

**How good are ideas identified by an automatic idea detection system?**

Journal:	<i>Creativity and Innovation Management</i>
Manuscript ID	Draft
Manuscript Type:	Article
Keywords:	Idea quality, Text mining, Machine learning, Natural language processing, Big data, Online communities, Support vector machines, Crowdsourcing

SCHOLARONE™  
Manuscripts

Review

# How good are ideas identified by an automatic idea detection system?

## Abstract

Online communities are an attractive source of potential ideas for products and process'. Recent advances in machine learning have made it possible to screen the vast amounts of information in online communities and automatically detect user-contributed ideas. However, it is still uncertain whether the ideas identified by such a system will also be regarded as sufficiently novel, feasible and valuable by firms who might decide to develop them further. A validation study is reported in which 200 posts were extracted from an online community using the automatic idea detection system by Christensen, Nørskov, Frederiksen and Scholderer (2016; DOI: 10.1111/caim.12202). Two company professionals evaluated the posts in terms of idea content and idea quality. The results suggest that the automatic idea detection system is sufficiently valid to be deployed for the harvesting and initial screening of innovation ideas and that the profile of the identified ideas (in terms of novelty, feasibility and value) follows the same pattern identified in studies of user ideation in general.

## Introduction

The digitalisation of business life is progressing: more and more tasks can be solved by automated systems. Whilst in the past, these were predominantly tasks of a mundane and repetitive nature, recent advances in artificial intelligence have also made it possible to solve complex problems. A common problem during the introduction of such systems is that they can be intransparent to their prospective users. Whilst the traditional business processes they are intended to rationalise have often been in use for many years and are implicitly trusted by management and staff, newly introduced automated systems lack such a track record. Scepticism and reactance can be the consequence.

To earn the trust of prospective users, automated systems have to enable superior performance. Benchmarked against the traditional business processes they are intended to rationalise, they should lead to increases in effectiveness or efficiency. This is easily demonstrated in application areas such as sales forecasting or inventory control where commonly accepted and routinely measured performance criteria exist. Such criteria rarely exist in more complex and creative areas such as innovation management. The aim of the research presented here is to show how the performance of automated systems in such areas can be evaluated. We will demonstrate this in the context of a particular type of task: the automated detection of ideas for product and process innovations in the contributions to an online developer forum.

### *Online communities as idea reservoirs*

Firms need a continuous stream of ideas to fuel their innovation processes (Van de Ven, 1986; Ekvall, 1997; Vandenbosch, Saatcioglu, & Fay, 2006; van den Ende, Frederiksen, & Prencipe, 2015). Ideas do not have to originate from the creative mind of the firm's employees but can also originate from the users of its products, services and technologies

1  
2  
3 (Kristensson, Gustafsson, & Archer, 2004; Magnusson, 2009; von Hippel, Ogawa, & PJ de  
4  
5 Jong, 2011; Poetz & Schreier, 2012; Majchrzak & Malhotra, 2013; Magnusson, Wästlund, &  
6  
7 Netz, 2014). Online communities where users exchange experiences and discuss potential  
8  
9 improvements are a particularly rich reservoir of ideas for product and process innovations.  
10

11  
12 Prominent examples include the user communities hosted by Dell (di Gangi, Wasko,  
13  
14 & Hooker, 2010; Poetz & Schreier, 2012), Lego (Antorini, 2007; Antorini, Muñiz, &  
15  
16 Askildsen, 2012; Nørskov, Antorini, & Jensen, 2015), Propellerhead (Jeppesen &  
17  
18 Frederiksen, 2006) and IBM (Mahr & Lievens, 2012). Firm-hosted communities such as these  
19  
20 have the advantage that the hosting firm can retain a certain degree of control. The  
21  
22 communities are typically based on software that allows registered users to post ideas,  
23  
24 comment on and vote for ideas posted by other users in a highly structured manner. The  
25  
26 downside of this approach is that it requires an extensive base of committed product users or  
27  
28 firm-loyal customers who have an intrinsic interest in suggesting ideas to the firm.  
29  
30  
31

32  
33 However, users do not only gather in firm-hosted communities. A vast amount of  
34  
35 online communities exists that are firm-free (Füller, Bartl, Ernst, & Mühlbacher, 2006; Füller,  
36  
37 Jawecki, & Mühlbacher, 2007). The most prominent cases include open-source software  
38  
39 development communities such as those responsible for the Linux kernel, R and Python.  
40  
41 These are examples of firm-free “products” that have been developed in a distributed manner,  
42  
43 utilising online collaboration platforms such as GitHub and Sourceforge. The fact that the  
44  
45 resulting products are now perfectly able to compete with their commercial counterparts (such  
46  
47 as the products ranges of the SAS Institute or Microsoft) is a clear demonstration of the  
48  
49 potential of such communities (von Krogh, Spaeth, & Lakhani, 2003; von Krogh & von  
50  
51 Hippel, 2006)  
52  
53  
54  
55  
56  
57  
58  
59  
60

