
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 should 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

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

2. Background

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

2.1. Prioritization Methods

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

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

Ratio scale priorities have several advantages over ordinal scale priorities. Ratio scale shows not just the order of items but also relative distance between them. This enables the priority of a group of items to be calculated by summing up the priorities of individual items [4]. It is possible to say that one item or set of items has higher priority than another set of items. Supposing stakeholders have to choose between several low priority items and one item with higher priority; with ordinal scale, the item with highest priority will always be selected first. However, if priorities are given on a ratio scale, it is possible that lower priority items will be selected if their cumulative priority is higher. Knowing the relative importance of sets of prioritization items helps in software release planning. Ratio scale allows the combining of multiple priority factors by calculating ratios between them. One example of this is the cost-value ratio that shows which requirements 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 analysed 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 analysed using the Kruskal-
108 Wallis test. In [10] the results of [9] are further analysed 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 analysed using various techniques. We will next describe in more
116 detail one of the more common methods to manage prioritization issues —
117 cumulative voting — which has been used in software engineering for some
118 time, but has its roots in corporate governance and biology.

119 *2.3. Cumulative Voting*

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

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

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

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

147 *2.3.1. Benefits and Drawbacks of Cumulative Voting*

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

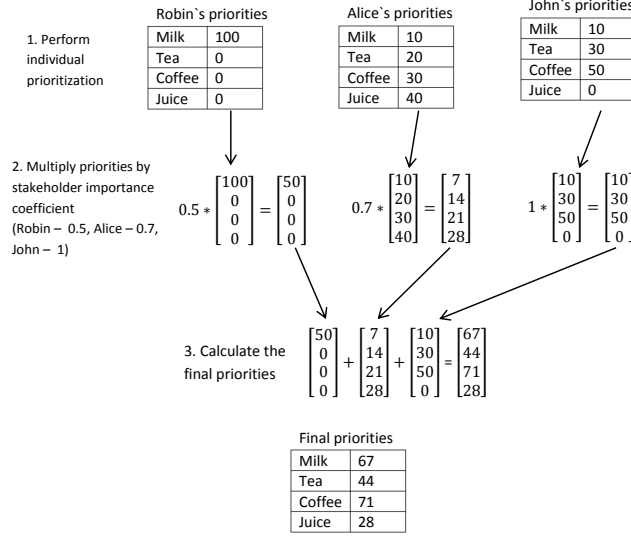


Figure 1: Example of CV with several stakeholders.

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

156 2.3.2. Example of Cumulative Voting with Several Stakeholders

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

167 2.3.3. Stakeholder Bias

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

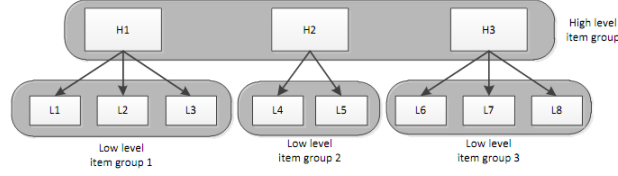


Figure 2: Example of prioritization item hierarchy.

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

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

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

2.3.4. Scalability of Cumulative Voting, Hierarchical Cumulative Voting

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

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

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

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

202 2.3.5. *Compensation Factors*

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

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

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

220 2.3.6. *HCV Execution*

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

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

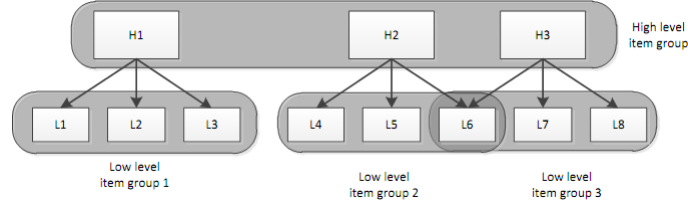


Figure 3: Overlapping prioritization item hierarchy example.

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

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

It is possible that one low level item is child of more than one high level requirement and, thus, belongs to two or more low level requirement groups (see Figure 3). Such requirements participate in the standard HCV prioritization process and are prioritized two or more times with each group they belong to. At the end of the prioritization they receive several priority values. These values can 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.

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

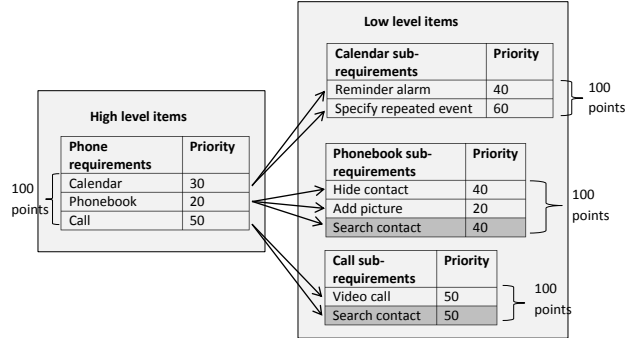


Figure 4: Example of hierarchical cumulative voting, requirement hierarchy

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

requirement is computed by multiplying the priority of the requirement with the priority of its parent high level requirement and the compensation factor. The compensation factor in this particular case is the number of elements in a group, two for the ‘calendar’ and ‘call’ sub-requirements and three for the ‘phonebook’ sub-requirement.

2.4. Compositional Data Analysis

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

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

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

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

287 Another property of compositional data items is *subcompositional coher-*
288 *ence*. Consider that two compositions are analysed. One composition is a
289 subcomposition of the other. *Subcompositional coherence* means that the re-
290 sults of the analysis are the same for the common parts of the compositions
291 [20]. This property is important for the analysis of HCV results. Statements
292 that are made regarding each smaller group of prioritization items are also
293 true for all items prioritized with HCV.

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

301 When doing analysis of compositional data one must take into account
302 that compositional data special type of data and should be analysed differ-
303 ently than ordinary data. Ordinary unconstrained variables are free to take
304 any positive or negative values, whereas, compositional data values can only
305 be positive and have a constrained maximum value. Moreover, components
306 of compositional data vectors are not independent from each other. The fact
307 that an item is assigned 70 priority points means that the next item can take
308 only values between 0 and 30. Hence, there is a negative correlation between
309 the items.

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

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

314 where Y is the product of lineal combination and a_i is any real number.
315 Now, since the sum of priorities assigned in CV must add up to 100 (or any

316 other constant number) at least one linear combination of X is not normally
 317 distributed because it always adds up to 100:

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

318 In our opinion, the above indicates, quite strongly, that CV results do
 319 not follow a multivariate normal distribution and, hence, it follows that they
 320 should be analysed using non-parametric statistical tests [21].

321 2.4.1. Problem of Zeroes

322 Compositional data analysis requires that log-ratios between any com-
 323 ponents in a vector can be computed. But computing a log-ratio with a
 324 zero value is, in this case, meaningless. This is a problem since CV allows
 325 stakeholders to assign zero priorities to some prioritization items (we would
 326 even strongly argue that this is very common).

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

333 Before compositional data analysis can be applied to CV results, we
 334 should first remove zeroes in the data. One approach can be to forbid stake-
 335 holders to assign zero priorities. This approach is used in e.g. [7]. But this
 336 can add some unnecessary complexity to the prioritization process and, ex-
 337 plicitly, delimits an expert's freedom. In [10] the authors propose the use
 338 of a multiplicative replacement strategy (as defined in [22]) for CV result
 339 analysis.

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

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

340 where δ_j is the imputed value and c is the constant sum constraint.
 341 In order for the total sum of components to stay constant, the equation
 342 subtracts some value from the items with a priority higher than zero. More
 343 is subtracted from components with higher values than from components
 344 with lower values (and the value of the imputed δ_j is arbitrary).

345 2.4.2. Isometric log-ratio transformation

346 In order to apply standard statistical methods to compositional data
 347 it should be transformed to remove the inherent correlation of the values.
 348 Compositional data analysis proposes special transformations that change
 349 the compositional data values to unconstrained real values. One such trans-
 350 formation is isometric log-ratio (*ilr*) transformation (as proposed by [21, 23]):

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

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

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

355 3. Related Work

356 A systematic review of requirements prioritization methods is presented
 357 in [24]. The study focuses on prioritization method comparison and selects
 358 eight relevant studies. Two of the studies use CV. These studies are also
 359 revealed by the systematic literature review conducted as part of this study.
 360 Khan [24] concludes that there is little research on requirements prioritization
 361 and studies usually deal with a small number of requirements.

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

365 4. Methodology

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

368 4.1. Selection of Research Methods

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

371 One way of collecting this knowledge is to conduct an empirical study. A
 372 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 [25]), while a case study can be used to receive qualitative feedback on the use of CV [26].

Knowledge on the empirical use of CV can also be obtained from existing studies. This may be done by means of a systematic literature review. Several studies have used CV in industry as well as in academic settings. Nevertheless, there are no studies that provide an overview of the current state of the practice in this field (as reported by research studies). Therefore, before continuing with the refinement of CV and conducting new empirical studies (i.e. case study or experiment), a systematic literature review would be required.

