
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].

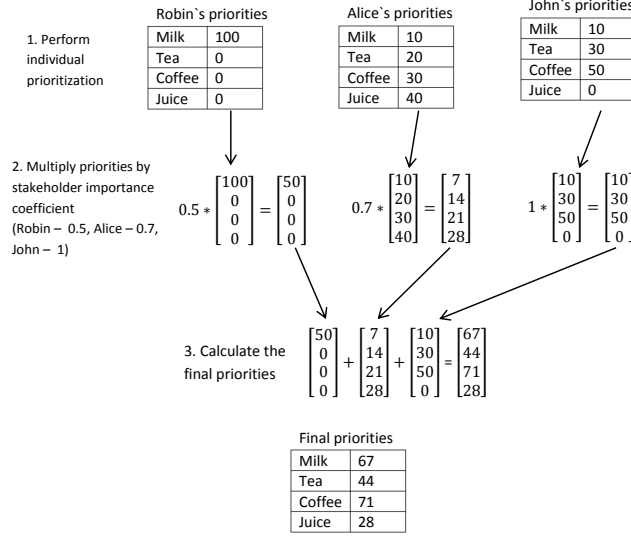


Figure 1: Example of CV with several stakeholders.

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

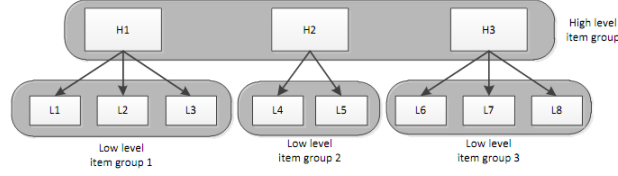


Figure 2: Example of prioritization item hierarchy.

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

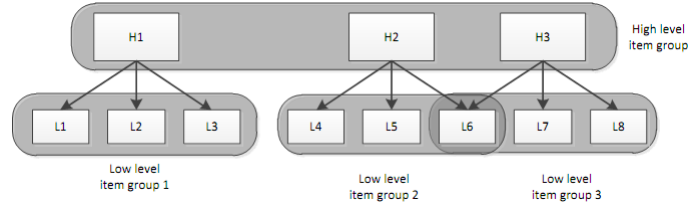


Figure 3: Overlapping Prioritization Item Hierarchy Example

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. Apply the compensation factor to all low level requirements as described in Section 2.3.5.
6. 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.

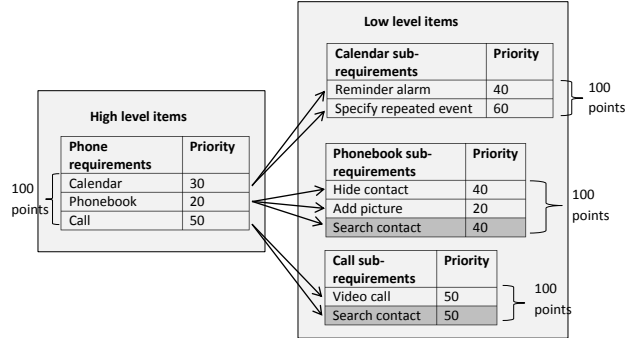


Figure 4: Example of Hierarchical Cumulative Voting, Requirement Hierarchy

Table 1: Example of Hierarchical Cumulative Voting

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

2.4. Compositional Data Analysis

CV results can be seen as a special type of data—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 lower relative priority. This also
 293 holds true the other way around; if the item is put in a set with low priority
 294 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.

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

$$Y = a_1X_1 + a_2X_2 + \dots + a_nX_n \quad (2)$$

302 where Y is the product of lineal combination and a_i is any real number.
 303 Since the sum of priorities assigned in CV must add up to 100 (or any
 304 other constant number) at least one linear combination of X is not normally
 305 distributed because it must always add up to 100:

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

306 In our opinion, the above shows that CV results do not follow a multivari-
 307 ate normal distribution and, hence, means that they must not be analysed
 308 using parametric statistical tests [20].

309 2.4.1. Problem of Zeroes

310 Compositional data analysis requires that ratios between any compo-
 311 nents in a vector can be computed. Computing a ratio with a zero value
 312 is meaningless. This causes a problem because CV allows stakeholders to
 313 assign zero priorities to some prioritization items. There are two types of
 314 zeroes in compositional data: essential and rounded.

315 Essential zeroes mean that a data component is not present. Rounded
 316 zeroes mean that the component is present but its value is very low. We can
 317 assume that zeroes in CV results are rounded because the priority of an item
 318 is a completely abstract notion and the instrument for measuring priority is
 319 human judgement [10].

320 Before compositional data analysis can be applied to CV results, we must
 321 first remove zeroes in the data. One approach can be to forbid stakeholders
 322 to assign zero priorities. This approach is used in e.g. [7]. But this can
 323 add some unnecessary complexity to the prioritization process. In [10] the
 324 authors propose the use of a multiplicative replacement strategy (as defined
 325 in [21]) for CV result analysis.

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

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

328 where δ_j is the imputed value and c is the constant sum constraint (the
 329 same as k in equation (1)). In order for the total sum of components to stay
 330 constant the equation subtracts some value from the items with a priority
 331 higher than zero. More is subtracted from components with higher values
 332 than from the components with lower values and the value of the imputed
 333 δ_j is arbitrary.