1  
2  
3 The problem with firm-free communities is that they, unlike most firm-hosted  
4 communities, are usually *not* based on a crowdsourcing architecture that would enable easy  
5 harvesting and collaborative filtering of the community-generated ideas. Assigning employees  
6 to manual monitoring of community contributions is often the only viable solution if firms  
7 want to benefit from the ideas generated in firm-free communities. This is time-consuming  
8 and expensive; online communities may contain several hundred thousand posts and  
9 comments. The sheer amount of information in which the ideas are hidden is a practical  
10 barrier to finding the ideas and utilising them for innovation (Lin, Hsieh, & Chuang, 2009;  
11 Thorleuchter & Van den Poel, 2013).  
12  
13  
14  
15  
16  
17  
18  
19  
20  
21  
22

### 23 *Automatic idea detection*

24  
25  
26  
27 A new and efficient way of solving the needle-in-a-haystack problem is to use classifi-  
28 cation algorithms that can screen arbitrary amounts of community posts and comments and  
29 identify those that are likely to contain ideas. Using natural language processing and machine  
30 learning methods, Christensen, Nørskov, Frederiksen, & Scholderer (2016) develop such an  
31 algorithm and demonstrate its classification performance and efficiency for the case of ex-  
32 tracting new product ideas from an online community related to Lego. Christensen et al.  
33 (Submitted manuscript) show that the same principles can be applied to extract ideas for in-  
34 novations from a community related to craft brewing.  
35  
36  
37  
38  
39  
40  
41  
42  
43  
44

45 The authors argue that their method is applicable across different technological areas  
46 and product categories because most people use a specific set of words and expressions when  
47 they communicate ideas to each other. Since the presence of such linguistic markers can easi-  
48 ly be detected in a given post or comment, it can also be exploited for the screening of arbi-  
49 trarily large collections of posts, comments or other types of semi- or unstructured text. Im-  
50 plemented as a screening tool in a firm's R&D or marketing department, it can significantly  
51  
52  
53  
54  
55  
56  
57  
58  
59  
60

1  
2  
3 reduce the labour costs that would arise if R&D staff were assigned to manual monitoring of  
4  
5 community activity.  
6  
7

### 8 *How good are automatically detected ideas?* 9

10  
11 A crucial question is whether the ideas detected by such an automated system would  
12  
13 also be seen as sufficiently novel, feasible and valuable by the R&D or marketing staff who  
14  
15 would have to decide if the ideas should be taken further (e.g., developed into concepts or  
16  
17 prototypes). Ideas identified by the Christensen et al. (2016) method, for example, have not  
18  
19 yet been evaluated by company-internal R&D or marketing staff. The aim of the present paper  
20  
21 is to fill this gap. Specifically, we would like to contribute in two respects to the literature:  
22  
23

- 24  
25 • Our first contribution is to assess whether ideas from an online community, identified by  
26  
27 an artificial intelligence system such as the one described by Christensen et al. (2016),  
28  
29 will also be perceived as ideas by company-internal staff.  
30  
31
- 32  
33 • Our second contribution is to investigate if the ideas that are detected by the system will  
34  
35 also be perceived as *good* ideas by company-internal staff.  
36  
37

38  
39 These issues reflect potential acceptance problems that were in the innovation  
40  
41 literature initially seen as general barriers for the uptake of user-contributed ideas by  
42  
43 companies. Since then, many studies have demonstrated that user-contributed ideas can often  
44  
45 compete with the ideas generated by company-internal staff (e.g., see (Kristensson et al.,  
46  
47 2004; Magnusson, 2009; Poetz & Schreier, 2012; Magnusson et al., 2014) and therefore  
48  
49 deserve to be given a fair chance. As a consequence, dedicated crowdsourcing systems have  
50  
51 gained widespread acceptance in the business community. Our study extends this question to  
52  
53 the *mode* of idea harvesting: can user-contributed ideas identified by an artificial intelligence  
54  
55 system reach sufficient recognition among company professionals? An online community  
56  
57  
58  
59  
60

1  
2  
3 related to craft brewing was used as the idea base for our study. Employees of Norwegian  
4  
5 craft brewery *Nøgne Ø* evaluated the automatically extracted ideas.  
6  
7  
8  
9  
10  
11  
12  
13  
14  
15  
16  
17  
18  
19  
20  
21  
22  
23  
24  
25  
26  
27  
28  
29  
30  
31  
32  
33  
34  
35  
36  
37  
38  
39  
40  
41  
42  
43  
44  
45  
46  
47  
48  
49  
50  
51  
52  
53  
54  
55  
56  
57  
58  
59  
60

For Peer Review

## Method

### *Machine learning for idea detection*

The machine learning system we employed is described in detail in Christensen et al. (2016) and Christensen et al. (Submitted manuscript). Although the technical properties of the system are not the central focus of the present paper, we will give a brief description of the system and how it was employed in our study. The machine learning system takes as input idea texts and non-idea texts that have been identified by human raters. The texts used for this study originate from *alt.beer.home-brewing*, a Usenet-based online community related to craft brewing. In this community people from all over the world discuss brewing-related issues. We expected ideas for product and process to be available in this community. At the time the texts were extracted, the community contained altogether 10582 posts. 3000 of these were extracted for the development of the training of the system. Those that contained ideas were identified by via crowdsourcing, using the *CrowdFlower* platform (a service similar to Amazon's *Mechanical Turk*). Five raters were assigned to each text and instructed to label the text as an idea text if it contained at least one idea.