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

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

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

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

4.2. Research Questions

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

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

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

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

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

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

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

420 5. Systematic Literature Review

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

423 Table 2 presents an overview of activities performed during the system-
424 atic literature review. The review protocol was developed by one researcher
425 and evaluated by another researcher. Studies were searched for in two itera-
426 tions. The first search was performed by using databases. The second search
427 was performed using snowball sampling [27] (snowball sampling examines the
428 references of primary studies revealed by the first search). References that
429 are relevant to the review, i.e. they pass the selection criteria, are then added
430 to the set of primary studies.

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

436 To ensure the quality of the review, the quality evaluation and data ex-
437 traction was performed independently by two researchers. Inter-rater anal-
438 ysis was performed using Krippendorff’s Alpha statistics [28, 29].

439 5.1. Data Sources and Search Strategy

440 This SLR was designed based on the guidelines by Kitchenham [30]. First
441 a trial search in electronic databases was conducted. In order to scale the
442 review to a manageable, yet sufficient size, databases were searched with
443 different search strings. Relevant papers that were found during the trial

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	

search were used to extract additional search strings. The trial search revealed that the number of studies that use CV is not very large. Therefore, we decided to include not only software engineering studies but also studies in other research areas, such as forestry or corporate governance, since one key aspect we intended to investigate was analysis methods for CV.

Since CV is frequently used in studies without mentioning this in the abstract, full text search in databases is preferable. Unfortunately not all databases support full text search. Full text search was performed in the IEEE Xplore and Springer Link databases. In ACM Digital Library, Inspec/Compendex, ISI Web of Knowledge, and SCOPUS only metadata was searched. Search strings consisting of a Boolean expression (A or B or C or D or E or F or G), where:

Search strings contained only synonyms of CV and they did not limit the research area to software engineering. The search was performed independently using each of the search strings in each database. All search results were combined and documented using reference management software. The quality of the search strings and the selection of electronic databases were validated against a previously known core set of papers—[3, 31, 10, 32]—checking that all papers from the core set were found by the search.

5.2. Study Selection

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

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

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

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	

478 5.3. *Quality Evaluation*

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

482 Study quality obviously depends on the correctness of the study process
483 including planning, operation, analysis and interpretation of the results (is
484 the study right?) The correctness of the process can be measured by evalu-
485 ating the description of the study or replicating the study. Thus, to gain the
486 trust of industry practitioners and other researchers, the process of the study
487 should be rigorously described. In short, the description has to facilitate the
488 replication of the study as well as the presentation of limitations and validity
489 threats.

490 Even the most correct and rigorously described study is useless if it does
491 not contribute to the industry or research community (is it the right study?)
492 The topic of the research ought to address important goals and issues. The
493 findings of the study should also be significant, i.e. there is a high probability
494 of the results of the study are true. The significance of the findings depends
495 on how realistic the study is, the correctness of the process and the results
496 of the study, as well as the statistical significance of the findings.

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

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

516 The **type** of the research study can be considered as a criterion for the
517 evaluation of study realism. Reference [33] 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 [34] 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 settings, whereas, experiments are usually performed in academia using students as subjects. Therefore, [34] argue that case studies are more realistic than formal experiments. Obviously the effect of study design on the study realism may be interpreted in different ways. Therefore, we will not use this parameter in our quality evaluation.

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

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

5.3.1. Rating of the Studies

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

Realism of the studies is rated in three aspects: subjects, setting, and scale. The subjects and setting is rated according to Table 5. The total rating of study realism is determined by summing up the ratings of the two aspects. For instance, if a study is conducted with industry professionals as subjects in an academic context the study will receive rating 1 (out of 2 maximal points).

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

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

563 The availability of raw results of CV is rated separately because it is
564 especially important for our purposes (and for most other researchers in
565 order to replicate a study). The data availability rating criteria is given in
566 Table 6. If the data of a study are not available it is not possible to validate
567 the results of the study and, hence, the credibility of the findings is lower.
568 Ideally the data collected in the study should be presented directly in the
569 paper. An alternative may be to make the data freely available online and
570 reference the online source.

571 The quality of the research methodology of a paper is rated according to
572 a checklist presented in Appendix C. The checklist is based on guidelines
573 for presenting research studies as presented in [35, 36] and the guidelines for
574 quality evaluation of research studies presented in [34, 30]. Evaluation is done
575 with regard to the rigor of the description and correctness of the research
576 process and reasoning. Checklist items represent issues that research studies
577 should implement and present in research paper. The checklist also contains
578 item descriptions or questions that are used to evaluate the quality. Each
579 item in the checklist is rated according to criteria presented in Table 7. The
580 final rating of correctness of the research process of a study is computed by
581 summing up the ratings assigned to all items in the checklist.

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

586 As a result of the rating each study was assigned four rating values on an
587 ordinal scale. In order to perform a more advanced analysis of the quality
588 evaluation results these ratings were then converted into ratio scale ranks.
589 For each study, the number of studies that have received lower ratings is

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.

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.

counted. The resulting number is the rank of the study; thereby, the quality of a study is expressed as four rank values.

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

5.4. Data Extraction

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

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

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

Data extraction was documented with the help of spreadsheet software. Extracted data items are available from [37].

6. Equality of Cumulative Votes

In the previous section we described the execution of the systematic literature review. In order to perform a more thorough analysis later we here present the design of ECV before presenting the results of the systematic literature review. For the results of the evaluation of ECV please see Section 7.3 (ECV is implemented in the *R* programming language [38] and the code can be found at [39].)

In CV stakeholders may assign similar or equal values to several prioritization items. As a result the difference between the items is small. The variation in priorities is caused not only by the difference between prioritization items but also by human error and lack of information for decision making. For instance, people tend to simplify the task of prioritization by assigning rounded values to items or giving equal values to several items [40].

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

6.1. Testing Equality of Two Items

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

The first approach uses statistical tests to evaluate differences between e.g. two population means, in order to determine that two items are different. Populations in this case consist of priorities assigned by all stakeholders to a particular prioritization item. The number of stakeholders that perform the

631 prioritization is frequently small. Hence, the size of the sample is very often
632 too small for statistical tests to detect a significant difference in the tests,
633 thus, identify too many equal items to make any useful conclusions.

634 ECV, in contrast, uses the second approach. It finds items that are
635 similar and the rest of the items are considered different. This method tests
636 the probability of the difference between the means of two items being smaller
637 than the given value. In short, ECV tests the probability of the means of two
638 prioritization items differing by less than 25%. If the probability is higher
639 than 70% the items are considered equal.

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

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

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

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

654 where a is an element of matrix A , and k and l are the columns that repre-
655 sent items that are tested for equality, "*" denotes all rows of corresponding
656 column.

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

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

660 where c_i is the i^{th} element of C and b_{i1} and b_{i2} are the first and second
661 elements in the i^{th} row of B . Each value c_i represents a log-ratio between
662 values of columns k and l . The mean of the values of C can be interpreted as

an average log-ratio between the items that expresses the difference between the items.

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

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

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

Otherwise a non-parametric estimation of the distribution function could be performed. In our case, the CV result data obtained from the primary studies identified by the systematic review, were tested for normality using the Anderson-Darling test. Before applying the test the data was transformed using methods of compositional data analysis. To compute the test we used method *adtestWrapper* from *R* language library *robCompositions*.

The tests we performed indicated, quite strongly, that in most of the prioritization cases the data is not normally distributed. Hence, our recommendation is that, in general, a non-parametric approach should be used to determine the probability density function, and one such, common, approach would be to use the kernel density estimation. (In our implementation of ECV in the *R* programming language, kernel density estimation is performed using the package *ks*.)

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

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

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

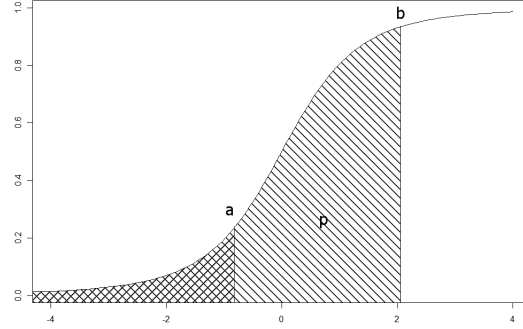


Figure 5: Cumulative distribution function of the log-ratio c_i between the items k and l (area p denotes probability that c_i is between $\frac{3}{4}$ and $\frac{4}{3}$.)

Table 10: Example of equality table

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

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

6.2. Grouping Prioritization Items

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

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

7. Results

This section presents the results of this study including the systematic literature review and the application of ECV on industry and academic data collected from the primary studies. Data extracted from primary studies and the results of the quality evaluation are available in [37].