334 2.4.2. Isometric Log-Ratio Transformation

335 In order to be able to apply standard statistical methods to composi-
 336 tional data it must be transformed to remove the inherent correlation of
 337 the values. Compositional data analysis proposes special transformations
 338 that change the compositional data values to unconstrained real values. One

such transformation is isometric log-ratio (*ilr*) transformation (as proposed by [20, 22]):

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

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 paper 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. This will answer RQ 1 and RQ 2.

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]). Case studies 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.

368 Several studies have used CV empirically in industrial as well as in academic
369 settings. Nevertheless, there are no studies that provide an overview of the
370 current state of the practice in this field. Therefore, before continuing with
371 the refinement of CV and conducting new empirical studies (i.e. case study
372 or experiment), a systematic literature review is required.

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

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

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

390 In short, this study consists of two parts: a systematic literature review
391 of CV and an evaluation of ECV.

392 *4.2. Research Questions*

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

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

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

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

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

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

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

409 5. Systematic Literature Review

410 This section presents the design of the systematic review. For the results
411 of the execution please see Section 7.1 and 7.2.

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

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

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

428 5.1. Data Sources and Search Strategy

429 This systematic literature review is designed based on the guidelines by
430 Kitchenham [29]. First a trial search in electronic databases was conducted.
431 In order to scale the review to a manageable, yet sufficient size, databases
432 were searched with different search strings. Relevant papers that were found
433 during the trial search were used to extract additional search strings. The
434 trial search revealed that the number of studies that use CV is not very large.
435 Therefore, we decided to include not only software engineering studies but
436 also studies in other research areas, such as forestry or corporate governance.

437 Since CV is frequently used in studies without mentioning this in the
438 abstract, full text search in databases is preferable. Unfortunately not all
439 databases support full text search. Full text search was performed in the
440 IEEE Xplore and Springer Link databases. In ACM Digital Library, In-
441 spec/Compendex, ISI Web of Knowledge, and SCOPUS only metadata was

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	

442 searched. Search strings consist of a Boolean expression: (A or B or C or D or
443 E or F or G), where:

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

451 5.2. Study Selection

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

454 Papers were selected in two phases: selection based on metadata and
455 selection based on full text.

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

466 5.3. Quality Evaluation

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

470 5.3.1. Is the Study Right?

471 Study quality obviously depends on the correctness of the study process
472 including planning, operation, analysis and interpretation of the results. The
473 correctness of the process can be measured by evaluating the description of

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	

the study or replicating the study. Thus, to gain the trust of industry practitioners and other researchers, the process of the study must be rigorously described. In short, the description must facilitate replication of the study as well as the presentation of limitations and validity threats.

5.3.2. *Is it the Right Study?*

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

Realism of a study depends on the context, scale, and subjects of the study. The study should be conducted in a **setting** that is similar or equal to the setting in which the findings of the study are intended to be used. Hence, studies that are conducted in an industrial setting are more valuable. The **subjects** of a study should be similar to the people who are supposed to use the findings of the study. The subjects ought to have appropriate work experience, role in the organization, skills, cultural background, motivation, and so forth. The **scale** of a study refers to the size of the study objects. In the case of this systematic review the scale of a study is measured as the number of prioritization items. Study in academia may have a large number of prioritization items. At the same time, an industrial study, with professionals as subjects, may involve a smaller number of prioritization items.

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

The **type** of the research study can be considered as a criterion for the evaluation of study realism. [32] suggest that study designs that are more rigorous (e.g. experiments) are more realistic than observational studies (e.g. case study) due to a higher level of control. On the other hand [33] rate study designs based on other criteria, i.e. how frequently each type of study design is used in an industrial or academic setting. If a study design is used more in an industrial setting, then it is considered more realistic. For instance, in software engineering case studies are frequently used in industrial

513 settings, whereas, experiments are usually performed in academia using stu-
514 dents as subjects. Therefore, [33] argue that case studies are more realistic
515 than formal experiments. Obviously the effect of study design on the study
516 realism may be interpreted in different ways. Therefore, we will not use this
517 parameter in our quality evaluation.

518 The statistical significance of the results of a study can be used to evalu-
519 ate the significance of the study findings. This measure will not be used, be-
520 cause the studies that are evaluated belong to very different research areas.
521 Thus, the significance levels of the findings of the studies are not directly
522 comparable. Additionally, sometimes no result is more interesting than a
523 significant result. If study results does not conform to the expectations of
524 researchers, this may reveal important gaps in existing knowledge. Never-
525 theless, the evaluation of the correctness of the study process verifies that
526 the statistical analysis is performed and significance levels are reported.

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

538 5.3.3. *Rating of the Studies*

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

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

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

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.

552 is counted only once.

553 The availability of raw results of CV is rated separately because it is
554 especially important for our purposes. The data availability rating criteria is
555 given in Table 6. If the results of a study are not available it is not possible
556 to validate the results of the study and, hence, the credibility of the findings
557 is lower. Ideally the data collected in the study should be presented directly
558 in the paper. An alternative may be to make the data freely available online
559 and reference the online source.

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

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

575 As a result of the rating each study was assigned four rating values in

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

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 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 the 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.

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

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

593 6. Equality of Cumulative Votes

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

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

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

612 6.1. Testing Equality of Two Items

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

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

623 ECV, in contrast, uses the second approach. It finds items that are
624 similar and the rest of the items are considered different. This method tests
625 the probability of the difference between the means of two items being smaller
626 than the given value. In short, ECV tests the probability of the means of two

627 prioritization items differing by less than 25%. If the probability is higher
 628 than 70% the items are considered equal.