Before the texts could be used for machine learning, several text pre-processing steps were performed. In this process the raw text content was turned into a row-column format, where each text was represented as a row and each term (i.e., each unique word or expression) as a column. In this process, all numbers, punctuation marks and stop words were removed. Uni-grams, bi-grams and tri-grams were generated. All terms that did not occur in at least 0.2% of the texts were omitted from the analysis. This process resulted in a dataset consisting of 10514 terms representing 10582 texts.

The 3000 training texts were separated from the remaining 10582 texts. From the 3000 training texts, we excluded all texts where not all five CrowdFlower raters had agreed on the



1  
2  
3 class membership. After excluding these, the new training set contained 1393 texts. 405 of the  
4  
5 texts were idea texts and 988 were non-idea texts. The training texts were partitioned at  
6  
7 random into three separate data sets: a training set (consisting of 70% of the texts), a  
8  
9 validation set (15% of the texts) and a hold-out or test set (15% of the texts). Such a partition  
10  
11 is essential for the tuning of the machine learning system (in the validation set) and for an  
12  
13 unbiased evaluation of its performance in the context of previously unseen data (hold-out set).  
14  
15 Based on the training set, validation set and hold-out, the automatic idea detection system was  
16  
17 trained and tested. The system was based on a linear support vector machine classifier (SVM;  
18  
19 for details, see Christensen et al., 2016). Performance statistics are reported in Table 1.  
20  
21  
22  
23

24 --- Table 1 ---

25  
26 From the remaining 7582 texts which had not been involved in the training, validation  
27  
28 and testing of the system in the study by Christensen et al. (Submitted manuscript), another  
29  
30 200 were extracted for the present study. Using the SVM classifier, the texts were scored as to  
31  
32 how likely they were to contain an idea. A histogram of the resulting posterior probabilities is  
33  
34 shown in Figure 1. These 200 texts were then used in the present study as the idea and non-  
35  
36 idea texts to be classified and rated by two brewing professionals.  
37  
38  
39

40 --- Figure 1 ---

41  
42  
43 *Measuring idea quality*

44  
45  
46  
47 The perceived quality of an idea can depend on the perspective of the person  
48  
49 evaluating the idea. This topic has received much attention in the creativity and innovation  
50  
51 management literature. In principle, idea quality could be measured on a “good idea” to “bad  
52  
53 idea” scale, but in most research it is decomposed into several attributes that represent  
54  
55 conceptually distinct dimensions of quality. Dean, Hender, Rodgers, & Santanen (2006)  
56  
57  
58  
59  
60

1  
2  
3 provide a comprehensive review of the idea quality literature published between 1990 and  
4  
5 2005. Based on the altogether 90 identified studies, they suggest that four dimensions of idea  
6  
7 quality can be distinguished: novelty, workability, relevance and specificity. An idea is novel  
8  
9 if it contains something that is new. An idea is workable if it is easy to implement and does  
10  
11 not violate known constraints. An idea is relevant if it satisfies pre-defined goals. An idea is  
12  
13 specific if it has been worked out in detail.  
14  
15

16  
17 Comparable sets of sub-dimensions have been suggested in the user innovation  
18  
19 literature. Kristensson, Gustafsson and Archer (2004) compared the ideation performance of  
20  
21 ordinary users, expert users and professionals. They used three quality attributes: originality  
22  
23 (comparable to the novelty dimension suggested by Dean et al., 2006), realisability  
24  
25 (comparable to the feasibility dimension) and value (comparable to the relevance dimension).  
26  
27 In a similar study, Magnusson (2009) compared the ideation performance of professionals,  
28  
29 technically skilled users, ordinary users, consulting users and creativity-trained ordinary users.  
30  
31 He used the quality attributes originality (comparable to novelty), producibility (comparable  
32  
33 to feasibility) and user-value (comparable to relevance). Using the same attributes,  
34  
35 Magnusson et al. (2014) compared technically skilled users with technically naïve users.  
36  
37 Poetz & Schreier (2012) compared the ideas of users and professionals in terms of the  
38  
39 attributes novelty, feasibility and customer benefit (comparable to value). Based on the four  
40  
41 studies that have a product user ideation focus, we chose novelty, feasibility and value as the  
42  
43 quality attributes for our study.  
44  
45  
46  
47

#### 48 49 *Procedure*

50  
51  
52 We established contact with Norwegian craft brewery *Nøgne Ø*. The brewery was  
53  
54 founded in 2002 by two Norwegian home brewers and is nowadays part of Norwegian  
55  
56 brewery group Hansa Borg Bryggerier. In 2015, *Nøgne Ø* produced 30 different styles of ales  
57  
58  
59  
60

1  
2  
3 and exported to more than 40 markets. Two company professionals were recruited as expert  
4  
5 raters. Expert 1 was 29 years old, female and had a business school background. Her  
6  
7 responsibilities at *Nøgne Ø* were sales and logistics. At the time the study was conducted, she  
8  
9 had been working for the brewery for 12 years. Expert 2 was 40 years old, male and had an  
10  
11 engineering background. His responsibilities at *Nøgne Ø* were related to marketing and the  
12  
13 web shop. At the time the study was conducted, he had been working for the brewery for 4.5  
14  
15 years.  
16