718 *7.1. State of Practice in Empirical Studies that use CV or Analyze the Re-*
719 *sults of CV (RQ 1)*

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

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

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

738 To perform snowball sampling we examined the references of primary
739 studies that were found during the database search. References were used
740 to search for the papers in the Google and Google Scholar search engines.
741 Studies that were found in the search and passed the study selection criteria
742 were added to the set of primary studies.

743 After the primary studies were selected, data extraction and quality eval-
744 uation was performed by two researchers. One researcher examined all stud-
745 ies while the second researcher did quality evaluation and data extraction for
746 10% of the studies. The studies were randomly selected. Inter-rater agree-
747 ment were calculated by means of Krippendorff’s alpha coefficient. Agree-
748 ment for data extraction results was 0.86 and agreement for the quality
749 evaluation was 0.73. According to [29] it is common to require agreement
750 above 0.8 and the lowest acceptable agreement is 0.667. Therefore, we con-
751 clude that the agreement calculated for this study is sufficient. Ratings of
752 the study setting, correctness, research data availability, and number of pri-
753 oritization items are presented in Figure 6.

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

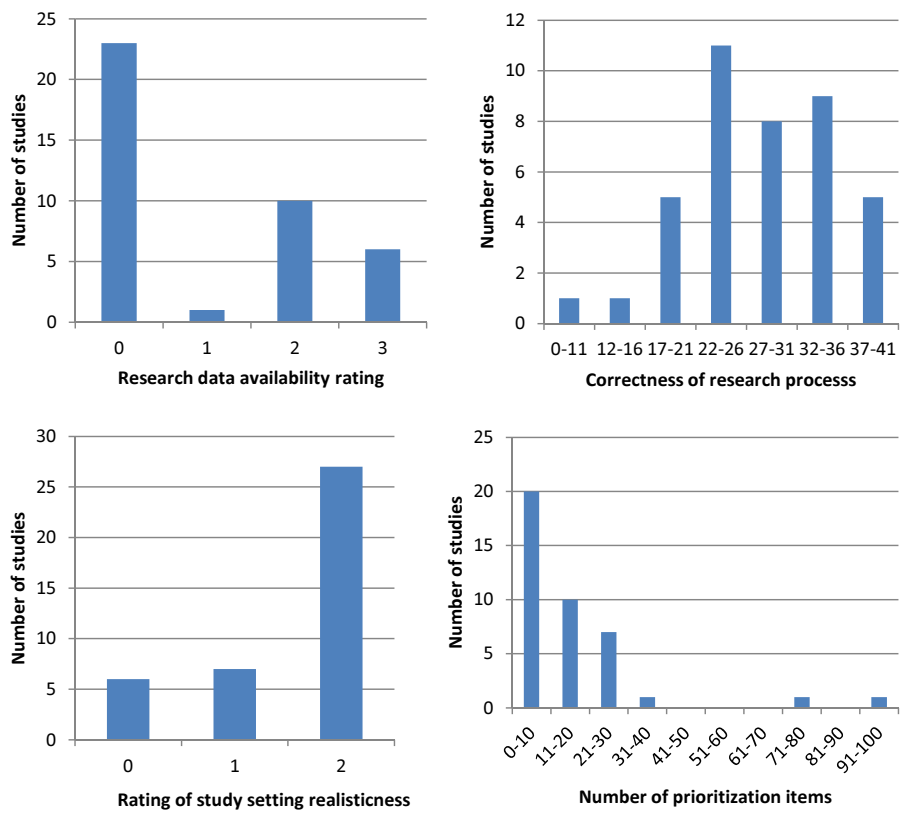
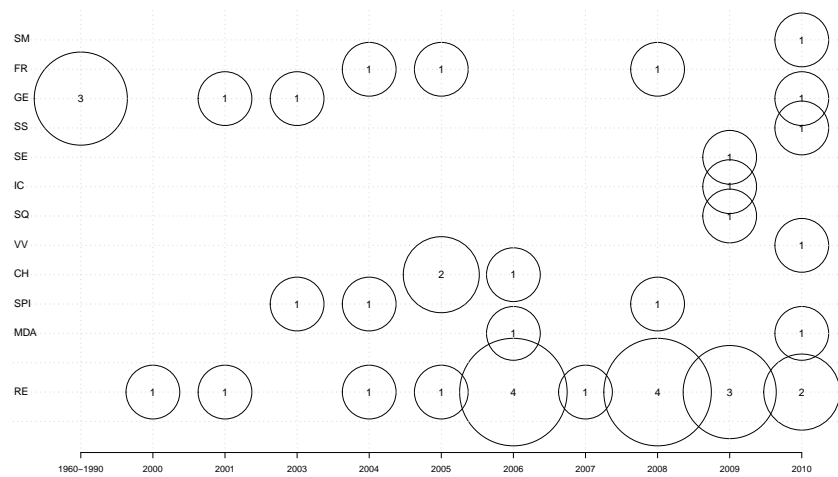


Figure 6: Study quality ratings



MDA - model driven software development	FR - forestry
CH - change impact analysis in software engineering	GE - government elections
RE - requirements engineering and software release planning	SS - software security
IC - intellectual capital in software company	SQ - software quality
SPI - software process improvement	SM - software metrics
V&V - software verification and validation	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 2009 [41]	36	2	2	17
Berander 2009 [17]	41	2	0	29
Barney 2009 [42]	40	2	2	5
Barney 2009 [8]	31	2	2	27
Barney 2008 [43]	34	2	2	14
Laukkanen 2005 [44]	22	3	2	30
Hu 2006 [45]	34	2	1	14
Feldt 2010 [46]	24	3	2	8
Regnell 2001 [32]	21	3	2	91
Svahnberg 2008 [47]	34	1	1	7

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

766 *7.2. CV Result Analysis Methods Identified by RQ 1 (RQ 2)*

767 The papers identified in the review use various CV result analysis meth-
768 ods. The main goals for CV result analysis are presented in Table 13 and
769 a summary of methods used in the primary studies can be found in Section
770 Appendix B.

771 In order to present prioritization results many studies use charts or tables.
772 These charts and tables show the average priority of each prioritization item
773 that is computed from priorities assigned by all stakeholders. In [48] a table
774 of five items with highest total priority is presented. [49] shows tables with
775 min , max , \tilde{x} , \bar{x} and σ of priorities assigned by different stakeholders to a
776 particular prioritization item. Finally, in [50, 49] error bars are added to the
777 chart of final priorities (denoting σ of priorities).

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

781 Several papers are interested in the difference between priorities from dif-
782 ferent prioritization perspectives (e.g. current and ideal situation) or stake-
783 holder groups (e.g. software developers and management). Pearson or Spear-
784 man correlation coefficients are commonly used to determine what the level of
785 similarity is between all priorities from two perspectives. Whereas, Wilcoxon,
786 Kruskal-Wallis, Nemenyi-Damico-Wolfe-Dunn tests and the χ^2 statistic are
787 used to detect if there is a significant difference in the value of one prioritiza-
788 tion item from two or more perspectives. In addition, PCA is used to detect
789 if there are distinct groups of stakeholders with common priorities [7, 10, 51].

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

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

800 In [52] distribution, disagreement, and satisfaction charts are presented.
801 The distribution chart shows how the final value of a prioritization item is
802 constructed from priorities assigned by different stakeholders. This chart
803 shows how much each stakeholder has contributed to the final value of a
804 prioritization item. The disagreement chart shows the level of agreement be-
805 tween different stakeholders on the value of a particular prioritization item.

Table 13: Goals for CV result analysis.

Purpose of the method	Name
Show the final priority of each prioritization item. Stakeholder priorities are combined into one value.	Chart or table of final priorities
Difference between priorities assigned by different perspectives (status quo, ideal situation) or different stakeholder groups (developers, management) [10]	Bi-plot
detect stakeholder groups with similar priorities [10]	Bi-plot
show the relative number of issues that have received high, medium, or low priority [10]	Ternary plot
detect stakeholder groups with common priorities [10]	PCA
how the final value of prioritization item is constructed from priorities assigned by different stakeholder. This chart shows how much each stakeholder has contributed to the final value of prioritization item [52]	Distribution chart
the level of agreement between different stakeholders on value of particular prioritization item [52]	Disagreement chart
satisfaction of a stakeholder with the prioritization results by the calculating correlation between the final priorities and priorities assigned by a stakeholder [52]	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

806 The satisfaction chart shows stakeholder satisfaction with prioritization re-
807 sults by calculating the correlation between final priorities and priorities
808 assigned by a stakeholder.

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

816 7.2.1. Problems with Data Analysis in Primary Studies

817 A few primary studies, as revealed by the systematic review, have prob-
818 lems with the data analysis. These studies disregard the compositional na-
819 ture of CV results.