629 The input to ECV is an $n \times p$ matrix A that contains the raw results of
 630 the prioritization. The columns of the matrix represent prioritization items
 631 while rows represent stakeholders. ECV performs the following operations
 632 for the priorities of each two prioritization items:

- 633 1. Replace zeroes in CV results.
- 634 2. Transform the data using *ilr* transformation.
- 635 3. Determine distribution function using kernel density estimation.
- 636 4. Use the distribution function to find the probability that the difference
 637 between two prioritization items is smaller than 25%.
- 638 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 :

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

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

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

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

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

648 After the data is transformed into log-ratios statistical test can be ap-
 649 plied. The purpose of the test is to determine what the probability is of the
 650 relative difference between two prioritization items k and l being less than
 651 25%. This means determining the probability of the ratio k/l between the
 652 items k and l as being in the range of $\frac{3}{4}$ to $\frac{4}{3}$. Or in terms of log-ratios
 653 it means determining the probability of $ilr(k, l)$ being between $ilr(3, 4)$ and
 654 $ilr(4, 3)$. Hence, the objective of the test is to determine the probability of
 655 the sample mean (i.e. mean value of C) laying between the two values.

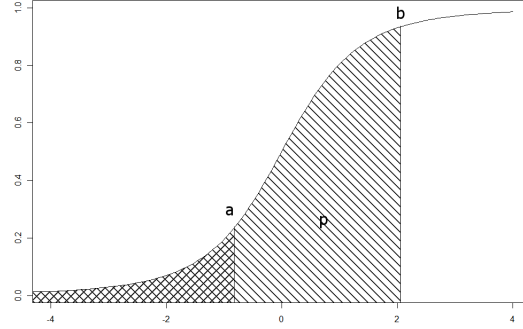


Figure 5: Cumulative distribution function of the ratio k/l between the items k and l (area p denotes probability that k/l is between $\frac{3}{4}$ and $\frac{4}{3}$.)

656 The probability that the mean takes a particular value can be expressed
657 in form of a cumulative distribution function. The probability of the mean
658 being between two values a and b (where a is smaller than b) can be deter-
659 mined by subtracting the probability of the mean being smaller than a from
660 probability of the mean being smaller than b .

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

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

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

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

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

678 After both prioritization items are tested for equality it may be conve-

Table 10: Example of equality table

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

nient to display the equality of different items in the form of a table. Please see Table 10 for an example.

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. 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 Analyse 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 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.

706 The review process was facilitated by the reference management software
707 Mendeley. All results of the study selection are available online and can be
708 obtained by contacting the authors of this paper. For each study we specify
709 keywords and databases that were used to find the study. If a study has
710 been excluded, the exclusion criteria are provided.

711 The number of papers revealed by each search string and database is
712 presented in Table 11. It should be noted that several papers were found
713 by more than one search string or in more than one database. Table 11
714 shows that the search string ‘cumulative voting’ was the most frequently
715 used in research community to denote CV. Therefore, researchers should use
716 or reference this term when talking about CV.

717 To perform snowball sampling we examined the references of primary
718 studies that were found during the database search. References were used
719 to search for the papers in the Google and Google Scholar search engines.
720 Studies that were found in the search and passed the study selection criteria
721 were added to the set of primary studies.

722 After the primary studies were selected, data extraction and quality eval-
723 uation was performed by two researchers. One researcher examined all stud-
724 ies while the second researcher did quality evaluation and data extraction for
725 10% of the studies. The studies were randomly selected. Inter-rater agree-
726 ment were calculated by means of Krippendorff’s alpha coefficient. Agree-
727 ment for data extraction results is 0,86 and agreement for the quality evalu-
728 ation is 0,73. According to Krippendorff [28] it is common to require agree-
729 ment above 0,8 and the lowest acceptable agreement is 0,667. Therefore,
730 we conclude that agreement calculated for this study is sufficient. Ratings
731 of the study setting, correctness, research data availability, and number of
732 prioritization items are presented in Figure 6.

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

737 Figure 7 shows a bubble chart of the distribution of studies over research
738 areas and time. The figure shows that CV was first applied some time
739 ago in research of government elections. Nowadays, though, CV has been
740 adopted in a wide range of software engineering areas. Most frequently in
741 requirements engineering and software release planning. Eight studies use
742 CV as a research method while the remaining 32 studies use it as industry
743 practice.