17  
18 The experts evaluated the 200 texts one-by-one and independently from each other.  
19  
20 First, the experts were instructed to read the respective text carefully. Then, they were asked  
21  
22 “Please evaluate if you think that the text contains one or more ideas” and to respond on a  
23  
24 binary “yes” versus “no” scale. If the expert had responded “yes”, three rating scales were  
25  
26 presented on which the expert was asked to evaluate the quality of the idea in terms of the  
27  
28 three attributes novelty, feasibility and value. The scales were horizontally aligned ranging  
29  
30 from very low (1) to very high (10). The instruction for the novelty attribute was: “Please  
31  
32 evaluate the novelty of the idea(s) in the text (by this we mean: to what degree does the idea  
33  
34 suggest something new)”. The instruction for the feasibility attribute was: “Please evaluate the  
35  
36 feasibility of the idea(s) in the text (by this we mean: to what degree is it possible to  
37  
38 implement the idea)”. The instruction for the value attribute was: “Please evaluate the value  
39  
40 of the idea(s) in the text (by this we mean: to what degree does the idea solve the underlying  
41  
42 problem)”.  
43  
44  
45  
46  
47

### 48 *Inter-rater reliability*

49  
50 To assess the inter-rater reliability of the idea/non-idea classification task, we calculated  
51  
52 Cohen’s kappa, normalised for differences between raters in their marginal distributions  
53  
54 (Cohen, 1960; Landis & Koch, 1977; von Eye & von Eye, 2008). The normalised version of  
55  
56 kappa takes on values between 0 and 1 where a value of 0 stands for chance-level agreement  
57  
58  
59  
60

1  
2  
3 and a value of 1 for the theoretical maximum of agreement, given the marginal distributions  
4 of the raters. Expert 1 identified 41 texts as containing ideas and 159 as not containing ideas.  
5  
6 Expert 2 identified 87 texts as containing ideas and 113 as not containing ideas. They agreed  
7  
8 on 35 texts as containing ideas and 107 as not containing ideas (see Table 2 for examples).  
9  
10 These counts correspond to a normalised kappa of 0.74, suggesting that there was substantial  
11  
12 agreement between the two experts as to whether a given text did or did not contain an idea.  
13  
14  
15

16  
17 --- Table 2 ---  
18

19  
20 To assess the inter-rater reliability of the idea quality rating task, we calculated  
21  
22 reliability measures based on generalisability theory (Cronbach, Gleser, Nanda, &  
23  
24 Rajaratnam, 1972; Brennan, 2001). Only the 69 texts which the machine learning classifier  
25  
26 had classified as an idea and which at least one of the brewery professionals had identified as  
27  
28 an idea were included in the analysis. The design was a two-facet crossed design with tasks  
29  
30 (the three quality attributes) and raters (the two brewery professionals) treated as fixed effects.  
31  
32 The reliability (generalisability coefficient) of the averaged rating of a randomly picked idea  
33  
34 text on the three attributes by the two raters was  $E\rho^2 = 0.71$ .  
35  
36  
37  
38  
39  
40  
41  
42  
43  
44  
45  
46  
47  
48  
49  
50  
51  
52  
53  
54  
55  
56  
57  
58  
59  
60

## Results

### *Presence of ideas*

Since our two company professionals had not perfectly agreed with each other on the presence or absence of ideas in the texts, we defined two validation criteria: a lenient criterion (Boolean OR: at least one professional had identified the respective text as containing an idea) and a strict criterion (Boolean AND: both professionals had identified the respective text as containing an idea).

Using the lenient criterion as a gold standard (where 47% of the 200 texts would be defined as true idea texts), the automatic idea detection system performed well. The classifier agreed with the company professionals in 77% of the cases as to whether a text did or did not contain an idea (accuracy). 75% of the texts which the classifier had identified as idea texts were also identified as idea texts by the company professionals (precision, also referred to as positive predictive value in the literature). The classifier correctly identified as idea texts 74% of the texts the professionals had identified as ideas (recall, also referred to as sensitivity or true positive rate in the literature). Since precision and recall always represent a trade-off, we also calculated their harmonic mean, the  $F_1$  measure, as a compromise. Using the lenient criterion, it reached a very respectable value of  $F_1 = 0.75$ . Classification accuracy statistics are reported in Table 3.

Using the strict criterion as a gold standard (where only 18% of the 200 texts would be defined as containing ideas), the automatic idea classification system still agreed with the company professionals in 67% of the cases as to whether a text did or did not contain an idea (accuracy). Due to the much stricter criterion as to what defined an idea text, the precision of the classifier was lower: only 33% of the texts which the classifier had identified as idea texts were also identified as idea texts by the company professionals. For the same reason, recall

1  
2  
3 was higher: the classifier correctly identified as idea texts 86% of the texts the professionals  
4 had identified as ideas. The  $F_1$  measure, as a compromise between precision and recall,  
5 reached a value of 0.47.  
6  
7  
8

9  
10 Taken together, the criterion validity of the automatic idea detection system can be  
11 regarded as satisfactory as long as it is used for the screening of potential ideas. Deployed in a  
12 company as a tool for filtering out candidate ideas for product and process innovations, it may  
13 significantly reduce the time and effort that would otherwise have to be spent by company  
14 staff on manual screening and preliminary evaluation of a number of user contributions in  
15 potentially relevant online fora.  
16  
17  
18  
19  
20  
21  
22