820 In [51, 7] standard PCA is performed without applying log-ratio trans-
821 formations to compositional data. According to [53], this is likely to be
822 inadequate and in [54], a more appropriate method for performing PCA of
823 compositional data is shown.

824 The normality of compositional data is defined in [55]. It is stated that

Table 14: Identified groups of equal items.

Paper identifier & Description	Type of CV	Pairs of equal items	Groups of equal items
Barney 2009 [42] 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 2009 [17] Software requirements for course management system	uncomp. & comp. HCV	(3:2, 3:3)	(3:2, 3:3)
Svahnberg 2008 [47] The view of academia researchers on the requirements understandability criteria	CV	(Development, Verification & Validation) (Development, Product Planning 1)	(Development, Product Planning 1)

825 it is convenient to transform compositional data using isometric log-ratio
826 transformation before the tests for normality can be applied. [48] violates
827 this requirement by applying the Shapiro-Wilk test for normality to untrans-
828 formed compositional data.

829 The Kruskal-Wallis test is used in [48] to analyze compositional data.
830 The test is used to evaluate the difference between three organization levels.
831 The Kruskal-Wallis test assumes that variables within each sample are in-
832 dependent [56]. However, values within compositional data vectors are not
833 independent (as described in Section 2.4). Hence, we claim the Kruskal-
834 Wallis test to be somewhat misused in [48].

835 7.3. Identifying Prioritization Items with Equal Priority Using ECV (RQ 3)

836 This section presents the results of applying ECV to the industrial and
837 academic CV data as found through the systematic literature review. Six
838 primary studies included the raw prioritization results in the paper itself or

839 referenced online sources where the data was available. To collect the data
840 from the remaining 34 papers, the authors of all papers were contacted.

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

846 In short, ECV was applied to 27 CV prioritization cases from 14 stud-
847 ies. In the cases of HCV, ECV was applied two times to the same data
848 to test both compensated and uncompensated priorities. Equal items were
849 detected in three prioritization cases. A summary of the results is presented
850 in Table 14.

851 In [47] a prioritization of requirement understandability criteria is pre-
852 sented. One of the main findings of paper [47] is that from an academic
853 viewpoint Development and Verification and Validation are more important
854 than other criteria. ECV adds new knowledge to these results. It shows that
855 Development and Verification and Validation are equally important, i.e. it is
856 not true that either one of the criteria is more important.

857 A prioritization of software requirements for an academic course man-
858 agement system is presented in [17]. ECV detected that two features—
859 Assignment Submission and Assignment Feedback—have the same priority.
860 If the system is developed in several releases Assignment Submission and As-
861 signment Feedback features can be freely interchanged between the releases
862 and, hence, in this way ECV simplifies release planning.

863 In [42] software product investments are prioritized with HCV. The re-
864 sults of ECV was different for uncompensated and compensated HCV. When
865 compensated HCV was used ECV detected equal items that belonged to dif-
866 ferent high level prioritization groups (*A*, *B* and *C*), indicating that ECV
867 provided a more fine-grained view. In the case of uncompensated HCV, on
868 the other hand, all equal items belonged to one high level prioritization group
869 (group *B*).

870 8. Discussion and Conclusions

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

873 CV has been applied in various areas, but most frequently in requirements
874 prioritization and release planning, and quite often also as part of research
875 methodologies. A large part of the studies have been conducted in Sweden,
876 at Ericsson AB. One can see a slight increase in the interest in CV. During

877 the last five years there have been more studies that use CV than between,
878 say, year 2000–2005.

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

884 *8.1. Implications for Practitioners*

885 The results of this study provide decision support for industry practi-
886 tioners. We believe that a collection of state of the practice studies help
887 the adoption of CV prioritization method. (Top studies are summarized in
888 Table 12.) In addition, a set of CV analysis methods enables comprehen-
889 sive understanding of the prioritization results. (The analysis methods are
890 presented in Table 13.) One of the most common goals of CV analysis is to
891 display of the prioritization results and, thus, to show the difference between
892 several prioritization perspectives.

893 Additionally, we present ECV—a novel method for CV analysis. Priori-
894 tization often results in the assignment of similar priorities to several prior-
895 itization items. CV results contain both ‘real priorities’ and random errors.
896 Due to random errors, equal prioritization items may receive different pri-
897 orities. ECV identifies such items. It allows stakeholders to disregard the
898 random part of the CV results. Thus, ECV simplifies the understanding of
899 the prioritization results.

900 ECV identifies prioritization items with similar priority and tests whether
901 these items can be considered equal. In this case, ECV can be used in
902 software release planning. For example, let us suppose that a set of software
903 requirements are prioritized with regard to the implementation costs. First
904 of all, ECV can then detect items with equal cost. Second, the equal items
905 can be freely swapped between the releases. Finally, the decision to allocate
906 a requirement to a particular release can be made based on another criteria,
907 such as risk or business value.

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

911 *8.2. Implications for Academia*

912 In the systematic review 36% of papers were revealed by the snowball
913 sampling. That is a considerable amount. Several studies do not mention the
914 name of the prioritization method (i.e. cumulative voting or hundred dollar

915 test). Others are not available through selected databases because they are
916 conference publications or theses. It shows, in our opinion, that snowball
917 sampling ought to be used in all systematic literature reviews.

918 CV results are a special type of data—compositional data. Standard
919 statistical analysis methods that assume the independence of the samples
920 cannot be applied to CV results. In [57] methods for the analysis of com-
921 positional data have been presented. The systematic review conducted as a
922 part of this study revealed that 22 studies analyze CV results; yet, only one
923 study uses compositional data analysis methods, i.e. [10]. None of the stud-
924 ies, including [10], present methods for detecting items with equal priority
925 in CV results. Hence, ECV is, in this respect, a unique method.

926 The small use of compositional data analysis is really not surprising, since
927 literature describing CV does not state that the results are compositional
928 data. Standard statistical analysis methods may produce useful results for
929 compositional data. However, there are cases when they are misleading or
930 even faulty. Section 7.2.1 contains evidence of inappropriate use of statistical
931 methods by several papers.

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

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

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

- 940 1. All systematic literature reviews should include snowball sampling.
- 941 2. Researchers can improve their statistical analysis of CV results using
942 compositional data analysis methods collected and developed by this
943 study.
- 944 3. When CV or any other ratio scale prioritization method is taught,
945 compositional data analysis should also be presented as part of the
946 solution.

947 8.3. *Validity Threats*

948 The validity of the systematic review is mainly limited by the chosen
949 databases, the design of the review, and human judgement in study selection
950 and data extraction.

951 To mitigate the threats we use the most popular databases in the field
952 of software engineering. In the beginning of the systematic review a re-
953 view protocol was developed, peer-reviewed, and revised. Search strategy

954 was validated against a set of previously known papers obtained from other
955 researchers. One of many terms used to name cumulative voting is ‘\$100
956 method’. We were not able to search for this term because none of the cho-
957 sen databases support search for special characters like ‘\$’ and the search
958 string ‘100 method’ yields hundreds of thousands of results. To increase the
959 likelihood of discovering relevant studies snowball sampling was extensively
960 used.

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

964 The large number of studies identified by snowball sampling (15 out of
965 40 studies) may be caused by faulty design or by faulty execution of the
966 search in the databases. There are several reasons why the studies revealed
967 by snowball sampling are not revealed by the search in databases. Reason
968 for each study is given in Table Appendix A.2. Based on these reasons we
969 argue that snowball sampling does not indicate any problems with the design
970 of the search in the databases.

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

980 The quality of the data extraction and quality evaluation was validated
981 using inter-rater agreement analysis. In our case, 10% of the studies were
982 rated by two researchers and Krippendorff’s alpha was calculated. The agree-
983 ment for the data extraction results was 0.86 and the agreement for the
984 quality evaluation was 0.73 (indicating a credible level of quality).

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

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

992 To mitigate the second threat we visually inspected the results of the
993 application of ECV on the real world data from the primary studies. We

concluded that items identified by ECV can be considered equal.

CV results used in the evaluation of ECV were tested for normality. The tests indicated that CV results do not have multivariate normal distribution. Therefore, the design of ECV was based on a non-parametric statistical test.

