this study are:

1. Maximize the number of items that have a group.
2. Maximize the number of items in each group.

7. Results

This section presents the results of this study including the systematic literature review and the application of ECV on industry and academic data collected from the primary studies. Data extracted from primary studies and the results of the quality evaluation are available in [36].

7.1. State of Practice in Empirical Studies that use CV or Analyze the Results of CV (RQ 1)

The study search resulted in 634 unique studies. The search in databases revealed 180 papers, while an additional 454 papers were discovered using snowball sampling. The study selection resulted in 40 primary studies. Hence, 94% of the studies were excluded by the selection criteria. Snowball sampling revealed 15 or 36% out of all primary studies. The study selection criteria and the number of papers excluded by each criterion are shown in Tables 3 and 4. In total 163 of 634 studies were excluded because full text was not available.

All results of the study selection are available online and can be obtained by contacting the authors of this paper. For each study we specify keywords and databases that were used to find the study. If a study has been excluded, the exclusion criteria are provided.

718 The number of papers revealed by each search string and database is
719 presented in Table 11. It should be noted that several papers were found
720 by more than one search string or in more than one database. Table 11
721 shows that the search string ‘cumulative voting’ was the most frequently
722 used in research community to denote CV. Therefore, researchers should use
723 or reference this term when discussing CV.

724 To perform snowball sampling we examined the references of primary
725 studies that were found during the database search. References were used
726 to search for the papers in the Google and Google Scholar search engines.
727 Studies that were found in the search and passed the study selection criteria
728 were added to the set of primary studies.

729 After the primary studies were selected, data extraction and quality eval-
730 uation was performed by two researchers. One researcher examined all stud-
731 ies while the second researcher did quality evaluation and data extraction for
732 10% of the studies. The studies were randomly selected. Inter-rater agree-
733 ment were calculated by means of Krippendorff’s alpha coefficient. Agree-
734 ment for data extraction results was 0.86 and agreement for the quality
735 evaluation was 0.73. According to Krippendorff [28] it is common to re-
736 quire agreement above 0.8 and the lowest acceptable agreement is 0.667.
737 Therefore, we conclude that the agreement calculated for this study is suf-
738 ficient. Ratings of the study setting, correctness, research data availability,
739 and number of prioritization items are presented in Figure 6.

740 Table 12 shows the studies with the highest quality according to our cri-
741 teria. These studies show a high level of rigor in a realistic setting. Moreover,
742 authors of the studies manifest confidence by providing raw data for further
743 use and evaluation.

744 Figure 7 shows a bubble chart of the distribution of studies over research
745 areas and time. The figure shows that CV was first applied some time
746 ago in research of government elections. Nowadays, though, CV has been
747 adopted in a wide range of software engineering areas. Most frequently in
748 requirements engineering and software release planning. Eight studies use
749 CV as a research method while the remaining 32 studies report on using CV
750 in industry.

751

752 7.2. CV Result Analysis Methods Identified by RQ 1 (RQ 2)

753 The papers identified in the review use various CV result analysis meth-
754 ods. The main goals for CV result analysis are presented in Table 13 and

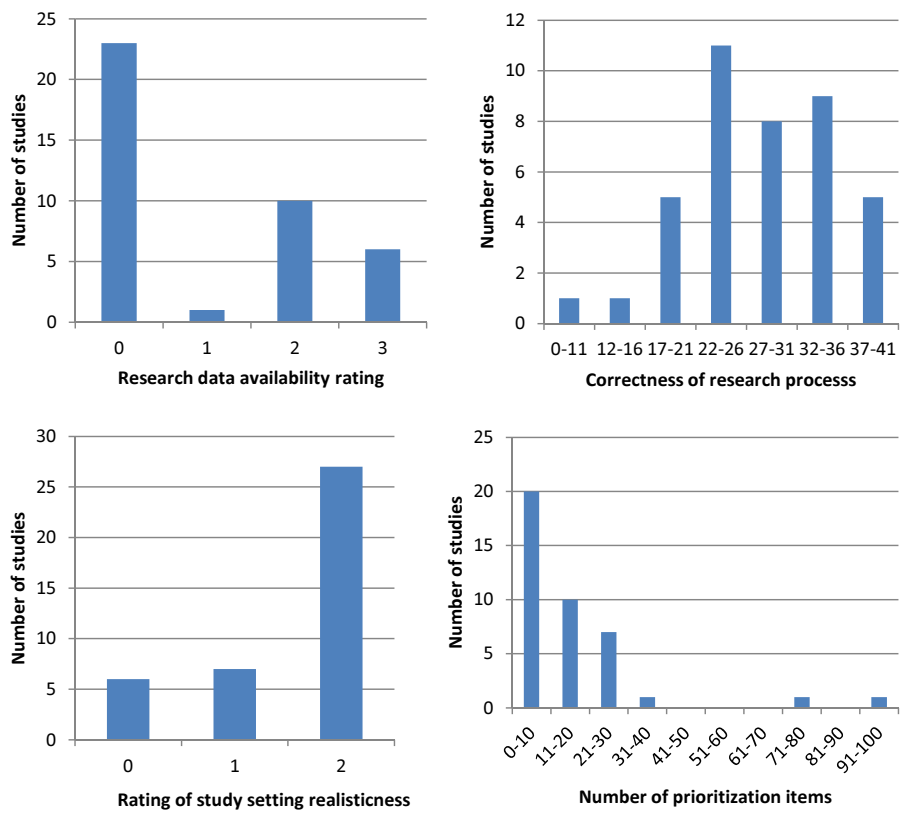
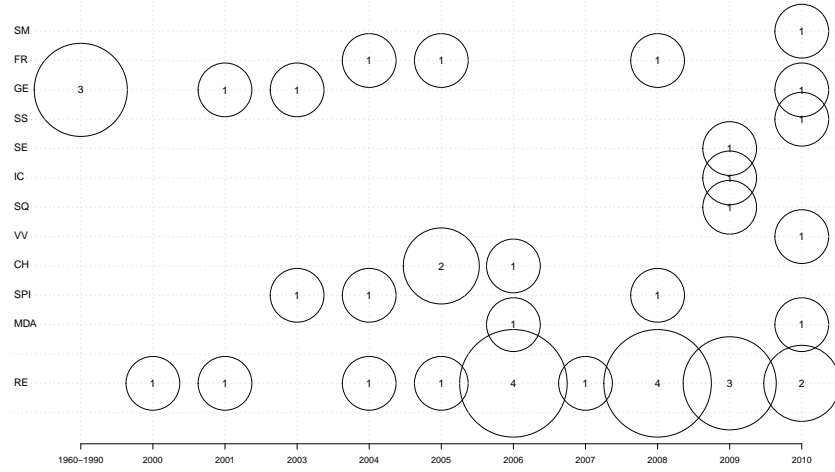