23  
24 --- Table 3 ---  
25

26  
27 *Quality of automatically detected ideas*  
28

29  
30 Figure 2 shows the distribution of the quality ratings of the ideas (i.e., those texts that  
31 had been identified as ideas by the automatic idea detection system and which had been also  
32 been identified as ideas by at least one of the two company professionals). For texts which  
33 both company professionals had classified as an idea, the values on the novelty, feasibility  
34 and value attributes are the averaged ratings of both company professionals. For texts which  
35 only one of the company professionals had identified as an idea, the values are the ratings  
36 given by that professional. The overall quality values were calculated as unweighted averages  
37 of the ratings on the novelty, feasibility and value attributes.  
38  
39  
40  
41  
42  
43  
44  
45  
46  
47  
48

49 --- Figure 2 ---  
50

51  
52 The distribution of the novelty ratings was concentrated in the lower range of the  
53 response scale (which had a minimum of 1 and a maximum of 10), the distribution of the  
54 feasibility ratings in the upper range of the response scale, and the distributions of the value  
55  
56  
57  
58  
59  
60

1  
2  
3 ratings and overall quality in the middle of the response scale. The results suggest that, on  
4  
5 average, the ideas which the automatic idea detection system extracted from the  
6  
7 *alt.beer.home-brewing* community appeared rather feasible to brewery professionals, were not  
8  
9 particularly novel, but had medium value and medium overall idea quality.  
10  
11  
12  
13  
14  
15  
16  
17  
18  
19  
20  
21  
22  
23  
24  
25  
26  
27  
28  
29  
30  
31  
32  
33  
34  
35  
36  
37  
38  
39  
40  
41  
42  
43  
44  
45  
46  
47  
48  
49  
50  
51  
52  
53  
54  
55  
56  
57  
58  
59  
60

For Peer Review

## Discussion and conclusion

The first aim of the present study was to investigate if ideas for product and process innovations detected by an artificial intelligence system (in this case, the one developed by Christensen et al., 2016) would also be regarded as ideas by company-internal staff who will be responsible for taking the ideas further in the innovation process. Our results suggest that this is to a considerable extent the case: the performance of the system can be regarded as sufficient for an initial screening of potential ideas. Deployed in a company as a tool for selecting candidate ideas for product and process innovations, it can significantly reduce the time and effort that would otherwise have to be spent by company staff on wading through a large number of user contributions in potentially relevant online communities.

The exact level of criterion-related validity that our system could achieve depended on several factors. The most important of these are (a) the definition of the “gold standard” against which the predictions are validated and (b) the cut-off used for transforming the continuous posterior probability score generated by the system into a binary prediction. In our analysis, we used two of the possible gold standards: a lenient criterion (at least one of the company professionals had rated the respective text as containing an idea) and a strict criterion (both company professionals had rated the text as containing an idea). The lenient criterion led to an implied base rate of 47% for the target event (i.e., the probability that a randomly chosen text from among the 200 used in the present study would contain an idea), whereas the strict criterion reduced the implied base rate to 18%. It is not possible to define on purely statistical grounds what the right base rate should be. This is complicated by the fact that the two company professionals who served as experts in our study did not have the same base rates in their individual classifications: Expert 1 appeared to use a more conservative



1  
2  
3 standard of judgment, rating 21% of the 200 texts as containing ideas, whilst Expert 2  
4  
5 appeared to use a more liberal standard, rating 44% of the texts as containing ideas.  
6  
7

8 Since the two experts also differed in terms of their functional responsibilities in the  
9  
10 company, it might not even be appropriate to look for perfect agreement—after all, a  
11  
12 company's ability to integrate different functional perspectives is one of the strongest  
13  
14 predictors of innovation success (e.g., see Evanschitzky, Eisend, Calantone, & Jiang, 2012).  
15  
16 Whether it makes more sense for a given company to use a stricter or more lenient criterion  
17  
18 for further filtering of the automatically identified ideas may depend more on strategy and  
19  
20 available resources: a lenient criterion may be more appropriate if a company wants to cast its  
21  
22 net wide and thereby reduce the risk of missing certain ideas which might not yet be able to  
23  
24 achieve full cross-functional consensus. However, the company would also have to be  
25  
26 prepared to assign the necessary resources for dealing with the larger number of ideas that  
27  
28 would enter the innovation funnel. If, on the other hand, a company wants to limit its resource  
29  
30 expenditure and focus on ideas that can already in the early phases achieve cross-functional  
31  
32 consensus, a stricter criterion would be appropriate.  
33  
34  
35  
36  
37

38 A similar objective can be achieved by tuning the cut-off value of the SVM classifier  
39  
40 underlying the Christensen et al. (2016) system. The algorithm yields a posterior probability  
41  
42 score that is continuous on the (0,1) interval. A traditional way of transforming the posterior  
43  
44 probability score into a binary classification is use the value 0.50 as a cut-off such that a text  
45  
46 is classified as an idea text if the probability that the text contains an idea, given the support  
47  
48 vectors, is larger than 0.50, and classified as a non-idea text otherwise. However, the  
49  
50 traditional way of setting the cut-off value may not always be the most useful way. Another  
51  
52 heuristic that is typically more useful is to set the cut-off equal to one minus the base rate of  
53  
54 the target even, either on the posterior probability scale or on the empirical percentile scale.  
55  
56 This heuristic would match the prior probability of classifying a text as an idea to the base rate  
57  
58  
59  
60