744

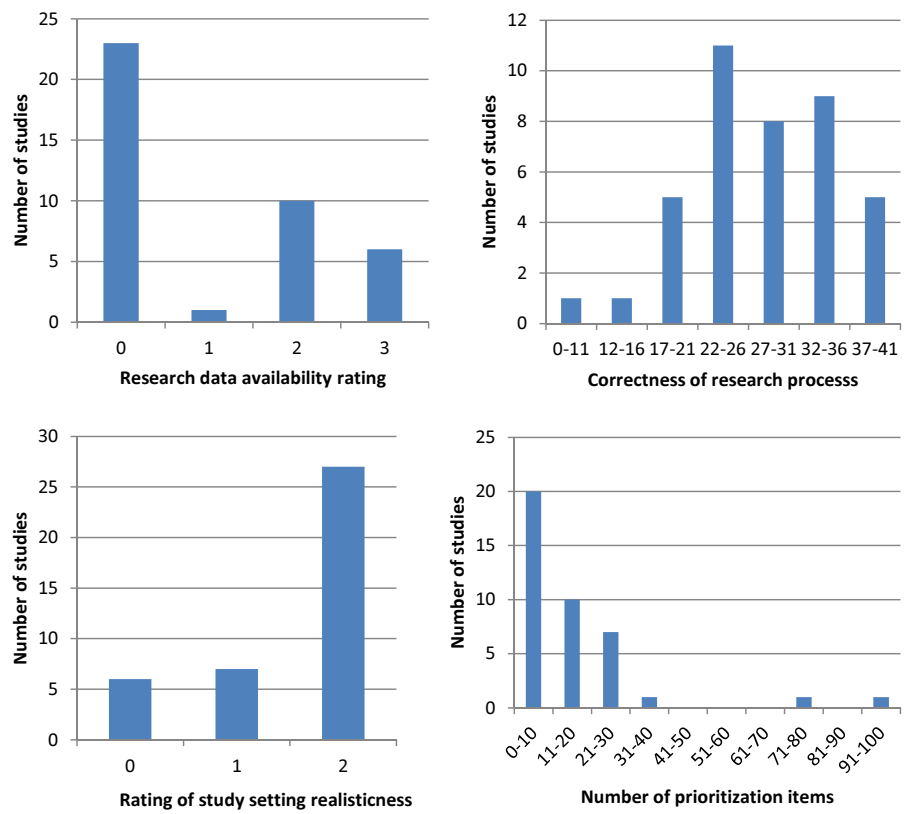
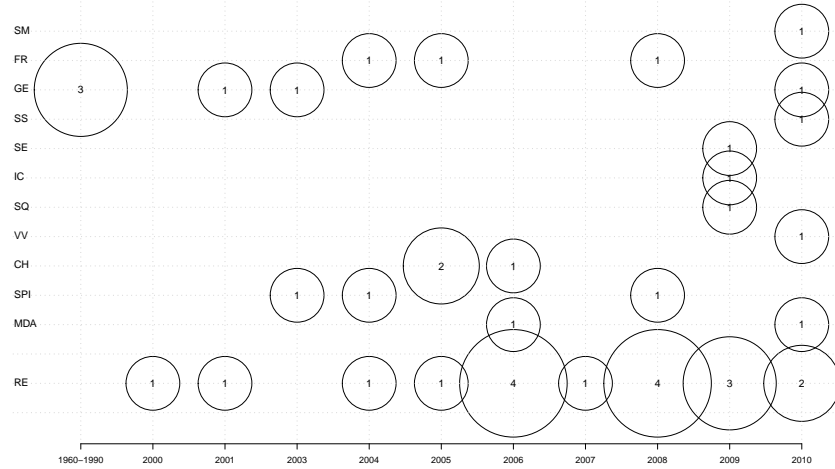


Figure 6: Study quality ratings



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
 VV - 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

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

The papers identified in the systematic review use various CV result analysis methods. The main goals for CV result analysis are presented in Table 13 and a summary of methods used in the primary studies can be found in Appendix B.

In order to reflect prioritization results many studies use charts or tables. These charts and tables show the average priority of each prioritization item that is computed from priorities assigned by all stakeholders. In [47] a table of five items with highest total priority is presented. [48] shows tables with minimum, maximum, median, mean, and standard deviation of priorities assigned by different stakeholders to a particular prioritization item. Finally, in [49, 48] error bars are added to the chart of final priorities (denoting the standard deviation of priorities).

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

Several papers are interested in the difference between priorities from different prioritization perspectives (e.g. current and ideal situation) or stakeholder groups (e.g. software developers and management). Pearson or Spearman correlation coefficients are commonly used to determine what the level of similarity between all priorities from two perspectives is. Whereas, Wilcoxon, Kruskal Wallis, Nemenyi-Damico-Wolfe-Dunn tests and the χ^2 statistic are used to detect if there is a significant difference in the value of one prioritization item from two or more perspectives. In addition, PCA is used to detect if there are distinct groups of stakeholders with common priorities [10, 7, 50].

In some cases, a stakeholder may assign equal priority to several prioritization items or leave several items unrated, e.g. the stakeholder may not have

carefully considered all prioritization items. Hence, the difference between the items may have been unnoticed.

In [4] the scalability of prioritization is measured using two charts. The first chart shows the average percentages of items given a non-zero value. The second chart shows average percentages of divergence of values. If a stakeholder assigns equal priorities to many prioritization items the divergence of values is low. Unfortunately it is unclear from [4] how the average percentage of divergence is calculated.

In [51] distribution, disagreement, and satisfaction charts are presented. The distribution chart shows how the final value of a prioritization item is constructed from priorities assigned by different stakeholders. This chart shows how much each stakeholder has contributed to the final value of a prioritization item. The disagreement chart shows the level of agreement between different stakeholders on the value of a particular prioritization item. The satisfaction chart shows stakeholder satisfaction with prioritization results by calculating the correlation between final priorities and priorities assigned by a stakeholder.

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

7.2.1. Problems with Compositional Data Analysis in Primary Studies

A few primary studies, as revealed by the systematic review, have problems with the analysis of compositional data.

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

The normality of compositional data is defined in [54]. It is stated that compositional data must first be transformed using isometric log-ratio transformation before the tests for normality can be applied. [47] violates this requirement by applying the Shapiro-Wilk test for normality to untransformed compositional data.

The Kruskal-Wallis test is used in [47] to analyse compositional data. The test is used to evaluate the difference between three organization levels. The Kruskal-Wallis test assumes that variables within each sample are in-

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]	Biplot
detect stakeholder groups with similar priorities [10]	Biplot
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

811 dependent [55]. However, values within compositional data vectors are not
812 independent (as described in Section 2.4). Hence, we claim the Kruskal-
813 Wallis test to be somewhat misused in [47].