Figure 6: Study quality ratings



MDA - model driven software development
 CH - change impact analysis in software engineering
 RE - requirements engineering and software release planning
 IC - intellectual capital in software company
 SPI - software process improvement
 V&V - software verification and validation
 FR - forestry
 GE - government elections
 SS - software security
 SQ - software quality
 SM - software metrics
 SE - software engineering in general

Figure 7: Distribution of studies over time.

Table 11: Number of papers found in the databases.

database	search strings							unique papers found	primary studies selected
	"100 point method"	"100 dollar method"	"100 dollar test"	"hundred point method"	"hundred dollar method"	"hundred dollar test"	"cumulative voting"		
ACM	2	0	0	1	2	3	31	34	7
IEEE	3	2	0	1	2	6	38	46	11
Inspec/Compendex	1	0	0	1	1	1	22	14	7
ISI web of science	0	0	0	0	1	1	15	16	6
SCOPUS	2	0	0	0	1	2	24	25	9
Springer	2	0	2	0	2	2	89	95	6
unique papers found	6	2	2	1	4	11	165	180	
primary studies selected	1	2	1	1	2	4	18		25

Table 12: Top ranked studies.

	Correctness of research process	Research data availability	Study setting	Number of prioritization items
Barney et al. [40]	36	2	2	17
Berander and Svahnberg [17]	41	2	0	29
Barney et al. [41]	40	2	2	5
Barney and Wohlin [8]	31	2	2	27
Barney et al. [42]	34	2	2	14
Laukkanen et al. [43]	22	3	2	30
Hu [44]	34	2	1	14
Feldt et al. [45]	24	3	2	8
Regnell et al. [31]	21	3	2	91
Svahnberg et al. [46]	34	1	1	7

755 a summary of methods used in the primary studies can be found in Section
756 Appendix B.

757 In order to present prioritization results many studies use charts or tables.
758 These charts and tables show the average priority of each prioritization item
759 that is computed from priorities assigned by all stakeholders. In [47] a table
760 of five items with highest total priority is presented. [48] shows tables with
761 min , max , \hat{x} , \bar{x} and σ of priorities assigned by different stakeholders to a
762 particular prioritization item. Finally, in [49, 48] error bars are added to the
763 chart of final priorities (denoting σ of priorities).

764 In a few cases final priorities are presented in the form of ranks and
765 CV results are degraded from ratio to ordinal scale. This is done when the
766 interest lies only in the order of final priorities.

767 Several papers are interested in the difference between priorities from dif-
768 ferent prioritization perspectives (e.g. current and ideal situation) or stake-
769 holder groups (e.g. software developers and management). Pearson or Spear-
770 man correlation coefficients are commonly used to determine what the level of
771 similarity is between all priorities from two perspectives. Whereas, Wilcoxon,
772 Kruskal-Wallis, Nemenyi-Damico-Wolfe-Dunn tests and the χ^2 statistic are
773 used to detect if there is a significant difference in the value of one prioritiza-
774 tion item from two or more perspectives. In addition, PCA is used to detect
775 if there are distinct groups of stakeholders with common priorities [7, 10, 50].

776 In some cases, a stakeholder may assign equal priority to several prioritiza-
777 tion items or leave several items unrated, e.g. the stakeholder may not have
778 carefully considered all prioritization items. Hence, the difference between
779 the items may have been unnoticed.

780 In [4] the scalability of prioritization is measured using two charts. The
781 first chart shows the average percentages of items given a non-zero value.

Table 13: Goals for CV result analysis.

Purpose of the method	Name
Show the final priority of each prioritization item. Stakeholder priorities are combined into one value.	Chart or table of final priorities
Difference between priorities assigned by different perspectives (status quo, ideal situation) or different stakeholder groups (developers, management) [10]	Bi-plot
detect stakeholder groups with similar priorities [10]	Bi-plot
show the relative number of issues that have received high, medium, or low priority [10]	Ternary plot
detect stakeholder groups with common priorities [10]	PCA
how the final value of prioritization item is constructed from priorities assigned by different stakeholder. This chart shows how much each stakeholder has contributed to the final value of prioritization item [51]	Distribution chart
the level of agreement between different stakeholders on value of particular prioritization item [51]	Disagreement chart
satisfaction of a stakeholder with the prioritization results by the calculating correlation between the final priorities and priorities assigned by a stakeholder [51]	Satisfaction chart
percentage of the divergence of the priorities assigned by a stakeholder [4]	average percentage of divergence
average percentage of items given a non-zero value [4]	
detect equal prioritization items (presented in this paper)	ECV

782 The second chart shows average percentages of divergence of values. If a
783 stakeholder assigns equal priorities to many prioritization items the diver-
784 gence of values is low. Unfortunately it is unclear from [4] how the average
785 percentage of divergence is calculated.

786 In [51] distribution, disagreement, and satisfaction charts are presented.
787 The distribution chart shows how the final value of a prioritization item is
788 constructed from priorities assigned by different stakeholders. This chart
789 shows how much each stakeholder has contributed to the final value of a
790 prioritization item. The disagreement chart shows the level of agreement be-
791 tween different stakeholders on the value of a particular prioritization item.
792 The satisfaction chart shows stakeholder satisfaction with prioritization re-
793 sults by calculating the correlation between final priorities and priorities
794 assigned by a stakeholder.