1  
2  
3 of the event. A third way of setting the cut-off is to estimate how many additional ideas a  
4  
5 company would be able to absorb into its innovation funnel and to use an appropriate absolute  
6  
7 cut-off, selecting the right number of ideas from the top of the posterior probability ranking.  
8  
9

10  
11 The second aim of the present study was to investigate if the automatic idea detection  
12  
13 system developed by Christensen et al. (2016) would extract *good* ideas from the online  
14  
15 community that served as an example here. For the online community under investigation, our  
16  
17 answer is a qualified yes: the distribution of the overall idea quality score, calculated as the  
18  
19 average rating of each idea on the three quality attributes (novelty, feasibility, value) by the  
20  
21 two company professionals, was concentrated in the middle of the response scale (mean = 4.8,  
22  
23 25<sup>th</sup> percentile = 3.8, 50<sup>th</sup> percentile = 5, 75<sup>th</sup> percentile = 5.7) and ranged from a minimum of  
24  
25 1 (the lower end of the response scale) to a maximum of 8 (two points below the maximum of  
26  
27 the response scale). Overall, the ideas extracted by the automatic detection system appear to  
28  
29 have made a reasonable impression on the company professionals.  
30  
31  
32

33  
34 An interesting detail is that the identified ideas tended to be regarded as more feasible  
35  
36 and valuable by our company professionals than they were regarded as novel. This finding  
37  
38 reflects results obtained by Kristensson et al. (2004) and Magnusson (2009) for user ideation  
39  
40 in general. However, as already observed, agreement between our experts was not perfect  
41  
42 here either. As an example, consider the text shown in Table 4: a community member  
43  
44 suggests a new mead recipe. Overall, the idea was rated as one of the best by the two  
45  
46 company professionals. Expert 1 assigned a rating of 2 on the novelty attribute, 7 on  
47  
48 feasibility and 4 on value. Expert 2 rated it 9 on novelty, 9 on feasibility and 9 on value. In the  
49  
50 additional, qualitative responses we obtained from the two professionals, it became clear that  
51  
52 Expert 1 evaluated the idea in terms of its quality as an idea for process innovation whereas  
53  
54 Expert 2 evaluated it in terms of its quality as idea for product innovation. Different  
55  
56 perspectives, either due to the functional specialisation of our company professionals or due  
57  
58  
59  
60

1  
2  
3 to their different levels of experience with the product category, seem to have led to different  
4  
5 standards of judgment.  
6

7  
8 --- Table 4 ---  
9

10 The results presented here are an evaluation of a particular automatic idea detection  
11 system (the one developed by Christensen et al., 2016) to a particular case (the craft brewing  
12 community *alt.beer.home-brewing*), evaluated from the point of view of two brewing  
13 professionals connected to a particular craft brewing company (*Nøgne Ø*). Naturally, this  
14 poses limits to the generalisability of our findings. The ideas detected by an automated system  
15 can only be as good as the ideas voiced by the users in the online community under  
16 investigation. Furthermore, the 200 texts we selected for evaluation were only a sample and  
17 therefore unlikely to reflect the whole range of ideas discussed in the community. It is an open  
18 question whether similar results will be achieved when automatic idea detection systems are  
19 applied to other technology domains or product categories.  
20  
21  
22  
23  
24  
25  
26  
27  
28  
29  
30  
31

32 This question can only be answered by follow-up research. However, we do believe  
33 that we have demonstrated the potential of automatic idea identification systems: they can be  
34 a powerful technique for the harvesting and initial screening of user ideas from online fora  
35 that do not conform, and are not limited to, the highly restrictive architecture and user basis of  
36 dedicated crowdsourcing systems. We hope that studies such as ours can also make a  
37 contribution to a wider discussion: which business tasks of a more complex nature can  
38 credibly be solved by artificial intelligence-based systems? We are convinced that the answer  
39 does not only lie in what is technically possible but also in what is acceptable to the  
40 prospective users of the information generated by such systems. More user evaluations of the  
41 performance of artificial intelligence-based systems are needed.  
42  
43  
44  
45  
46  
47  
48  
49  
50  
51  
52  
53  
54  
55  
56  
57  
58  
59  
60

## References

- 1  
2  
3  
4  
5  
6 Antorini, Y. M. (2007). *Brand Community Innovation: An Intrinsic Case Study of the Adult*  
7  
8 *Fans of LEGO Community*. Copenhagen Business School, Frederiksberg: Center for  
9  
10 Europaforskning.
- 11  
12 Antorini, Y. M., Muñiz, J., Albert M., & Askildsen, T. (2012). Collaborating With Customer  
13  
14 Communities: Lessons from the Lego Group. *MIT Sloan Management Review*, 53(3),  
15  
16 73–95.
- 17  
18  
19 Brennan, R. L. (2001). *Generalizability theory*. New York: Springer-Verlag.
- 20  
21 Christensen, K., Liland, K. H., Kvall, K., Risvik, E., Biancolillo, A., Scholderer, J., ... Næs,  
22  
23 T. (Submitted manuscript). Mining online community data: The nature of ideas in  
24  
25 online communities. *Food Quality and Preference*.
- 26  
27  
28 Christensen, K., Nørskov, S., Frederiksen, L., & Scholderer, J. (2016). In search of new  
29  
30 product ideas: Identifying ideas in online communities by machine and text mining.  
31  
32 *Creativity and Innovation Management (Available Online)*.
- 33  
34 Cohen, J. (1960). A coefficient of agreement for nominal scales. *Educational and*  
35  
36 *Psychological Measurement*, 20(1), 37–46.
- 37  
38  
39 Cronbach, L. J., Gleser, G. C., Nanda, H., & Rajaratnam, N. (1972). *The dependability of*  
40  
41 *behavioral measurements*. New York: Wiley.
- 42  
43 Dean, D. L., Hender, J. M., Rodgers, T. L., & Santanen, E. L. (2006). Identifying quality,  
44  
45 novel, and creative Ideas: Constructs and scales for idea evaluation. *Journal of the*  
46  
47 *Association for Information Systems*, 7(1), 646–698.
- 48  
49  
50 di Gangi, P. M., Wasko, M. M., & Hooker, R. E. (2010). Getting customers' ideas to work for  
51  
52 you: Learning from Dell how to succeed with online user innovation communities.  
53  
54 *MIS Quarterly Executive*, 9(4), 213–228.
- 55  
56  
57  
58  
59  
60