814 7.3. Identifying Prioritization Items with Equal Priority Using ECV (RQ 3).

815 This section presents the results of applying ECV to the industrial and
816 academic CV result data as found through the systematic literature review.
817 Six primary studies included the raw prioritization results in the paper itself
818 or referenced online sources where the data was available. To collect the data
819 from the remaining 34 papers, the author(s) of all papers were contacted.

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

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

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)

830 In [46] a prioritization of requirement understandability criteria is pre-
831 sented. ECV shows that from the viewpoint of academia researchers, devel-
832 opment have the same importance as product planning (i.e. making strategic
833 product planning decisions: release planning, choosing which requirements
834 to dismiss).

835 A prioritization of software requirements for an academic course man-
836 agement system is presented in [17]. ECV detected that two features—
837 Assignment Submission and Assignment Feedback—have the same priority.

838 In [41] software product investments are prioritized with HCV. The re-
839 sults of ECV was different for uncompensated and compensated HCV results.
840 When compensated HCV was used ECV detected equal items that belong to
841 different high level prioritization groups (*A*, *B* and *C*). Whereas, in case of
842 uncompensated HCV all equal items belong to one high level prioritization
843 group (group *B*).

844 8. Discussion and Conclusions

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

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

853 Overall, studies that use CV or analyse the results of CV have high qual-
854 ity in terms of correctness of research process and study realism. However,
855 very few studies present prioritization of more than 30 items and the avail-
856 ability of research data is somewhat limited. In our particular case we were
857 able to obtain data from 43% of studies.

858 8.1. Implications for Practitioners

859 The results of this study provide decision support for industry practi-
860 tioners. We believe that a collection of state of the practice studies help
861 the adoption of CV prioritization method. (Top studies are summarized in
862 Table 12.) In addition, a set of CV analysis methods enables comprehensive
863 understanding of the prioritization results. (The analysis methods are pre-
864 sented in Table 13.) One of the most common goals of CV analysis are to
865 display of the prioritization results and, thus, to show the difference between
866 several prioritization perspectives.

867 Additionally, we present ECV—a novel method for CV analysis. Prior-
868 itization often results in the assignment of similar priorities to several prior-
869 itization items. ECV identifies prioritization items with similar priority and
870 tests whether these items can be considered equal. In this case, ECV can
871 be used in software release planning. For example, let us suppose that a set
872 of software requirements are prioritized with regard to the implementation
873 cost. First of all, ECV can then detect items with equal cost. Second, the
874 equal items can be freely swapped between the releases. Finally, the deci-
875 sion to allocate a requirement to a particular release can be made based on
876 another criteria, such as risk or business value.

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

880 *8.2. Implications for Academia*

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

887 CV results are a special type of data—compositional data. Standard
888 statistical analysis methods that assume the independence of the samples
889 cannot be applied to CV results. In [56] methods for the analysis of com-
890 positional data have been presented. The systematic review conducted as a
891 part of this study revealed that 22 studies analyse the results of CV. Yet,
892 only one study uses compositional data analysis methods, i.e. [10].

893 The small use of compositional data analysis is really not surprising,
894 because literature describing CV does not state that the results are com-
895 positional data. Standard statistical analysis methods may produce useful
896 results for compositional data. However, there are cases when they are mis-
897 leading or even faulty. Section 7.2.1 contains evidence of inappropriate use
898 of statistical methods by several papers.

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

902 Since CV is associated with compositional data, it might be tempting to
903 choose another requirements prioritization method. However, it would not
904 solve the problem, because any ratio scale prioritization, for instance AHP,
905 contains compositional data.

906 The principal implications for the academia are the following:

- 907 1. All systematic literature reviews should include snowball sampling.
- 908 2. Researchers can improve their statistical analysis of CV results using
909 compositional data analysis methods collected and developed by this
910 study.
- 911 3. When CV or any other ratio scale prioritization method is taught,
912 compositional data analysis should also be presented as part of the
913 solution.

914 8.3. *Validity Threats*

915 The validity of the systematic review is limited by the chosen databases,
916 the design of the review, and human judgement in study selection and data
917 extraction.

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

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

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

938 Four studies are not found because they are not available through databases
939 used in this systematic review. Out of them one is master thesis, two are
940 conference publications and one is a publication in the area of forestry. Seven
941 studies do not mention the name of the prioritization method (i.e. hundred
942 dollar method or cumulative voting). Only phrases like "distribution of a

943 predefined amount of fictitious money (\$100,000) over the items to be prior-
944 itized" or "1000 points" allowed us to identify that CV is used. One paper
945 used previously unknown name for CV - 100-point technique.

946 The quality of the data extraction and quality evaluation was validated
947 using inter-rater agreement analysis. In our case, 10% of the studies were
948 rated by two researchers and Krippendorff's alpha was calculated. The agree-
949 ment for the data extraction results was 0.86 and the agreement for the
950 quality evaluation was 0.73 (indicating a credible level of quality).

951 The failure to obtain raw results of several CV studies may be due to
952 several reasons, e.g. the authors of the primary studies might be unwilling
953 to communicate the data because of lack of motivation or spare time. In
954 our case we found that we were able to minimize this threat by searching for
955 the researchers through various channels, e.g. Google search, LinkedIn and
956 Facebook.

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

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

964 To mitigate the second threat we visually inspected the results of the ap-
965 plication of ECV on the real world data. We concluded that items identified
966 by ECV can be considered equal.