795 The use of bi-plots and ternary plots are proposed in [10]. A bi-plot shows
796 final priorities and stakeholder viewpoints in a two dimensional plane while a
797 ternary plot shows prioritization items inside a triangle. Ternary plots show
798 how many low, medium or high priorities are assigned to a prioritization
799 item. The corners of the triangle represent high, medium, and low priority,
800 e.g. if a prioritization item has received mostly high priority values then it
801 is shown closer to the high priority corner.

802 *7.2.1. Problems with Compositional Data Analysis in Primary Studies*

803 A few primary studies, as revealed by the systematic review, have prob-
804 lems with the analysis of compositional data.

805 In [50, 7] standard PCA is performed without applying log-ratio trans-
806 formations to compositional data. According to [52], this is likely to be
807 inadequate and in [53], a more appropriate method for performing PCA of
808 compositional data is shown.

809 The normality of compositional data is defined in [54]. It is stated that
810 compositional data must first be transformed using isometric log-ratio trans-
811 formation before the tests for normality can be applied. [47] violates this re-
812 quirement by applying the Shapiro-Wilk test for normality to untransformed
813 compositional data.

814 The Kruskal-Wallis test is used in [47] to analyze compositional data.
815 The test is used to evaluate the difference between three organization levels.
816 The Kruskal-Wallis test assumes that variables within each sample are in-
817 dependent [55]. However, values within compositional data vectors are not
818 independent (as described in Section 2.4). Hence, we claim the Kruskal-
819 Wallis test to be somewhat misused in [47].

820 *7.3. Identifying Prioritization Items with Equal Priority Using ECV (RQ 3)*

821 This section presents the results of applying ECV to the industrial and
822 academic CV data as found through the systematic literature review. Six
823 primary studies included the raw prioritization results in the paper itself or
824 referenced online sources where the data was available. To collect the data
825 from the remaining 34 papers, the authors of all papers were contacted.

826 First, the email addresses provided in the papers were used. If no answer
827 was received authors were searched for using Google, Facebook and LinkedIn.
828 Authors from 11 papers provided us with data to be used in the evaluation
829 of ECV. However, due to confidentiality reasons we can not publish this data
830 directly and instead urge interested parties to contact the authors directly.

831 In short, ECV was applied to 27 CV prioritization cases from 14 studies.
832 In the cases of HCV, ECV was applied two times to the same data to test
833 both compensated and uncompensated priorities. Equal items were detected
834 in three prioritization cases. A summary of the results of is presented in
835 Table 14.

836 In [46] a prioritization of requirement understandability criteria is pre-
837 sented. ECV shows that from the viewpoint of academia researchers, devel-
838 opment have the same importance as product planning (i.e. making strate-
839 gic product planning decisions and perform release planning while choosing
840 which requirements to dismiss).

Table 14: Identified groups of equal items.

Paper identifier & Description	Type of CV	Pairs of equal items	Groups of equal items
Barney et al. [41] Perceived priorities of software product investments in an ideal situation	comp. HCV	(A2, B4) (B4, B5) (B4, C1) (B5, B15) (B6, B7) (B7, B8) (B14, B15) (B14, B18) (B17, B18)	(A2, B4) (B4, C1) (B5, B15) (B6, B7) (B14, B15) (B17, B18)
	uncomp. HCV	(B4, B5) (B4, B8) (B5, B15) (B6, B7) (B7, B12) (B14, B15) (B14, B18) (B16, B17) (B12, B13)	(B4, B5) (B5, B15) (B6, B7) (B14, B15) (B16, B17) (B12, B13)
Berander and Svahnberg [17] Software requirements for course management system	uncomp. & comp. HCV	(3:2, 3:3)	(3:2, 3:3)
Svahnberg et al. [46] The view of academia researchers on the requirements understandability criteria	CV	(Development, Verification & Validation) (Development, Product Planning 1)	(Development, Product Planning 1)

841 A prioritization of software requirements for an academic course man-
842 agement system is presented in [17]. ECV detected that two features—
843 Assignment Submission and Assignment Feedback—have the same priority.
844 In [41] software product investments are prioritized with HCV. The re-
845 sults of ECV was different for uncompensated and compensated HCV results.
846 When compensated HCV was used ECV detected equal items that belong to
847 different high level prioritization groups (A , B and C). Whereas, in case of
848 uncompensated HCV all equal items belong to one high level prioritization
849 group (group B).

850 8. Discussion and Conclusions

851 This section discusses the results of the systematic review and evaluation
852 of ECV conducted as part of this study.

853 CV has been applied in various areas, but most frequently in requirements
854 prioritization and release planning, and quite often also as part of research
855 methodologies. A large part of the studies have been conducted in Sweden,
856 at Ericsson AB. One can see a slight increase in the interest in CV. During
857 the last five years there have been more studies that use CV than between,
858 say, year 2000–2005.

859 Overall, studies that use CV or analyze the results of CV have a high
860 quality in terms of correctness of research process and study realism. How-
861 ever, very few studies present prioritization of more than 30 items and the
862 availability of research data is somewhat limited. In our particular case we
863 were able to obtain data from 43% of the primary studies.

864 8.1. Implications for Practitioners

865 The results of this study provide decision support for industry practi-
866 tioners. We believe that a collection of state of the practice studies help
867 the adoption of CV prioritization method. (Top studies are summarized in
868 Table 12.) In addition, a set of CV analysis methods enables comprehen-
869 sive understanding of the prioritization results. (The analysis methods are
870 presented in Table 13.) One of the most common goals of CV analysis is to
871 display of the prioritization results and, thus, to show the difference between
872 several prioritization perspectives.

873 Additionally, we present ECV—a novel method for CV analysis. Priori-
874 tization often results in the assignment of similar priorities to several prior-
875 itization items. ECV identifies prioritization items with similar priority and
876 tests whether these items can be considered equal. In this case, ECV can
877 be used in software release planning. For example, let us suppose that a set

878 of software requirements are prioritized with regard to the implementation
879 costs. First of all, ECV can then detect items with equal cost. Second, the
880 equal items can be freely swapped between the releases. Finally, the deci-
881 sion to allocate a requirement to a particular release can be made based on
882 another criteria, such as risk or business value.

883 ECV has been successfully applied on a considerable amount of CV data
884 and, additionally, has also detected equal items in different groups of HCV
885 hierarchies.