- 1  
2  
3 Ekvall, G. (1997). Organizational conditions and levels of creativity. *Creativity and*  
4  
5 *Innovation Management*, 6(4), 11.  
6  
7 Evanschitzky, H., Eisend, M., Calantone, R. J., & Jiang, Y. (2012). Success factors of product  
8  
9 innovation: An updated meta-analysis. *Journal of Product Innovation Management*,  
10  
11 29(S1), 21–37.  
12  
13 Füller, J., Bartl, M., Ernst, H., & Mühlbacher, H. (2006). Community based innovation: How  
14  
15 to integrate members of virtual communities into new product development.  
16  
17 *Electronic Commerce Research*, 6(1), 57–73.  
18  
19 Füller, J., Jawecki, G., & Mühlbacher, H. (2007). Innovation creation by online basketball  
20  
21 communities. *Journal of Business Research*, 60(1), 60–71.  
22  
23 Jeppesen, L. B., & Frederiksen, L. (2006). Why do users contribute to firm-hosted user  
24  
25 communities? The case of computer-controlled music instruments. *Organization*  
26  
27 *Science*, 17(1), 45–63.  
28  
29 Kristensson, P., Gustafsson, A., & Archer, T. (2004). Harnessing the creative potential among  
30  
31 users\*. *Journal of Product Innovation Management*, 21(1), 4–14.  
32  
33 Landis, J. R., & Koch, G. G. (1977). The Measurement of Observer Agreement for  
34  
35 Categorical Data. *Biometrics*, 33(1), 159.  
36  
37 Lin, F.-R., Hsieh, L.-S., & Chuang, F.-T. (2009). Discovering genres of online discussion  
38  
39 threads via text mining. *Computers & Education*, 52(2), 481–495.  
40  
41 Magnusson, P. R. (2009). Exploring the Contributions of Involving Ordinary Users in  
42  
43 Ideation of Technology-Based Services\*. *Journal of Product Innovation Management*,  
44  
45 26(5), 578–593.  
46  
47 Magnusson, P. R., Wästlund, E., & Netz, J. (2014). Exploring Users' Appropriateness as a  
48  
49 Proxy for Experts When Screening New Product/Service Ideas: Exploring Users as a  
50  
51 Proxy for Expert Judges. *Journal of Product Innovation Management*, 33(1), 4–18.  
52  
53  
54  
55  
56  
57  
58  
59  
60

- 1  
2  
3 Mahr, D., & Lievens, A. (2012). Virtual lead user communities: Drivers of knowledge  
4  
5 creation for innovation. *Research Policy*, 41(1), 167–177.  
6  
7 Majchrzak, A., & Malhotra, A. (2013). Towards an information systems perspective and  
8  
9 research agenda on crowdsourcing for innovation. *The Journal of Strategic*  
10  
11 *Information Systems*, 22(4), 257–268.  
12  
13 Nørskov, S., Antorini, Y. M., & Jensen, M. B. (2015). Innovative brand community members  
14  
15 and their willingness to share ideas with companies. *International Journal of*  
16  
17 *Innovation Management*.  
18  
19 Poetz, M. K., & Schreier, M. (2012). The value of crowdsourcing: Can users really compete  
20  
21 with professionals in generating new product ideas? *Journal of Product Innovation*  
22  
23 *Management*, 29(2), 245–256.  
24  
25  
26  
27 Thorleuchter, D., & Van den Poel, D. (2013). Web mining based extraction of problem  
28  
29 solution ideas. *Expert Systems with Applications*, 40(10), 3961–3969.  
30  
31  
32 Van de Ven, A. (1986). Central problems in the management of innovation. *Management*  
33  
34 *Science*, 32(5), 590–607.  
35  
36 van den Ende, J., Frederiksen, L., & Prencipe, A. (2015). The Front End of Innovation:  
37  
38 Organizing Search for Ideas. *Journal of Product Innovation Management*, 32(4), 482–  
39  
40 487.  
41  
42  
43 Vandenbosch, B., Saatcioglu, A., & Fay, S. (2006). Idea management: a systemic view.  
44  
45 *Journal of Management Studies*, 43(2), 259–288.  
46  
47 von Eye, A., & von Eye, M. (2008). On the marginal dependency of Cohen's  $\kappa$ . *European*  
48  
49 *Psychologist*, 13(4), 305–315.  
50  
51  
52 von Hippel, E., Ogawa, S., & PJ de Jong, J. (2011). The age of the consumer-innovator.  
53  
54  
55 von Krogh, G., Spaeth, S., & Lakhani, K. R. (2003). Community, joining, and specialization  
56  
57 in open source software innovation: a case study. *Research Policy*, 32(7), 1217–1241.  
58  
59  
60