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

970 8.4. Future Research

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

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

982 *8.5. Conclusions*

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

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

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

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

999 **References**

- 1000 [1] P. Berander, A. Andrews, Requirements Prioritization, in: A. Aurum,
1001 C. Wohlin (Eds.), Engineering and Managing Software Requirements,
1002 Springer-Verlag, Berlin/Heidelberg, 2005, 2005, pp. 69–94.
- 1003 [2] D. Leffingwell, D. Widrig, Managing software requirements: A unified
1004 approach (1999) 118–119.
- 1005 [3] V. Ahl, An experimental comparison of five prioritization methods,
1006 Master’s Thesis, School of Engineering, Blekinge Institute of Technology
1007 (2005).
- 1008 [4] P. Berander, P. Jonsson, Hierarchical Cumulative Voting (HCV) prior-
1009 itization of requirements in hierarchies, 2006.
- 1010 [5] J. Karlsson, K. Ryan, A cost-value approach for prioritizing require-
1011 ments, IEEE Software 14 (1997) 67–74.
- 1012 [6] J. Karlsson, An evaluation of methods for prioritizing software require-
1013 ments, Information and Software Technology 39 (1998) 939–947.
- 1014 [7] F. Pettersson, M. Ivarsson, T. Gorschek, P. Öhman, A practitioner’s
1015 guide to light weight software process assessment and improvement plan-
1016 ning (2008).

- 1017 [8] S. Barney, C. Wohlin, Software Product Quality: Ensuring a Common
1018 Goal, in: Q. Wang, V. Garousi, R. Madachy, D. Pfahl (Eds.), Trust-
1019 worthy Software Development Processes, volume 5543 of *Lecture Notes*
1020 *in Computer Science*, Springer Berlin Heidelberg, Berlin, Heidelberg,
1021 2009, 2009, pp. 256–267.
- 1022 [9] P. Jönsson, C. Wohlin, A study on prioritisation of impact analysis
1023 issues: A comparison between perspectives, *Software Engineering Re-*
1024 *search and Practice in Sweden* (2005).
- 1025 [10] P. Chatzipetrou, L. Angelis, P. Rovegard, C. Wohlin, Prioritization of
1026 Issues and Requirements by Cumulative Voting: A Compositional Data
1027 Analysis Framework, 2010, pp. 361–370.
- 1028 [11] R. Engstrom, Cumulative Voting as a Remedy for Minority Vote Dilu-
1029 tion, *Local Government Election ...* (1999).
- 1030 [12] S. Bhagat, J. Brickley, Cumulative voting: The value of minority share-
1031 holder voting rights, *Journal of Law and Economics* (1984).
- 1032 [13] V. Hiltunen, J. Kangas, J. Pykalainen, Voting methods in strategic
1033 forest planning - Experiences from Metsähallitus, *Forest Policy and*
1034 *Economics* 10 (2008) 117–127.
- 1035 [14] P. Boldi, F. Bonchi, C. Castillo, S. Vigna, Voting in social networks,
1036 CIKM '09, ACM Press, New York, New York, USA, 2009.
- 1037 [15] H. Ayad, M. Kamel, Cumulative Voting Consensus Method for Parti-
1038 tions with Variable Number of Clusters, *Pattern Analysis and Machine*
1039 *Intelligence*, IEEE Transactions on 30 (2008) 160–173.
- 1040 [16] M. Svahnberg, A. Karasira, A Study on the Importance of Order in
1041 Requirements Prioritisation, IEEE, 2009.
- 1042 [17] P. Berander, M. Svahnberg, Evaluating two ways of calculating priorities
1043 in requirements hierarchies - An experiment on hierarchical cumulative
1044 voting, 2009.
- 1045 [18] T. Saaty, The analytic hierarchy process., McGraw-Hill, New York
1046 (1980).
- 1047 [19] S. Brenner, J. Schwalbach, Legal Institutions, Board Diligence, and
1048 Top Executive Pay, *Corporate Governance: An International Review*
1049 17 (2009) 1–12.

- 1050 [20] V. Pawlowsky-Glahn, J. J. Egozcue, Compositional data and their anal-
 1051 ysis: an introduction, Geological Society, London, Special Publications
 1052 264 (2006) 1–10.
- 1053 [21] J. Martin-Fernandez, C. Barceló-Vidal, V. Pawlowsky-Glahn, Dealing
 1054 with zeros and missing values in compositional data sets using nonpara-
 1055 metric imputation, *Mathematical Geology* 35 (2003) 253–278.
- 1056 [22] P. Filzmoser, K. Hron, Outlier detection for compositional data using
 1057 robust methods *Outlier Detection for Compositional Data Using Robust*
 1058 *Methods, Analysis and Applications* (2007).
- 1059 [23] K. Khan, A systematic review of software requirements prioritization,
 1060 Unpublished master’s thesis, Blekinge Institute of Technology, Ronneby,
 1061 Sweden (2006).
- 1062 [24] F. Zahedi, The analytic hierarchy process: a survey of the method and
 1063 its applications, *Interfaces* (1986) 96–108.
- 1064 [25] P. Runeson, M. Höst, Guidelines for conducting and reporting case
 1065 study research in software engineering, *Empirical Software Engineering*
 1066 14 (2008) 131–164.
- 1067 [26] L. Goodman, Snowball sampling, *The Annals of Mathematical Statis-*
 1068 *tics* (1961).
- 1069 [27] K. Krippendorff, Bivariate agreement coefficients for reliability of data,
 1070 *Sociological methodology* (1970).
- 1071 [28] K. Krippendorff, *Content analysis: An introduction to its methodology*
 1072 (2004).
- 1073 [29] B. Kitchenham, Guidelines for performing systematic literature reviews
 1074 in software engineering, *Engineering* (2007).
- 1075 [30] P. Berander, P. Jönsson, A goal question metric based approach for effi-
 1076 cient measurement framework definition, *ACM, Rio de Janeiro, Brazil,*
 1077 2006, pp. 316–325.
- 1078 [31] B. Regnell, M. Höst, J. och Dag, An industrial case study on distributed
 1079 prioritisation in market-driven requirements engineering for packaged
 1080 software, *Requirements ...* (2001).
- 1081 [32] B. Kitchenham, *Procedures for performing systematic reviews*, Keele,
 1082 UK, Keele University 33 (2004).