8.4. Future Research

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

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

8.5. Conclusions

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

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

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

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

References

- [1] P. Berander, A. Andrews, Requirements Prioritization, in: A. Aurum, C. Wohlin (Eds.), Engineering and Managing Software Requirements,

- 1030 Springer-Verlag, Berlin/Heidelberg, 2005, pp. 69–94. doi:10.1007/3-
1031 540-28244-0.
1032 URL [http://www.springerlink.com/index/10.1007/](http://www.springerlink.com/index/10.1007/3-540-28244-0)
1033 [3-540-28244-0](http://www.springerlink.com/index/10.1007/3-540-28244-0)
- 1034 [2] D. Leffingwell, D. Widrig, Managing software requirements: A unified
1035 approach (1999) 118–119.
1036 URL <http://portal.acm.org/citation.cfm?id=326459>
- 1037 [3] V. Ahl, An experimental comparison of five prioritization methods, Master's
1038 Thesis, School of Engineering, Blekinge Institute of Technology.
- 1039 [4] P. Berander, P. Jonsson, Hierarchical Cumulative Voting (HCV) prior-
1040 itization of requirements in hierarchies (2006).
1041 URL <http://dx.doi.org/10.1142/S0218194006003026>[http://www.](http://www.worldscinet.com/ijseke/16/1606/S0218194006003026.html)
1042 [worldscinet.com/ijseke/16/1606/S0218194006003026.html](http://www.worldscinet.com/ijseke/16/1606/S0218194006003026.html)
- 1043 [5] J. Karlsson, K. Ryan, A cost-value approach for prioritizing require-
1044 ments, IEEE Software 14 (5) (1997) 67–74. doi:10.1109/52.605933.
- 1045 [6] J. Karlsson, An evaluation of methods for prioritizing software require-
1046 ments, Information and Software Technology 39 (14-15) (1998) 939–947.
1047 doi:10.1016/S0950-5849(97)00053-0.
1048 URL [http://dx.doi.org/10.1016/S0950-5849\(97\)00053-0](http://dx.doi.org/10.1016/S0950-5849(97)00053-0)
- 1049 [7] F. Pettersson, M. Ivarsson, T. Gorschek, P. Öhman, A practitioner's
1050 guide to light weight software process assessment and improvement plan-
1051 ning.
1052 URL <http://portal.acm.org/citation.cfm?id=1363376.1363636>
- 1053 [8] S. Barney, C. Wohlin, Software Product Quality: Ensuring a Common
1054 Goal, in: Q. Wang, V. Garousi, R. Madachy, D. Pfahl (Eds.), Trust-
1055 worthy Software Development Processes, Vol. 5543 of Lecture Notes in
1056 Computer Science, Springer Berlin Heidelberg, Berlin, Heidelberg, 2009,
1057 pp. 256–267. doi:10.1007/978-3-642-01680-6.
1058 URL <http://www.springerlink.com/content/j140v26514t7276u/>
- 1059 [9] P. Jönsson, C. Wohlin, A study on prioritisation of impact analysis
1060 issues: A comparison between perspectives, Software Engineering Re-
1061 search and Practice in Sweden.
1062 URL <http://www.wohlin.eu/Articles/SERPS05.pdf>

- 1063 [10] P. Chatzipetrou, L. Angelis, P. Rovegard, C. Wohlin, Prioritization of
1064 Issues and Requirements by Cumulative Voting: A Compositional Data
1065 Analysis Framework, 2010, pp. 361–370. doi:10.1109/SEAA.2010.35.
- 1066 [11] R. Engstrom, Cumulative Voting as a Remedy for Minority Vote Dilu-
1067 tion, Local Government Election
- 1068 [12] S. Bhagat, J. Brickley, Cumulative voting: The value of minority share-
1069 holder voting rights, Journal of Law and Economics.
- 1070 [13] V. Hiltunen, J. Kangas, J. Pykalainen, Voting methods in strategic for-
1071 est planning - Experiences from Metsahallitus, Forest Policy and Eco-
1072 nomics 10 (3) (2008) 117–127.
- 1073 [14] P. Boldi, F. Bonchi, C. Castillo, S. Vigna, Voting in social net-
1074 works, CIKM '09, ACM Press, New York, New York, USA, 2009.
1075 doi:10.1145/1645953.1646052.
1076 URL <http://portal.acm.org/citation.cfm?doid=1645953.1646052>
- 1077 [15] H. Ayad, M. Kamel, Cumulative Voting Consensus Method for Par-
1078 titions with Variable Number of Clusters, Pattern Analysis and Ma-
1079 chine Intelligence, IEEE Transactions on 30 (1) (2008) 160–173.
1080 doi:10.1109/TPAMI.2007.1138.
- 1081 [16] M. Svahnberg, A. Karasira, A Study on the Importance of Order in
1082 Requirements Prioritisation, IEEE, 2009. doi:10.1109/IWSPM.2009.1.
1083 URL [http://ieeexplore.ieee.org/lpdocs/epic03/wrapper.htm?](http://ieeexplore.ieee.org/lpdocs/epic03/wrapper.htm?arnumber=5457322)
1084 [arnumber=5457322](http://ieeexplore.ieee.org/lpdocs/epic03/wrapper.htm?arnumber=5457322)
- 1085 [17] P. Berander, M. Svahnberg, Evaluating two ways of calculating priorities
1086 in requirements hierarchies - An experiment on hierarchical cumulative
1087 voting (2009).
- 1088 [18] T. Saaty, The analytic hierarchy process., McGraw-Hill, New York.
1089 URL [http://scholar.google.se/scholar?hl=en&q=analytic+](http://scholar.google.se/scholar?hl=en&q=analytic+hierarchy+process+mcgraw+1980&btnG=Search&as_sdt=0,5&as_ylo=&as_vis=0\#4)
1090 [hierarchy+process+mcgraw+1980&btnG=Search&as_sdt=0,5\](http://scholar.google.se/scholar?hl=en&q=analytic+hierarchy+process+mcgraw+1980&btnG=Search&as_sdt=0,5&as_ylo=&as_vis=0\#4)
1091 [&as_ylo=&as_vis=0\#4](http://scholar.google.se/scholar?hl=en&q=analytic+hierarchy+process+mcgraw+1980&btnG=Search&as_sdt=0,5&as_ylo=&as_vis=0\#4)
- 1092 [19] S. Brenner, J. Schwalbach, Legal Institutions, Board Diligence, and Top
1093 Executive Pay, Corporate Governance: An International Review 17 (1)
1094 (2009) 1–12. doi:10.1111/j.1467-8683.2008.00720.x.
1095 URL <http://doi.wiley.com/10.1111/j.1467-8683.2008.00720.x>

- 1096 [20] J. Aitchison, J. J. Egozcue, Compositional Data Analysis: Where Are
1097 We and Where Should We Be Heading?, *Mathematical Geology* 37 (7)
1098 (2005) 829–850. doi:10.1007/s11004-005-7383-7.
1099 URL [http://www.springerlink.com/index/10.1007/
1100 s11004-005-7383-7](http://www.springerlink.com/index/10.1007/s11004-005-7383-7)
- 1101 [21] V. Pawlowsky-Glahn, J. J. Egozcue, Compositional data and their
1102 analysis: an introduction, Geological Society, London, Special Publica-
1103 tions 264 (1) (2006) 1–10. doi:10.1144/GSL.SP.2006.264.01.01.
1104 URL [http://sp.lyellcollection.org/cgi/doi/10.1144/GSL.SP.
1105 2006.264.01.01](http://sp.lyellcollection.org/cgi/doi/10.1144/GSL.SP.2006.264.01.01)
- 1106 [22] J. Martin-Fernandez, C. Barceló-Vidal, V. Pawlowsky-Glahn, Dealing
1107 with zeros and missing values in compositional data sets using nonpara-
1108 metric imputation, *Mathematical Geology* 35 (3) (2003) 253–278.
1109 URL <http://www.springerlink.com/index/ku816485q4264772.pdf>
- 1110 [23] P. Filzmoser, K. Hron, Outlier detection for compositional data using
1111 robust methods Outlier Detection for Compositional Data Using Robust
1112 Methods, *Analysis and Applications* (April).
- 1113 [24] K. Khan, A systematic review of software requirements prioritization,
1114 Unpublished master’s thesis, Blekinge Institute of Technology, Ronneby,
1115 Sweden (October).
1116 URL [http://citeseerx.ist.psu.edu/viewdoc/download?doi=10.1.
1117 1.107.8608&rep=rep1&type=pdf](http://citeseerx.ist.psu.edu/viewdoc/download?doi=10.1.1.107.8608&rep=rep1&type=pdf)
- 1118 [25] F. Zahedi, The analytic hierarchy process: a survey of the method and
1119 its applications, *Interfaces* (1986) 96–108.
1120 URL <http://www.jstor.org/stable/25060854>
- 1121 [26] P. Runeson, M. Höst, Guidelines for conducting and reporting case
1122 study research in software engineering, *Empirical Software Engineering*
1123 14 (2) (2008) 131–164. doi:10.1007/s10664-008-9102-8.
1124 URL [http://www.springerlink.com/index/10.1007/
1125 s10664-008-9102-8](http://www.springerlink.com/index/10.1007/s10664-008-9102-8)
- 1126 [27] L. Goodman, Snowball sampling, *The Annals of Mathematical Statis-*
1127 *tics*.
1128 URL <http://www.jstor.org/stable/2237615>
- 1129 [28] K. Krippendorff, Bivariate agreement coefficients for reliability of data,
1130 *Sociological methodology*.