886 8.2. *Implications for Academia*

887 In the systematic review 36% of papers were revealed by the snowball
888 sampling. That is a considerable amount. Several studies do not mention the
889 name of the prioritization method (i.e. cumulative voting or hundred dollar
890 test). Others are not available through selected databases because they are
891 conference publications or theses. It shows, in our opinion, that snowball
892 sampling ought to be used in all systematic literature reviews.

893 CV results are a special type of data—compositional data. Standard
894 statistical analysis methods that assume the independence of the samples
895 cannot be applied to CV results. In [56] methods for the analysis of com-
896 positional data have been presented. The systematic review conducted as a
897 part of this study revealed that 22 studies analyze CV results; yet, only one
898 study uses compositional data analysis methods, i.e. [10].

899 The small use of compositional data analysis is really not surprising, since
900 literature describing CV does not state that the results are compositional
901 data. Standard statistical analysis methods may produce useful results for
902 compositional data. However, there are cases when they are misleading or
903 even faulty. Section 7.2.1 contains evidence of inappropriate use of statistical
904 methods by several papers.

905 This study has collected a set of compositional data analysis methods for
906 CV analysis (see Table 13). We believe that this could help researchers to
907 improve the analysis of CV results with appropriate methods.

908 Since CV is associated with compositional data, it might be tempting to
909 choose another requirements prioritization method. However, it would not
910 solve the problem *per se*, because any ratio scale prioritization, for instance
911 AHP, contains compositional data.

912 The principal implications for the academia are mainly the following:

- 913 1. All systematic literature reviews should include snowball sampling.
- 914 2. Researchers can improve their statistical analysis of CV results using
915 compositional data analysis methods collected and developed by this
916 study.

917 3. When CV or any other ratio scale prioritization method is taught,
918 compositional data analysis should also be presented as part of the
919 solution.

920 8.3. *Validity Threats*

921 The validity of the systematic review is mainly limited by the chosen
922 databases, the design of the review, and human judgement in study selection
923 and data extraction.

924 To mitigate the threats we use the most popular databases in the field
925 of software engineering. In the beginning of the systematic review a re-
926 view protocol was developed, peer-reviewed, and revised. Search strategy
927 was validated against a set of previously known papers obtained from other
928 researchers. One of many terms used to name cumulative voting is ‘\$100
929 method’. We were not able to search for this term because non of the cho-
930 sen databases support search for special characters like ‘\$’ and the search
931 string ‘100 method’ yields hundreds of thousands of results. To increase the
932 likelihood of discovering relevant studies snowball sampling was extensively
933 used.

934 To increase the validity of study selection, all included studies and 20
935 randomly selected excluded studies were examined by two researchers. There
936 were no disagreement on the inclusion/exclusion of the studies.

937 The large number of studies identified by snowball sampling (15 out of
938 40 studies) may be caused by faulty design or by faulty execution of the
939 search in the databases. There are several reasons why the studies revealed
940 by snowball sampling are not revealed by the search in databases. Reason
941 for each study is given in Table Appendix A.2. Based on these reasons we
942 argue that snowball sampling does not indicate any problems with the design
943 of the search in the databases.

944 Four studies were not found because they were not available through
945 databases used in this systematic review. Out of them one is a master
946 thesis, two are conference publications and one is a publication in the area
947 of forestry. Seven studies do not mention the name of the prioritization
948 method (i.e. hundred dollar method or cumulative voting). Only phrases
949 like “distribution of a predefined amount of fictitious money (\$100,000) over
950 the items to be prioritized” or “1,000 points” allowed us to identify that CV
951 was indeed used. One paper used a previously unknown name for CV, i.e.
952 the 100-point technique.

953 The quality of the data extraction and quality evaluation was validated
954 using inter-rater agreement analysis. In our case, 10% of the studies were

955 rated by two researchers and Krippendorff's alpha was calculated. The agree-
956 ment for the data extraction results was 0.86 and the agreement for the
957 quality evaluation was 0.73 (indicating a credible level of quality).

958 There are two main validity threats with ECV itself. First, ECV may not
959 detect prioritization items with equal priority. Second, ECV may produce a
960 false positive result. There may be a real difference between items that ECV
961 claims as being equal.

962 To mitigate the first threat ECV was applied on artificially created test
963 data with and without items with similar priority. ECV worked correctly in
964 both cases.

965 To mitigate the second threat we visually inspected the results of the
966 application of ECV on the real world data from the primary studies. We
967 concluded that items identified by ECV can be considered equal.

968 CV results used in the evaluation of ECV were tested for normality. The
969 tests indicated that CV results are not normally distributed. Therefore, the
970 design of ECV was based on a non-parametric statistical test.

971 8.4. Future Research

972 There are very few studies that apply CV on prioritization sets of more
973 than 30 items. However, in requirements engineering, industry practitioners
974 need to prioritize much larger numbers of software requirements. Therefore,
975 the state of art could benefit from the application of CV and HCV to large
976 prioritization sets.

977 The proposed method, ECV, has now been evaluated on existing research
978 data. To further evaluate the ECV, it could be applied in direct industry
979 practice and in prioritization cases with a larger number of prioritization
980 items. Additionally, compositional data analysis methods, as the ones iden-
981 tified by this paper, should be tried with other prioritization methods that
982 produce ratio scale results.

983 8.5. Conclusions

984 CV prioritization results are special type of data – compositional data.
985 Any analysis of CV results must take into account the compositional nature
986 of the CV results.

987 This study presents a systematic literature review of the empirical use
988 of CV. CV has been applied in various areas, but most frequently in re-
989 quirements prioritization and release planning. The review has resulted in
990 a collection of state of the practice studies and CV result analysis methods.
991 We believe that it can help the adoption of CV prioritization method.

992 In our case, snowball sampling was performed as a part of the review.
993 Since it revealed 36% out of all primary studies, we believe that in future
994 snowball sampling should be used in all systematic reviews.

995 Additionally, we present ECV—a novel method for CV analysis. As
996 suggested by our evaluation, ECV is able to detect prioritization items with
997 equal priority (i.e. items that have insignificant difference in priority). The
998 evaluation of ECV was based on the data obtained from the authors of the
999 primary studies.