- 1083 [33] M. Ivarsson, T. Gorschek, A method for evaluating rigor and industrial
1084 relevance of technology evaluations, *Empirical Software Engineering*
1085 (2010) 1–31.
- 1086 [34] C. Wohlin, P. Runeson, M. Höst, *Experimentation in software engineer-*
1087 *ing: an introduction*, Springer Netherlands, 2000.
- 1088 [35] A. Jedlitschka, D. Pfahl, Reporting guidelines for controlled experi-
1089 ments in software engineering, in: *2005 International Symposium on*
1090 *Empirical Software Engineering*, 2005., IEEE, 2005, p. 10.
- 1091 [36] K. Rinkevics, *Data Extraction and Quality Evaluation results*, 2011.
- 1092 [37] R. Ihaka, R. Gentleman, R: a language for data analysis and graphics,
1093 *Journal of computational and graphical statistics* (1996) 299–314.
- 1094 [38] K. Rinkevics, *ECV implementation source code*, 2011.
- 1095 [39] R. M. Groves, F. J. Fowler, M. P. Couper, J. M. Lepkowski, E. Singer,
1096 *Survey methodology*, John Wiley and Sons, 2009.
- 1097 [40] S. Barney, A. Aurum, C. Wohlin, The Relative Importance of Aspects
1098 of Intellectual Capital for Software Companies, in: *2009 35th Euromicro*
1099 *Conference on Software Engineering and Advanced Applications*, IEEE,
1100 2009, 2009, pp. 313–320.
- 1101 [41] S. Barney, C. Wohlin, A. Aurum, *Balancing software product invest-*
1102 *ments*, IEEE Computer Society, 2009, pp. 257–268.
- 1103 [42] S. Barney, A. Aurum, C. Wohlin, A product management challenge:
1104 Creating software product value through requirements selection, *Jour-*
1105 *nal of Systems Architecture* 54 (2008) 576–593.
- 1106 [43] S. Laukkanen, T. Palander, J. Kangas, A. Kangas, Evaluation of the
1107 multicriteria approval method for timber-harvesting group decision sup-
1108 port, *Silva Fennica* 39 (2005) 249–264.
- 1109 [44] G. Hu, Adding value to software requirements: An empirical study in
1110 the chinese software industry, *Seventeenth Australian Conference on*
1111 *...* (2006).
- 1112 [45] R. Feldt, R. Torkar, E. Ahmad, B. Raza, Challenges with Software
1113 Verification and Validation Activities in the Space Industry, IEEE, 2010.

- 1114 [46] M. Svahnberg, T. Gorschek, M. Eriksson, A. Borg, K. Sandahl,
1115 J. Börster, A. Loconsole, Perspectives on Requirements Understand-
1116 ability – For Whom Does the Teacher’s Bell Toll?, IEEE, 2008.
- 1117 [47] P. Jönsson, C. Wohlin, Understanding impact analysis: An em-
1118 pirical study to capture knowledge on different organisational levels,
1119 ... Conference on Software Engineering and Knowledge ... (2005).
- 1120 [48] L. a. Kuzniarz, Empirical extension of a classification framework for
1121 addressing consistency in model based development, Information and
1122 Software Technology (2010).
- 1123 [49] P. Rovegard, L. Angelis, C. Wohlin, An Empirical Study on Views of
1124 Importance of Change Impact Analysis Issues, Software Engineering,
1125 IEEE Transactions on 34 (2008) 516 –530.
- 1126 [50] C. Wohlin, A. Aurum, Criteria for selecting software requirements to
1127 create product value: An industrial empirical study, Value-Based Soft-
1128 ware Engineering (2006).
- 1129 [51] B. Regnell, M. Höst, J. Natt, Visualization of Agreement and Satisfac-
1130 tion in Distributed Prioritization of Market Requirements, Chart (2000)
1131 1–12.
- 1132 [52] J. Aitchison, Principal component analysis of compositional data,
1133 Biometrika 70 (1983) 57.
- 1134 [53] P. Filzmoser, K. Hron, C. Reimann, F. Sm, P. Filzmoser, K. Hron,
1135 C. Reimann, Principal component analysis for compositional data with
1136 outliers Principal component analysis for compositional data with out-
1137 liers, Analysis and Applications (2007).
- 1138 [54] V. Pawlowsky Glahn, J. Egozcue, R. Tolosana Delgado, Lecture notes
1139 on compositional data analysis, Interpretation A Journal Of Bible And
1140 Theology (2007).
- 1141 [55] W. Kruskal, W. Wallis, Use of ranks in one-criterion variance analysis,
1142 Journal of the American statistical Association 47 (1952) 583–621.
- 1143 [56] J. Aitchison, The statistical analysis of compositional data, Chapman
1144 & Hall, London, 1986.
- 1145 [57] D. Baca, K. Petersen, Prioritizing Countermeasures through the Coun-
1146 termeasure Method for Software Security (CM-Sec), in: M. Ali Babar,