1131 URL [http://scholar.google.se/scholar?hl=en&q=Bivariate+](http://scholar.google.se/scholar?hl=en&q=Bivariate+Agreement+Coefficients+for+Reliability+of+Data&btnG=Search&as_sdt=0,5&as_ylo=&as_vis=0\#0)
1132 [Agreement+Coefficients+for+Reliability+of+Data\&btnG=](http://scholar.google.se/scholar?hl=en&q=Bivariate+Agreement+Coefficients+for+Reliability+of+Data&btnG=Search&as_sdt=0,5&as_ylo=&as_vis=0\#0)
1133 [Search\&as_sdt=0,5\&as_ylo=&as_vis=0\#0](http://scholar.google.se/scholar?hl=en&q=Bivariate+Agreement+Coefficients+for+Reliability+of+Data&btnG=Search&as_sdt=0,5&as_ylo=&as_vis=0\#0)

1134 [29] K. Krippendorff, Content analysis: An introduction to its methodology.
1135 URL [http://scholar.google.se/scholar?hl=en&q=Krippendorff,](http://scholar.google.se/scholar?hl=en&q=Krippendorff,+K+2004&btnG=Search&as_sdt=0,5&as_ylo=&as_vis=0\#0)
1136 [+K+2004\&btnG=Search\&as_sdt=0,5\&as_ylo=&as_vis=0\#0](http://scholar.google.se/scholar?hl=en&q=Krippendorff,+K+2004&btnG=Search&as_sdt=0,5&as_ylo=&as_vis=0\#0)

1137 [30] B. Kitchenham, Guidelines for performing systematic literature reviews
1138 in software engineering, Engineering.

1139 [31] P. Berander, P. Jönsson, A goal question metric based approach for effi-
1140 cient measurement framework definition, ACM, Rio de Janeiro, Brazil,
1141 2006, pp. 316–325. doi:10.1145/1159733.1159781.

1142 [32] B. Regnell, M. Höst, J. och Dag, An industrial case study on distributed
1143 prioritisation in market-driven requirements engineering for packaged
1144 software, Requirements
1145 URL <http://www.springerlink.com/index/JG9G7KXALXYRT43B.pdf>

1146 [33] B. Kitchenham, Procedures for performing systematic reviews, Keele,
1147 UK, Keele University 33.

1148 [34] M. Ivarsson, T. Gorschek, A method for evaluating rigor and indus-
1149 trial relevance of technology evaluations, Empirical Software Engineer-
1150 ing (2010) 1–31.
1151 URL <http://www.springerlink.com/index/116531105174V25N.pdf>

1152 [35] C. Wohlin, P. Runeson, M. Höst, Experimentation in software engi-
1153 neering: an introduction, Springer Netherlands, 2000.
1154 URL [http://books.google.com/books?hl=en&lr=](http://books.google.com/books?hl=en&lr=&id=nG2UShV0wAEC&oi=fnd&pg=PR11&dq=Experimentation+in+software+engineering:+an+introduction&ots=9Gb9RW7j-l&sig=tKC8wLE4NShrt_XymaJq-7oKpRE)
1155 [&id=nG2UShV0wAEC\&oi=fnd\&pg=PR11\&dq=](http://books.google.com/books?hl=en&lr=&id=nG2UShV0wAEC&oi=fnd&pg=PR11&dq=Experimentation+in+software+engineering:+an+introduction&ots=9Gb9RW7j-l&sig=tKC8wLE4NShrt_XymaJq-7oKpRE)
1156 [Experimentation+in+software+engineering:+an+introduction\](http://books.google.com/books?hl=en&lr=&id=nG2UShV0wAEC&oi=fnd&pg=PR11&dq=Experimentation+in+software+engineering:+an+introduction&ots=9Gb9RW7j-l&sig=tKC8wLE4NShrt_XymaJq-7oKpRE)
1157 [&ots=9Gb9RW7j-l\&sig=tKC8wLE4NShrt_XymaJq-7oKpRE](http://books.google.com/books?hl=en&lr=&id=nG2UShV0wAEC&oi=fnd&pg=PR11&dq=Experimentation+in+software+engineering:+an+introduction&ots=9Gb9RW7j-l&sig=tKC8wLE4NShrt_XymaJq-7oKpRE)

1158 [36] A. Jedlitschka, D. Pfahl, Reporting guidelines for controlled experi-
1159 ments in software engineering, in: 2005 International Symposium on
1160 Empirical Software Engineering, 2005., IEEE, 2005, p. 10.
1161 URL [http://www.computer.org/portal/web/csdl/doi/10.1109/](http://www.computer.org/portal/web/csdl/doi/10.1109/ISESE.2005.1541818)
1162 [ISESE.2005.1541818](http://www.computer.org/portal/web/csdl/doi/10.1109/ISESE.2005.1541818)

1163 [37] K. Rinkevics, Data Extraction and Quality Evaluation results (2011).
1164 URL [http://rinkevic.wordpress.com/2011/11/26/](http://rinkevic.wordpress.com/2011/11/26/data-extraction-and-quality-evaluation-results/)
1165 [data-extraction-and-quality-evaluation-results/](http://rinkevic.wordpress.com/2011/11/26/data-extraction-and-quality-evaluation-results/)

- 1166 [38] R. Ihaka, R. Gentleman, R: a language for data analysis and graphics,
1167 Journal of computational and graphical statistics (1996) 299–314.
1168 URL <http://www.jstor.org/stable/1390807>
- 1169 [39] K. Rinkevics, ECV implementation source code (2011).
1170 URL [http://rinkevics.wordpress.com/2011/08/14/](http://rinkevics.wordpress.com/2011/08/14/ecv-implementation-in-r/)
1171 [ecv-implementation-in-r/](http://rinkevics.wordpress.com/2011/08/14/ecv-implementation-in-r/)
- 1172 [40] R. M. Groves, F. J. Fowler, M. P. Couper, J. M. Lepkowski, E. Singer,
1173 Survey methodology, John Wiley and Sons, 2009.
1174 URL <http://books.google.com/books?id=HXoSpXvo3s4C>
- 1175 [41] S. Barney, A. Aurum, C. Wohlin, The Relative Importance of Aspects
1176 of Intellectual Capital for Software Companies, in: 2009 35th Euromi-
1177 cro Conference on Software Engineering and Advanced Applications,
1178 IEEE, 2009, pp. 313–320. doi:10.1109/SEAA.2009.44.
1179 URL [http://ieeexplore.ieee.org/lpdocs/epic03/wrapper.htm?](http://ieeexplore.ieee.org/lpdocs/epic03/wrapper.htm?arnumber=5349937)
1180 [arnumber=5349937](http://ieeexplore.ieee.org/lpdocs/epic03/wrapper.htm?arnumber=5349937)
- 1181 [42] S. Barney, C. Wohlin, A. Aurum, Balancing software product invest-
1182 ments, IEEE Computer Society, 2009, pp. 257–268.
- 1183 [43] S. Barney, A. Aurum, C. Wohlin, A product management chal-
1184 lenge: Creating software product value through requirements
1185 selection, Journal of Systems Architecture 54 (6) (2008) 576–593.
1186 doi:10.1016/j.sysarc.2007.12.004.
1187 URL [http://linkinghub.elsevier.com/retrieve/pii/](http://linkinghub.elsevier.com/retrieve/pii/S1383762107001348)
1188 [S1383762107001348](http://linkinghub.elsevier.com/retrieve/pii/S1383762107001348)
- 1189 [44] S. Laukkanen, T. Palander, J. Kangas, A. Kangas, Evaluation of the
1190 multicriteria approval method for timber-harvesting group decision sup-
1191 port, Silva Fennica 39 (2) (2005) 249–264.
- 1192 [45] G. Hu, Adding value to software requirements: An empirical study in
1193 the chinese software industry, Seventeenth Australian Conference on
1194
1195 URL [http://citeseerx.ist.psu.edu/viewdoc/download?doi=10.1.](http://citeseerx.ist.psu.edu/viewdoc/download?doi=10.1.1.107.1945\&rep=rep1\&type=pdf)
1196 [1.107.1945\&rep=rep1\&type=pdf](http://citeseerx.ist.psu.edu/viewdoc/download?doi=10.1.1.107.1945\&rep=rep1\&type=pdf)
- 1197 [46] R. Feldt, R. Torkar, E. Ahmad, B. Raza, Challenges with Software
1198 Verification and Validation Activities in the Space Industry, IEEE,
1199 2010. doi:10.1109/ICST.2010.37.