1000 References

- 1001 [1] P. Berander, A. Andrews, Requirements Prioritization, in: A. Aurum,
1002 C. Wohlin (Eds.), Engineering and Managing Software Requirements,
1003 Springer-Verlag, Berlin/Heidelberg, 2005, 2005, pp. 69–94.
- 1004 [2] D. Leffingwell, D. Widrig, Managing software requirements: A unified
1005 approach (1999) 118–119.
- 1006 [3] V. Ahl, An experimental comparison of five prioritization methods,
1007 Master’s Thesis, School of Engineering, Blekinge Institute of Technology
1008 (2005).
- 1009 [4] P. Berander, P. Jonsson, Hierarchical Cumulative Voting (HCV) prior-
1010 itization of requirements in hierarchies, 2006.
- 1011 [5] J. Karlsson, K. Ryan, A cost-value approach for prioritizing require-
1012 ments, IEEE Software 14 (1997) 67–74.
- 1013 [6] J. Karlsson, An evaluation of methods for prioritizing software require-
1014 ments, Information and Software Technology 39 (1998) 939–947.
- 1015 [7] F. Pettersson, M. Ivarsson, T. Gorschek, P. Öhman, A practitioner’s
1016 guide to light weight software process assessment and improvement plan-
1017 ning (2008).
- 1018 [8] S. Barney, C. Wohlin, Software Product Quality: Ensuring a Common
1019 Goal, in: Q. Wang, V. Garousi, R. Madachy, D. Pfahl (Eds.), Trust-
1020 worthy Software Development Processes, volume 5543 of *Lecture Notes*
1021 *in Computer Science*, Springer Berlin Heidelberg, Berlin, Heidelberg,
1022 2009, 2009, pp. 256–267.

- 1023 [9] P. Jönsson, C. Wohlin, A study on prioritisation of impact analysis
1024 issues: A comparison between perspectives, *Software Engineering Re-*
1025 *search and Practice in Sweden* (2005).
- 1026 [10] P. Chatzipetrou, L. Angelis, P. Rovegard, C. Wohlin, Prioritization of
1027 Issues and Requirements by Cumulative Voting: A Compositional Data
1028 Analysis Framework, 2010, pp. 361–370.
- 1029 [11] R. Engstrom, Cumulative Voting as a Remedy for Minority Vote Dilu-
1030 tion, *Local Government Election ...* (1999).
- 1031 [12] S. Bhagat, J. Brickley, Cumulative voting: The value of minority share-
1032 holder voting rights, *Journal of Law and Economics* (1984).
- 1033 [13] V. Hiltunen, J. Kangas, J. Pykalainen, Voting methods in strategic
1034 forest planning - Experiences from Metsähallitus, *Forest Policy and*
1035 *Economics* 10 (2008) 117–127.
- 1036 [14] P. Boldi, F. Bonchi, C. Castillo, S. Vigna, Voting in social networks,
1037 *CIKM '09*, ACM Press, New York, New York, USA, 2009.
- 1038 [15] H. Ayad, M. Kamel, Cumulative Voting Consensus Method for Parti-
1039 tions with Variable Number of Clusters, *Pattern Analysis and Machine*
1040 *Intelligence*, *IEEE Transactions on* 30 (2008) 160–173.
- 1041 [16] M. Svahnberg, A. Karasira, A Study on the Importance of Order in
1042 Requirements Prioritisation, *IEEE*, 2009.
- 1043 [17] P. Berander, M. Svahnberg, Evaluating two ways of calculating priorities
1044 in requirements hierarchies - An experiment on hierarchical cumulative
1045 voting, 2009.
- 1046 [18] T. Saaty, The analytic hierarchy process., McGraw-Hill, New York
1047 (1980).
- 1048 [19] S. Brenner, J. Schwalbach, Legal Institutions, Board Diligence, and
1049 Top Executive Pay, *Corporate Governance: An International Review*
1050 17 (2009) 1–12.
- 1051 [20] V. Pawlowsky-Glahn, J. J. Egozcue, Compositional data and their anal-
1052 ysis: an introduction, *Geological Society, London, Special Publications*
1053 264 (2006) 1–10.

- 1054 [21] J. Martin-Fernandez, C. Barceló-Vidal, V. Pawlowsky-Glahn, Dealing
1055 with zeros and missing values in compositional data sets using nonpara-
1056 metric imputation, *Mathematical Geology* 35 (2003) 253–278.
- 1057 [22] P. Filzmoser, K. Hron, Outlier detection for compositional data using
1058 robust methods *Outlier Detection for Compositional Data Using Robust*
1059 *Methods, Analysis and Applications* (2007).
- 1060 [23] K. Khan, A systematic review of software requirements prioritization,
1061 Unpublished master’s thesis, Blekinge Institute of Technology, Ronneby,
1062 Sweden (2006).
- 1063 [24] F. Zahedi, The analytic hierarchy process: a survey of the method and
1064 its applications, *Interfaces* (1986) 96–108.
- 1065 [25] P. Runeson, M. Höst, Guidelines for conducting and reporting case
1066 study research in software engineering, *Empirical Software Engineering*
1067 14 (2008) 131–164.
- 1068 [26] L. Goodman, Snowball sampling, *The Annals of Mathematical Statis-*
1069 *tics* (1961).
- 1070 [27] K. Krippendorff, Bivariate agreement coefficients for reliability of data,
1071 *Sociological methodology* (1970).
- 1072 [28] K. Krippendorff, *Content analysis: An introduction to its methodology*
1073 (2004).
- 1074 [29] B. Kitchenham, Guidelines for performing systematic literature reviews
1075 in software engineering, *Engineering* (2007).
- 1076 [30] P. Berander, P. Jönsson, A goal question metric based approach for effi-
1077 cient measurement framework definition, *ACM*, Rio de Janeiro, Brazil,
1078 2006, pp. 316–325.
- 1079 [31] B. Regnell, M. Höst, J. och Dag, An industrial case study on distributed
1080 prioritisation in market-driven requirements engineering for packaged
1081 software, *Requirements ...* (2001).
- 1082 [32] B. Kitchenham, Procedures for performing systematic reviews, Keele,
1083 UK, Keele University 33 (2004).
- 1084 [33] M. Ivarsson, T. Gorschek, A method for evaluating rigor and industrial
1085 relevance of technology evaluations, *Empirical Software Engineering*
1086 (2010) 1–31.