- 1147 M. Vierimaa, M. Oivo (Eds.), Product-Focused Software Process Im-
 1148 provement, volume 6156 of *Lecture Notes in Computer Science*, Springer
 1149 Berlin / Heidelberg, 2010, 2010, pp. 176–190.
- 1150 [58] S. a. b. Bowler, Election systems and voter turnout: Experiments in
 1151 the United States, *Journal of Politics* 63 (2001) 902–915.
- 1152 [59] D. Brockington, A Low Information Theory of Ballot Position Effect,
 1153 *Political Behavior* 25 (2003) 1–27.
- 1154 [60] D. Cooper, A. Zillante, A comparison of cumulative voting and gener-
 1155 alized plurality voting, *Public Choice* (2010).
- 1156 [61] N. D. Fogelström, M. Svahnberg, T. Gorschek, Investigating Impact of
 1157 Business Risk on Requirements Selection Decisions, IEEE, 2009.
- 1158 [62] S. Hatton, Choosing the Right Prioritisation Method, in: Proceed-
 1159 ings of the 19th Australian Conference on Software Engineering, IEEE
 1160 Computer Society, Washington, 2008, pp. 517–526.
- 1161 [63] S. Hatton, Early prioritisation of goals, in: Proceedings of the 2007
 1162 conference on Advances in conceptual modeling: foundations and appli-
 1163 cations, ER’07, Springer-Verlag, Berlin, 2007, pp. 235–244.
- 1164 [64] V. Heikkilä, A. Jadallah, K. Rautiainen, G. Ruhe, Rigorous Support
 1165 for Flexible Planning of Product Releases - A Stakeholder-Centric Ap-
 1166 proach and Its Initial Evaluation, IEEE, 2010.
- 1167 [65] M. Staron, C. Wohlin, An Industrial Case Study on the Choice Be-
 1168 tween Language Customization Mechanisms, in: J. Münch, M. Vier-
 1169 imaa (Eds.), Product-Focused Software Process Improvement, volume
 1170 4034 of *Lecture Notes in Computer Science*, Springer Berlin / Heidel-
 1171 berg, 2006, 2006, pp. 177–191.
- 1172 [66] T. Touseef, C. Gancel, A structured goal based measurement framework
 1173 enabling traceability and prioritization, ... (ICET), 2010 6th Interna-
 1174 tional Conference on (2010).
- 1175 [67] P. Berander, C. Wohlin, Differences in views between development
 1176 roles in software process improvement-a quantitative comparison, in:
 1177 Proceedings 8th Conference on Empirical Assessment in Software Engi-
 1178 neering, 2004.

- 1179 [68] P. Berander, Using students as subjects in requirements prioritization,
 1180 Proceedings. 2004 International Symposium on Empirical Software En-
 1181 gineering, 2004. ISESE '04. (2004) 167–176.
- 1182 [69] P. Berander, C. Wohlin, Identification of Key Factors in Software Pro-
 1183 cess Management-A Case Study (2003).
- 1184 [70] R. L. Cole, D. a. Taebel, R. L. Engstrom, Cumulative Voting in a Munic-
 1185 ipal Election: A Note on Voter Reactions and Electoral Consequences,
 1186 The Western Political Quarterly 43 (1990) 191.
- 1187 [71] J. Kuklinski, Cumulative and Plurality Voting: An Analysis of Illinois’
 1188 Unique Electoral System, The Western Political Quarterly 26 (1973)
 1189 726–746.
- 1190 [72] S. Laukkanen, T. Palander, J. Kangas, Applying voting theory in par-
 1191 ticipatory decision support for sustainable timber harvesting, Canadian
 1192 Journal of Forest Research 34 (2004) 1511–1524.
- 1193 [73] J. Sawyer, D. MacRae, Game theory and cumulative voting in Illinois:
 1194 1902-1954, The American Political Science Review 56 (1962) 936–946.

1195 Appendix A. Primary Studies

1196 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 Jönsson [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]
	Investigating Impact of Business Risk on Requirements Selection Decisions	Fogelström et al. [61]
1197	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 multicriteria 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]

1198 *Appendix A.2. Primary studies revealed by snowball sampling*

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 "1000 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 "1000 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 organisational 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 1000 points to spend among the 13 criteria."

Appendix B. CV Result Analysis Methods

	Paper																					
	Svahnberg2008	Svahnberg2009	Starou2006	Petersson2008	Wohlin2006	Laakkonen2005a	Hu2006	Jonsson2005a	Kuzniarz2010	Rovgeard2008	Berander2006a	Berander2004a	Berander2006	Feldt2010	Barney2009b	Barney2008	Barney2009a	Barney2009	Jonsson2005	Chatzipetrou2010	Reguel2001	Reguel2000
analysis method																						
table that shows final priorities	x			x												x						
chart that shows final priorities	x			x	x	x	x									x						
table of top 5 prioritization items								x														
minimal, maximal, mean, median, and standard deviation of priorities assigned by different stakeholders									x	x												
bar chart of prioritization results showing mean priority and standard deviation of priorities									x	x												
Pearson correlation coefficient		x										x										
Nemenyi Damico Wolfe Dunn test														x								
Spearman's r															x			x				
Kruskal Wallis test								x														
Wilcoxon test							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 Chi-square 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
satisfaction chart				x																		x
biplot																					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 artefacts	
10. Measures, Data collection procedures	Who, when, and how does the measurements [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, generalisation	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		