- 1200 URL [http://ieeexplore.ieee.org/lpdocs/epic03/wrapper.htm?](http://ieeexplore.ieee.org/lpdocs/epic03/wrapper.htm?arnumber=5477080)
1201 [arnumber=5477080](http://ieeexplore.ieee.org/lpdocs/epic03/wrapper.htm?arnumber=5477080)
- 1202 [47] M. Svahnberg, T. Gorschek, M. Eriksson, A. Borg, K. Sandahl,
1203 J. Börstler, A. Loconsole, Perspectives on Requirements Understand-
1204 ability – For Whom Does the Teacher’s Bell Toll?, IEEE, 2008.
1205 doi:10.1109/REET.2008.4.
1206 URL [http://ieeexplore.ieee.org/lpdocs/epic03/wrapper.htm?](http://ieeexplore.ieee.org/lpdocs/epic03/wrapper.htm?arnumber=4797459)
1207 [arnumber=4797459](http://ieeexplore.ieee.org/lpdocs/epic03/wrapper.htm?arnumber=4797459)
- 1208 [48] P. Jönsson, C. Wohlin, Understanding impact analysis: An empir-
1209 ical study to capture knowledge on different organisational levels,
1210 ... Conference on Software Engineering and Knowledge
1211 URL <http://wohlin.eu/Articles/SEKE05.pdf>
- 1212 [49] L. a. Kuzniarz, Empirical extension of a classification framework for
1213 addressing consistency in model based development, Information and
1214 Software Technologydoi:10.1016/j.infsof.2010.10.004.
1215 URL [http://www.scopus.com/inward/record.url?](http://www.scopus.com/inward/record.url?eid=2-s2.0-78650489358&partnerID=40&md5=9a8d2b6e973700e4cd68106471759b10)
1216 [eid=2-s2.0-78650489358&partnerID=40&md5=](http://www.scopus.com/inward/record.url?eid=2-s2.0-78650489358&partnerID=40&md5=9a8d2b6e973700e4cd68106471759b10)
1217 [9a8d2b6e973700e4cd68106471759b10](http://www.scopus.com/inward/record.url?eid=2-s2.0-78650489358&partnerID=40&md5=9a8d2b6e973700e4cd68106471759b10)
- 1218 [50] P. Rovegard, L. Angelis, C. Wohlin, An Empirical Study on Views of
1219 Importance of Change Impact Analysis Issues, Software Engineering,
1220 IEEE Transactions on 34 (4) (2008) 516 –530. doi:10.1109/TSE.2008.32.
- 1221 [51] C. Wohlin, A. Aurum, Criteria for selecting software requirements to
1222 create product value: An industrial empirical study, Value-Based Soft-
1223 ware Engineering.
1224 URL <http://www.wohlin.eu/Articles/VBSE05.pdf>
- 1225 [52] B. Regnell, M. Höst, J. Natt, Visualization of Agreement and Satisfac-
1226 tion in Distributed Prioritization of Market Requirements, Chart (2000)
1227 1–12.
- 1228 [53] J. Aitchison, Principal component analysis of compositional data,
1229 Biometrika 70 (1) (1983) 57. doi:10.2307/2335943.
1230 URL <http://biomet.oxfordjournals.org/content/70/1/57.short>
- 1231 [54] P. Filzmoser, K. Hron, C. Reimann, F. Sm, P. Filzmoser, K. Hron,
1232 C. Reimann, Principal component analysis for compositional data with
1233 outliers Principal component analysis for compositional data with out-
1234 liers, Analysis and Applications (November).

- 1235 [55] V. Pawlowsky Glahn, J. Egozcue, R. Tolosana Delgado, Lecture notes
1236 on compositional data analysis, Interpretation A Journal Of Bible And
1237 Theology (July).
1238 URL <http://dugi-doc.udg.edu/handle/10256/297>
- 1239 [56] W. Kruskal, W. Wallis, Use of ranks in one-criterion variance analysis,
1240 Journal of the American statistical Association 47 (260) (1952) 583–621.
1241 URL <http://www.jstor.org/stable/2280779>
- 1242 [57] J. Aitchison, The statistical analysis of compositional data, Chapman
1243 & Hall, London, 1986.
- 1244 [58] D. Baca, K. Petersen, Prioritizing Countermeasures through the Counter-
1245 measure Method for Software Security (CM-Sec), in: M. Ali Babar,
1246 M. Vierimaa, M. Oivo (Eds.), Product-Focused Software Process Im-
1247 provement, Vol. 6156 of Lecture Notes in Computer Science, Springer
1248 Berlin / Heidelberg, 2010, pp. 176–190.
1249 URL http://dx.doi.org/10.1007/978-3-642-13792-1_15
- 1250 [59] S. a. b. Bowler, Election systems and voter turnout: Experiments in
1251 the United States, Journal of Politics 63 (3) (2001) 902–915.
1252 URL [http://www.scopus.com/inward/record.
1253 url?eid=2-s2.0-0035536318&partnerID=40&md5=
1254 517d9a827ee1af7860e2e4939693c4de](http://www.scopus.com/inward/record.url?eid=2-s2.0-0035536318&partnerID=40&md5=517d9a827ee1af7860e2e4939693c4de)
- 1255 [60] D. Brockington, A Low Information Theory of Ballot Position Effect,
1256 Political Behavior 25 (1) (2003) 1–27. doi:10.1023/A:1022946710610.
1257 URL <http://www.springerlink.com/content/x522750t32296220/>
- 1258 [61] D. Cooper, A. Zillante, A comparison of cumulative voting and gener-
1259 alized plurality voting, Public Choice doi:10.1007/s11127-010-9707-5.
1260 URL <http://www.springerlink.com/content/145774u78052x863/>
- 1261 [62] N. D. Fogelström, M. Svahnberg, T. Gorschek, Investigating Impact
1262 of Business Risk on Requirements Selection Decisions, IEEE, 2009.
1263 doi:10.1109/SEAA.2009.66.
1264 URL [http://ieeexplore.ieee.org/xpl/freeabs_all.jsp?
1265 arnumber=5349849](http://ieeexplore.ieee.org/xpl/freeabs/_all.jsp?arnumber=5349849)
- 1266 [63] S. Hatton, Choosing the Right Prioritisation Method, in: Proceedings
1267 of the 19th Australian Conference on Software Engineering, IEEE Com-
1268 puter Society, Washington, 2008, pp. 517–526.
1269 URL <http://portal.acm.org/citation.cfm?id=1395083.1395703>

- [64] S. Hatton, Early prioritisation of goals, in: Proceedings of the 2007 conference on Advances in conceptual modeling: foundations and applications, ER'07, Springer-Verlag, Berlin, 2007, pp. 235–244.
URL <http://portal.acm.org/citation.cfm?id=1784542.1784583>
- [65] V. Heikkilä, A. Jadallah, K. Rautiainen, G. Ruhe, Rigorous Support for Flexible Planning of Product Releases - A Stakeholder-Centric Approach and Its Initial Evaluation, IEEE, 2010. doi:10.1109/HICSS.2010.323.
URL <http://ieeexplore.ieee.org/lpdocs/epic03/wrapper.htm?arnumber=5428538>
- [66] M. Staron, C. Wohlin, An Industrial Case Study on the Choice Between Language Customization Mechanisms, in: J. Münch, M. Vierimaa (Eds.), Product-Focused Software Process Improvement, Vol. 4034 of Lecture Notes in Computer Science, Springer Berlin / Heidelberg, 2006, pp. 177–191.
URL http://dx.doi.org/10.1007/11767718_17
- [67] T. Touseef, C. Gancel, A structured goal based measurement framework enabling traceability and prioritization, ... (ICET), 2010 6th International Conference on.
URL http://ieeexplore.ieee.org/xpls/abs/_all.jsp?arnumber=5638475
- [68] P. Berander, C. Wohlin, Differences in views between development roles in software process improvement-a quantitative comparison, in: Proceedings 8th Conference on Empirical Assessment in Software Engineering, 2004.
URL <http://www.wohlin.eu/Articles/EASE04-2.pdf>
- [69] P. Berander, Using students as subjects in requirements prioritization, Proceedings. 2004 International Symposium on Empirical Software Engineering, 2004. ISESE '04. (2004) 167–176doi:10.1109/ISESE.2004.1334904.
URL <http://ieeexplore.ieee.org/lpdocs/epic03/wrapper.htm?arnumber=1334904>
- [70] P. Berander, C. Wohlin, Identification of Key Factors in Software Process Management-A Case Study.
URL <http://www.computer.org/portal/web/csdl/doi/10.1109/ISESE.2003.1237992>