- 1087 [34] C. Wohlin, P. Runeson, M. Höst, Experimentation in software engineer-
1088 ing: an introduction, Springer Netherlands, 2000.
- 1089 [35] A. Jedlitschka, D. Pfahl, Reporting guidelines for controlled experi-
1090 ments in software engineering, in: 2005 International Symposium on
1091 Empirical Software Engineering, 2005., IEEE, 2005, p. 10.
- 1092 [36] K. Rinkevics, Data Extraction and Quality Evaluation results, 2011.
- 1093 [37] R. Ihaka, R. Gentleman, R: a language for data analysis and graphics,
1094 Journal of computational and graphical statistics (1996) 299–314.
- 1095 [38] K. Rinkevics, ECV implementation source code, 2011.
- 1096 [39] R. M. Groves, F. J. Fowler, M. P. Couper, J. M. Lepkowski, E. Singer,
1097 Survey methodology, John Wiley and Sons, 2009.
- 1098 [40] S. Barney, A. Aurum, C. Wohlin, The Relative Importance of Aspects
1099 of Intellectual Capital for Software Companies, in: 2009 35th Euromicro
1100 Conference on Software Engineering and Advanced Applications, IEEE,
1101 2009, 2009, pp. 313–320.
- 1102 [41] S. Barney, C. Wohlin, A. Aurum, Balancing software product invest-
1103 ments, IEEE Computer Society, 2009, pp. 257–268.
- 1104 [42] S. Barney, A. Aurum, C. Wohlin, A product management challenge:
1105 Creating software product value through requirements selection, Jour-
1106 nal of Systems Architecture 54 (2008) 576–593.
- 1107 [43] S. Laukkanen, T. Palander, J. Kangas, A. Kangas, Evaluation of the
1108 multicriteria approval method for timber-harvesting group decision sup-
1109 port, Silva Fennica 39 (2005) 249–264.
- 1110 [44] G. Hu, Adding value to software requirements: An empirical study in
1111 the chinese software industry, Seventeenth Australian Conference on
1112 ... (2006).
- 1113 [45] R. Feldt, R. Torkar, E. Ahmad, B. Raza, Challenges with Software
1114 Verification and Validation Activities in the Space Industry, IEEE, 2010.
- 1115 [46] M. Svahnberg, T. Gorschek, M. Eriksson, A. Borg, K. Sandahl,
1116 J. Börstler, A. Loconsole, Perspectives on Requirements Understand-
1117 ability – For Whom Does the Teacher’s Bell Toll?, IEEE, 2008.

- 1118 [47] P. Jönsson, C. Wohlin, Understanding impact analysis: An em-
 1119 pirical study to capture knowledge on different organisational levels,
 1120 ... Conference on Software Engineering and Knowledge ... (2005).
- 1121 [48] L. a. Kuzniarz, Empirical extension of a classification framework for
 1122 addressing consistency in model based development, Information and
 1123 Software Technology (2010).
- 1124 [49] P. Rovegard, L. Angelis, C. Wohlin, An Empirical Study on Views of
 1125 Importance of Change Impact Analysis Issues, Software Engineering,
 1126 IEEE Transactions on 34 (2008) 516–530.
- 1127 [50] C. Wohlin, A. Aurum, Criteria for selecting software requirements to
 1128 create product value: An industrial empirical study, Value-Based Soft-
 1129 ware Engineering (2006).
- 1130 [51] B. Regnell, M. Höst, J. Natt, Visualization of Agreement and Satisfac-
 1131 tion in Distributed Prioritization of Market Requirements, Chart (2000)
 1132 1–12.
- 1133 [52] J. Aitchison, Principal component analysis of compositional data,
 1134 Biometrika 70 (1983) 57.
- 1135 [53] P. Filzmoser, K. Hron, C. Reimann, F. Sm, P. Filzmoser, K. Hron,
 1136 C. Reimann, Principal component analysis for compositional data with
 1137 outliers Principal component analysis for compositional data with out-
 1138 liers, Analysis and Applications (2007).
- 1139 [54] V. Pawlowsky Glahn, J. Egozcue, R. Tolosana Delgado, Lecture notes
 1140 on compositional data analysis, Interpretation A Journal Of Bible And
 1141 Theology (2007).
- 1142 [55] W. Kruskal, W. Wallis, Use of ranks in one-criterion variance analysis,
 1143 Journal of the American statistical Association 47 (1952) 583–621.
- 1144 [56] J. Aitchison, The statistical analysis of compositional data, Chapman
 1145 & Hall, London, 1986.
- 1146 [57] D. Baca, K. Petersen, Prioritizing Countermeasures through the Coun-
 1147 termeasure Method for Software Security (CM-Sec), in: M. Ali Babar,
 1148 M. Vierimaa, M. Oivo (Eds.), Product-Focused Software Process Im-
 1149 provement, volume 6156 of *Lecture Notes in Computer Science*, Springer
 1150 Berlin / Heidelberg, 2010, 2010, pp. 176–190.

- 1151 [58] S. a. b. Bowler, Election systems and voter turnout: Experiments in
1152 the United States, *Journal of Politics* 63 (2001) 902–915.
- 1153 [59] D. Brockington, A Low Information Theory of Ballot Position Effect,
1154 *Political Behavior* 25 (2003) 1–27.
- 1155 [60] D. Cooper, A. Zillante, A comparison of cumulative voting and gener-
1156 alized plurality voting, *Public Choice* (2010).
- 1157 [61] N. D. Fogelström, M. Svahnberg, T. Gorschek, Investigating Impact of
1158 Business Risk on Requirements Selection Decisions, *IEEE*, 2009.
- 1159 [62] S. Hatton, Choosing the Right Prioritisation Method, in: *Proceed-*
1160 *ings of the 19th Australian Conference on Software Engineering*, IEEE
1161 Computer Society, Washington, 2008, pp. 517–526.
- 1162 [63] S. Hatton, Early prioritisation of goals, in: *Proceedings of the 2007*
1163 *conference on Advances in conceptual modeling: foundations and appli-*
1164 *cations*, ER’07, Springer-Verlag, Berlin, 2007, pp. 235–244.
- 1165 [64] V. Heikkilä, A. Jadallah, K. Rautiainen, G. Ruhe, Rigorous Support
1166 for Flexible Planning of Product Releases - A Stakeholder-Centric Ap-
1167 proach and Its Initial Evaluation, *IEEE*, 2010.
- 1168 [65] M. Staron, C. Wohlin, An Industrial Case Study on the Choice Be-
1169 tween Language Customization Mechanisms, in: J. Münch, M. Vier-
1170 imaa (Eds.), *Product-Focused Software Process Improvement*, volume
1171 4034 of *Lecture Notes in Computer Science*, Springer Berlin / Heidel-
1172 berg, 2006, 2006, pp. 177–191.
- 1173 [66] T. Touseef, C. Gancel, A structured goal based measurement framework
1174 enabling traceability and prioritization, ... (ICET), 2010 6th Interna-
1175 tional Conference on (2010).
- 1176 [67] P. Berander, C. Wohlin, Differences in views between development
1177 roles in software process improvement-a quantitative comparison, in:
1178 *Proceedings 8th Conference on Empirical Assessment in Software Engi-*
1179 *neering*, 2004.
- 1180 [68] P. Berander, Using students as subjects in requirements prioritization,
1181 *Proceedings. 2004 International Symposium on Empirical Software En-*
1182 *gineering*, 2004. ISESE ’04. (2004) 167–176.

- 1183 [69] P. Berander, C. Wohlin, Identification of Key Factors in Software Pro-
1184 cess Management-A Case Study (2003).
- 1185 [70] R. L. Cole, D. a. Taebel, R. L. Engstrom, Cumulative Voting in a Munic-
1186 ipal Election: A Note on Voter Reactions and Electoral Consequences,
1187 The Western Political Quarterly 43 (1990) 191.
- 1188 [71] J. Kuklinski, Cumulative and Plurality Voting: An Analysis of Illinois’
1189 Unique Electoral System, The Western Political Quarterly 26 (1973)
1190 726–746.
- 1191 [72] S. Laukkanen, T. Palander, J. Kangas, Applying voting theory in par-
1192 ticipatory decision support for sustainable timber harvesting, Canadian
1193 Journal of Forest Research 34 (2004) 1511–1524.
- 1194 [73] J. Sawyer, D. MacRae, Game theory and cumulative voting in Illinois:
1195 1902-1954, The American Political Science Review 56 (1962) 936–946.

1196 Appendix A. Primary Studies

1197 Appendix A.1. Primary studies found during search in databases.

	Title	Reference
	Prioritizing countermeasures through the countermeasure method for software security (CM-Sec)	Baca and Petersen [57]
	The relative importance of aspects of intellectual capital for software companies	Barney et al. [40]
	Software product quality: Ensuring a common goal	Barney and Wohlin [8]
	Balancing software product investments	Barney et al. [41]
	Hierarchical cumulative voting (HCV) prioritization of requirements in hierarchies	Berander and Jonsson [4]
	A goal question metric based approach for efficient measurement framework definition	Berander and Jönsson [30]
	Evaluating two ways of calculating priorities in requirements hierarchies: An experiment on hierarchical cumulative voting	Berander and Svahnberg [17]
	Election systems and voter turnout: Experiments in the United States	Bowler [58]
	A low information theory of ballot position effect	Brockington [59]
	Prioritization of issues and requirements by cumulative Voting: A compositional data analysis framework	Chatzipetrou et al. [10]
	A comparison of cumulative voting and generalized plurality voting	Cooper and Zillante [60]
	Challenges with software verification and validation activities in the space industry	Feldt et al. [45]
1198	Investigating impact of business risk on requirements selection decisions	Fogelström et al. [61]
	Choosing the right prioritization method	Hatton [62]
	Early prioritization of goals	Hatton [63]
	Rigorous support for flexible planning of product releases: A stakeholder-centric approach and its initial evaluation	Heikkilä et al. [64]
	Voting methods in strategic forest planning: Experiences from Metsähallitus	Hiltunen et al. [13]
	Empirical extension of a classification framework for addressing consistency in model based development	Kuzniarz [48]
	Evaluation of the multi-criteria approval method for timber-harvesting group decision support	Laukkanen et al. [43]
	A practitioner’s guide to light weight software process assessment and improvement planning	Pettersson et al. [7]
	An empirical study on views of importance of change impact analysis issues	Rovegard et al. [49]
	An industrial case study on the choice between language customization mechanisms	Staron and Wohlin [65]
	Perspectives on requirements understandability—For whom does the teacher’s bell toll?	Svahnberg et al. [46]
	A study on the importance of order in requirements prioritization	Svahnberg and Karasira [16]
	A structured goal based measurement framework enabling traceability and prioritization	Touseef and Gancel [66]