- 1306 [71] R. L. Cole, D. a. Taebel, R. L. Engstrom, Cumulative Voting in a Munic-
 1307 ipal Election: A Note on Voter Reactions and Electoral Consequences,
 1308 The Western Political Quarterly 43 (1) (1990) 191. doi:10.2307/448513.
 1309 URL <http://www.jstor.org/stable/448513?origin=crossref>
- 1310 [72] J. Kuklinski, Cumulative and Plurality Voting: An Analysis of Illinois'
 1311 Unique Electoral System, The Western Political Quarterly 26 (4) (1973)
 1312 726–746.
 1313 URL <http://www.jstor.org/stable/447147>
- 1314 [73] S. Laukkanen, T. Palander, J. Kangas, Applying voting theory in par-
 1315 ticipatory decision support for sustainable timber harvesting, Canadian
 1316 Journal of Forest Research 34 (7) (2004) 1511–1524. doi:10.1139/x04-
 1317 044.
 1318 URL [http://article.pubs.nrc-cnrc.gc.ca/ppv/RPViewDoc?issn=](http://article.pubs.nrc-cnrc.gc.ca/ppv/RPViewDoc?issn=1208-6037&volume=34&issue=7&startPage=1511&ab=y)
 1319 [1208-6037&volume=34&issue=7&startPage=1511&ab=y](http://article.pubs.nrc-cnrc.gc.ca/ppv/RPViewDoc?issn=1208-6037&volume=34&issue=7&startPage=1511&ab=y)
- 1320 [74] J. Sawyer, D. MacRae, Game theory and cumulative voting in Illinois:
 1321 1902-1954, The American Political Science Review 56 (4) (1962) 936–
 1322 946.
 1323 URL <http://www.jstor.org/stable/1952795>

1324 Appendix A. Primary Studies

1325 Appendix A.1. Primary studies found during search in databases.

1326	Title	Reference
	Prioritizing countermeasures through the countermeasure method for software security (CM-Sec)	Baca 2010 [58]
	The relative importance of aspects of intellectual capital for software companies	Barney 2009 [41]
	Software product quality: Ensuring a common goal	Barney 2009 [8]
	Balancing software product investments	Barney 2009 [42]
	Hierarchical cumulative voting (HCV) prioritization of requirements in hierarchies	Berander 2006 [4]
	A goal question metric based approach for efficient measurement framework definition	Berander 2006 [31]
	Evaluating two ways of calculating priorities in requirements hierarchies: An experiment on hierarchical cumulative voting	Berander 2009 [17]
	Election systems and voter turnout: Experiments in the United States	Bowler 2001 [59]
	A low information theory of ballot position effect	Brockington 2003 [60]
	Prioritization of issues and requirements by cumulative Voting: A compositional data analysis framework	Chatzipetrou 2010 [10]
	A comparison of cumulative voting and generalized plurality voting	Cooper 2010 [61]
	Challenges with software verification and validation activities in the space industry	Feldt 2010 [46]
	Investigating impact of business risk on requirements selection decisions	Fogelstrom 2009 [62]
	Choosing the right prioritization method	Hatton 2008 [63]
	Early prioritization of goals	Hatton 2007 [64]
	Rigorous support for flexible planning of product releases: A stakeholder-centric approach and its initial evaluation	Heikkilä 2010 [65]
	Voting methods in strategic forest planning: Experiences from Metsähallitus	Hiltunen 2008 [13]
	Empirical extension of a classification framework for addressing consistency in model based development	Kuzniarz 2010 [49]
	Evaluation of the multi-criteria approval method for timber-harvesting group decision support	Laukkanen 2005 [44]
	A practitioner's guide to light weight software process assessment and improvement planning	Pettersson 2008 [7]
	An empirical study on views of importance of change impact analysis issues	Rovegard 2008 [50]
	An industrial case study on the choice between language customization mechanisms	Staron 2006 [66]
	Perspectives on requirements understandability—For whom does the teacher's bell toll?	Svahnberg 2008 [47]
	A study on the importance of order in requirements prioritization	Svahnberg 2009 [16]
	A structured goal based measurement framework enabling traceability and prioritization	Touseef 2010 [67]

1327 *Appendix A.2. Primary studies revealed by snowball sampling.*

1328

Reference	Title	Reason why the paper is not revealed by the search in databases
Ahl 2005 [3]	An experimental comparison of five prioritization methods	Selected databases does not contain the paper, master thesis at BTH
Barney 2008 [43]	A product management challenge: Creating software product value through requirements selection	Prioritization method name not mentioned, phrase "1,000 points" used instead.
Berander 2004 [68]	Differences in views between development roles in software process improvement—A quantitative comparison	Prioritization method name not mentioned, phrase "100 points" used instead.
Berander 2004 [69]	Using students as subjects in requirements prioritization	Unknown CV name: 100-point technique
Berander 2003 [70]	Identification of key factors in software process management: A case study	Prioritization method name not mentioned, phrase "100 points" used instead.
Cole 1990 [71]	Cumulative voting in a municipal election: A note on voter reactions and electoral consequences	Study published before year 2001.
Hu 2006 [45]	Adding value to software requirements: An empirical study in the chinese software industry	Prioritization method name not mentioned, phrase "1,000 points" used instead.
Jonsson 2005 [9]	A study on prioritization of impact analysis issues: A comparison between perspectives	Selected databases does not contain the paper.
Jonsson 2005 [48]	Understanding impact analysis: An empirical study to capture knowledge on different organizational levels	Selected databases does not contain the paper.
Kuklinski 1973 [72]	Cumulative and plurality voting: An analysis of Illinois' unique electoral system	Study published before year 2001.
Laukkanen 2004 [73]	Applying voting theory in participatory decision support for sustainable timber harvesting	Selected databases does not contain the paper.
Regnell 2001 [32]	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 2000 [52]	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."
Wohlin 2006 [74]	Game theory and cumulative voting in Illinois: 1902–1954	Study published before year 2001.
Wohlin 2006 [51]	Criteria for selecting software requirements to create product value: An industrial empirical study	Prioritization method name not mentioned: "The subjects had 1,000 points to spend among the 13 criteria."

Appendix B. CV Result Analysis Methods

	Paper																					
	Svalnberg2008	Svalnberg2009	Starr2006	Petersson2008	Wohlitz2006	Laukkonen2005a	Hu2006	Jonsson2005a	Kuzmar2010	Rowgard2008	Bernard2006a	Bernard2004a	Bernard2006	Feldt2010	Barney2009b	Barney2008	Barney2009a	Barney2009	Jonsson2005	Chatzipetrou2010	Reguel2001	Reguel2000
Analysis method																						
Table that shows final priorities				x																		
Chart that shows final priorities	x			x	x	x										x						
Table of top-5 prioritization items																						
min , max , \bar{x} , \bar{x} and σ of priorities assigned by different stakeholders									x	x												
Bar chart of prioritization results showing \bar{x} priority and σ of priorities									x	x												
Pearson correlation coefficient		x										x										
Nemenyi-Damico-Wolfe-Dunn														x								
Spearman's r															x							
Kruskal-Wallis								x								x		x				
Wilcoxon							x															
Correlation matrix		x														x		x				
Chart for comparing priorities from two perspectives, priorities are points in two dimensional plane, x - and y -axis represent two different perspectives										x									x			
Difference between priorities assigned by each two stakeholders using χ^2 -statistic										x												
Median ranks		x																				
CV results converted to priority ranks		x											x						x			
PCA				x	x															x		
Percentage of divergence of priorities assigned by a stakeholder											x											
Average percentage of items given non-zero value											x											
Distribution chart																					x	x
Disagreement chart				x																	x	x
Satisfaction chart			x																		x	x
Bi-plot																					x	
Ternary plot																					x	

Appendix C. Quality Evaluation Checklist

Item	Question or Description of the Item	Rating
1. Background, introduction	Introduce research area	
2. Problem statement, purpose	What is the problem [36]? Where does it occur [30]? Who has observed it [36]? Why is it important to be solved [36]?	
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[30]?	
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 [33].	
7.5. Balancing	Equal number of subjects should be assigned to each treatment [33].	
7.6. Blinding	Automated assignment of treatments to subjects [33] Automated distribution of study materials to subjects [33] Persons who grade the task results should not know which treatment was used [33] Analyst should not know which treatment group is which [33] Automated data collection from subjects [33]	
8. Subjects (participants)		
8.1. Population		
8.2. Sampling	How sampling is performed? What subjects are included and excluded? [30] 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[30]?	
8.4. Subject motivation	E.g. material benefits, course credits for students, etc.	
9. Objects	E.g. documents and other artifacts	
10. Measures, Data collection procedures	Who, when, and how to measure [30]? How is the measurement supported? Is it automated [30]? Are the measures used in the study the most relevant ones for answering the research questions [30]?	
11. Analysis procedure		
11.1. Data description	Do the numbers add up across different tables and subgroups [30]?	
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 [30]? If the study is concerned with differences among groups, are confidence limits given describing the magnitude of any observed differences [30]?	
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 [30]? Are negative findings presented [30]?	
14. Industry impact, inference, generalization	What implications does the report have for practice [30]? How and where the results can be used? Limitations under which findings are relevant [36]?	
15. Future work		