1199 *Appendix A.2. Primary studies revealed by snowball sampling.*

1200

Reference	Title	Reason why the paper is not revealed by the search in databases
Ahl [3]	An experimental comparison of five prioritization methods	Selected databases does not contain the paper, master thesis at BTH
Barney et al. [42]	A product management challenge: Creating software product value through requirements selection	Prioritization method name not mentioned, phrase "1,000 points" used instead.
Berander and Wohlin [67]	Differences in views between development roles in software process improvement—A quantitative comparison	Prioritization method name not mentioned, phrase "100 points" used instead.
Berander [68]	Using students as subjects in requirements prioritization	Unknown CV name: 100-point technique
Berander and Wohlin [69]	Identification of key factors in software process management: A case study	Prioritization method name not mentioned, phrase "100 points" used instead.
Cole et al. [70]	Cumulative voting in a municipal election: A note on voter reactions and electoral consequences	Study published before year 2001.
Hu [44]	Adding value to software requirements: An empirical study in the chinese software industry	Prioritization method name not mentioned, phrase "1,000 points" used instead.
Jönsson and Wohlin [9]	A study on prioritization of impact analysis issues: A comparison between perspectives	Selected databases does not contain the paper.
Jönsson and Wohlin [47]	Understanding impact analysis: An empirical study to capture knowledge on different organizational levels	Selected databases does not contain the paper.
Kuklinski [71]	Cumulative and plurality voting: An analysis of Illinois' unique electoral system	Study published before year 2001.
Laukkanen et al. [72]	Applying voting theory in participatory decision support for sustainable timber harvesting	Selected databases does not contain the paper.
Regnell et al. [31]	An industrial case study on distributed prioritization in market-driven requirements engineering for packaged software	Prioritization method name not mentioned: "distribution of a predefined amount of fictitious money (\$100,000) over the items to be prioritized."
Regnell et al. [51]	Visualization of agreement and satisfaction in distributed prioritization of market requirements	Prioritization method name not mentioned: "distribution of a predefined amount of fictitious money (\$100,000) over the items to be prioritized."
Sawyer and MacRae [73]	Game theory and cumulative voting in Illinois: 1902–1954	Study published before year 2001.
Wohlin and Aurum [50]	Criteria for selecting software requirements to create product value: An industrial empirical study	Prioritization method name not mentioned: "The subjects had 1,000 points to spend among the 13 criteria."

Appendix B. CV Result Analysis Methods

	Paper																					
	Svahnberg2008	Svahnberg2009	Staron2006	Pettersson2008	Wohlin2006	Laukkainen2005a	Hu2006	Jonsson2005a	Kuzniarz2010	Rowigard2008	Bernardie2006a	Berander2004a	Bernardie2006	Feldt2010	Barney2009b	Barney2008	Barney2009a	Barney2009	Jonsson2005	Chatzipetrou2010	Regnell2001	Regnell2000
Analysis method																						
Table that shows final priorities	x			x												x						
Chart that shows final priorities	x			x	x	x																
Table of top-5 prioritization items																						
min , max , \bar{x} , \bar{x} and σ of priorities assigned by different stakeholders									x													
Bar chart of prioritization results showing \bar{x} priority and σ of priorities									x	x												
Pearson correlation coefficient		x										x										
Nemenyi Damico Wolfe Dunn														x								
Spearman's r															x							
Kruskal-Wallis								x							x		x					
Wilcoxon							x															
Correlation matrix		x													x		x					
Chart for comparing priorities from two perspectives, priorities are points in two dimensional plane, x - and y -axis represent two different perspectives										x									x			
Difference between priorities assigned by each two stakeholders using χ^2 -statistic										x												
Median ranks		x																				
CV results converted to priority ranks		x											x						x			
PCA				x	x															x		
Percentage of divergence of priorities assigned by a stakeholder											x											
Average percentage of items given non-zero value											x											
Distribution chart																				x	x	
Disagreement chart				x																x	x	
Satisfaction chart			x																	x	x	
Bi-plot																				x		
Ternary plot																				x		

Appendix C. Quality Evaluation Checklist

	Item	Question or Description of the Item	Rating
1.	Background, introduction	Introduce research area	
2.	Problem statement, purpose	What is the problem [35]? Where does it occur [35]? Who has observed it [35]? Why is it important to be solved [35]?	
3.	Context, independent variables (aka. environment, setting)	Study location, time constraints, application domain, organization, tools, market, process (e.g. software development methodology), size of project, product that is being developed	
4.	Related work	Other existing work, alternative technologies, solutions, and studies	
5.	Goals and Hypotheses	Null hypothesis and one or more alternative hypotheses for each goal	
6.	Research questions		
7.	Design, Research methods		
7.1.	Design	Description of each step of the study	
7.2.	Control group	If there is a control group, are participants similar to the treatment group participants in terms of variables that may affect study outcomes[29]?	
7.3.	Randomization	Random selection of participants and objects Random assignment of treatment and objects to participants Random order of treatments in case of paired design. If each participant is assigned two treatments A and B, then part of participants perform A first and the other part start with B	
7.4.	Blocking	Group participants of the study into homogeneous groups called blocks (e.g. students in one course, database developers in one company) and implement the study design within each block independently. The idea is that variability of independent variables (e.g. experience and knowledge of subjects) is smaller within a group. That helps measuring changes in dependent variables [32].	
7.5.	Balancing	Equal number of subjects should be assigned to each treatment [32].	
7.6.	Blinding	Automated assignment of treatments to subjects [32] Automated distribution of study materials to subjects [32] Persons who grade the task results should not know which treatment was used [32] Analyst should not know which treatment group is which [32] Automated data collection from subjects [32]	
8.	Subjects (participants)		
8.1.	Population		
8.2.	Sampling	How sampling is performed? What subjects are included and excluded? [29] What is the type of the sampling (e.g. convenience, random)? Is the sample(selected participants) representative of the population?	
8.3.	"Drop outs" and response rate	Are reasons given for refusal to participate[29]?	
8.4.	Subject motivation	E.g. material benefits, course credits for students, etc.	
9.	Objects	E.g. documents and other artifacts	
10.	Measures, Data collection procedures	Who, when, and how to measure [29]? How is the measurement supported? Is it automated [29]? Are the measures used in the study the most relevant ones for answering the research questions [29]?	
11.	Analysis procedure		
11.1.	Data description	Do the numbers add up across different tables and subgroups [29]?	
11.2.	Data types (continuous, ordinal, categorical)		
11.3.	Scoring systems		
11.4.	Data set reduction, outliers		
11.5.	Statistical methods	Are the assumptions of statistical methods met? What statistical programs are used?	
11.6.	Statistical significance	If statistical tests are used to determine differences, is the actual p-value given [29]? If the study is concerned with differences among groups, are confidence limits given describing the magnitude of any observed differences [29]?	
12.	Validity threats	Threats, implications of the threats, and threat mitigation	
12.1.	Side-effects during study execution	Deviations from the plan, solutions for the deviations	
13.	Most important findings	Are all study questions answered [29]? Are negative findings presented [29]?	
14.	Industry impact, inference, generalization	What implications does the report have for practice [29]? How and where the results can be used? Limitations under which findings are relevant [35]?	
15.	Future work		