1  
2  
3 von Krogh, G., & von Hippel, E. (2006). The Promise of Research on Open Source Software.  
4

5 *Management Science*, 52(7), 975–983.  
6  
7  
8  
9  
10  
11  
12  
13  
14  
15  
16  
17  
18  
19  
20  
21  
22  
23  
24  
25  
26  
27  
28  
29  
30  
31  
32  
33  
34  
35  
36  
37  
38  
39  
40  
41  
42  
43  
44  
45  
46  
47  
48  
49  
50  
51  
52  
53  
54  
55  
56  
57  
58  
59  
60

For Peer Review

1  
2  
3  
4  
5  
6  
7  
8  
9  
10  
11  
12  
13  
14  
15  
16  
17  
18  
19  
20  
21  
22  
23  
24  
25  
26  
27  
28  
29  
30  
31  
32  
33  
34  
35  
36  
37  
38  
39  
40  
41  
42  
43  
44  
45  
46  
47  
48  
49  
50  
51  
52  
53  
54  
55  
56  
57  
58  
59  
60

Table 1 - Performance of the automatic idea detection system used by Christensen et al. (Submitted Manuscript)

Partition	True positives (TP)	True negatives (TN)	False positives (FP)	False negatives (FN)	Classification accuracy	Precision	Recall	$F_1$
Validation set	27%	70%	1%	2%	0.97	0.97	0.92	0.94
Hold-out set	25%	70%	1%	3%	0.96	0.96	0.88	0.92

For Peer Review



Table 2 - Example of an idea text and a non-idea text on which both raters agreed

Idea text	Non-idea text
'Buckwheat has been used as an adjunct for a long time in a few beers. It also is used to make gluten free beers. It has a high gelatinisation temp so need to be boiled first. Extract potential is about 1.032. Can be used lightly roasted to add colour to gluten free beers, or use Kasha (a roasted buchwheat). I think Rogues make a buckwheat ale'	'Thanks for the help. My internet is screwy or I would have replied sooner. I re- pitched and it is going crazy. a load off my mind! now i can concentrate on getting another cider and a wit going. Anyone have any suggestions for a good belgian style ale I ike duvel? I am an extract with specialty grains level brewer, so whole grain is out for now. Thanks again for all the help!'

For Peer Review

Table 3 - Presence of ideas: classification accuracy of the automatic idea detection system, validated against the judgments of two company professionals

Validation criterion	True positives (TP)	True negatives (TN)	False positives (FP)	False negatives (FN)	Classification accuracy	Precision	Recall	$F_1$
Lenient criterion:								
Classified as idea by Expert 1 OR Expert 2	35%	42%	12%	12%	0.77	0.75	0.74	0.75
Strict criterion:								
Classified as idea by Expert 1 AND Expert 2	15%	52%	31%	3%	0.67	0.33	0.86	0.47

Table 4 - Idea text identified by classifier, Expert 1 and Expert 2

I've made several batches. Below is my recipe The love of my life I love Mead as you can probably tell. Please note, this is Mead but I do not use any water. I use apple juice as the base. You can use water but I find the apple juice makes it a bit nicer for those of you who love apples and like a high alcohol content. No citric acid needed. This is called Apple- Honey Melonomel Meade You will need... 1 Package Red Star wine yeast 4 Gallons apple juice from concentrate 2-5 pounds of pure honey, the more the better. This shit is expensive though. 1 cup table sugar 5 Fuji apples Siphon hose, any small tube will work. A 5 gallon carboy or tub 1 balloon Step one, crush your apples or use a blender. Step two, boil apples in large pot with apple juice. Step three, set aside to cool Step Four, boil honey in large pot of apple juice Step five, set aside to cool. Step six, dump mixture into large 5 gallon carboy and add activated yeast. Step six, allow the mead to ferment for 3-4 weeks, once fermentation begins to slow prime with table sugar by diluting the 1 cup of table sugar in 1/2 gallon of apple juice then pour this directly into the carboy. A balloon can be placed over the mouth of the carboy to monitor the fermentation. Simply pierce a small hole in the balloon to allow CO2 to escape. Once the Meade has cleared (meaning you can read a newspaper through it) transfer it into a secondary (Save the sediment for use as the Yeast in your next batch of Meade) and let it clarify for 2-3 weeks. After this bottle the meade and let fermentation finish off. Total process about 70 days and its ready to drink. This will burn going down but is smooth as a whistle. Enjoy....'

1  
2  
3  
4  
5  
6  
7  
8  
9  
10  
11  
12  
13  
14  
15  
16  
17  
18  
19  
20  
21  
22  
23  
24  
25  
26  
27  
28  
29  
30  
31  
32  
33  
34  
35  
36  
37  
38  
39  
40  
41  
42  
43  
44  
45  
46  
47  
48  
49  
50  
51  
52  
53  
54  
55  
56  
57  
58  
59  
60

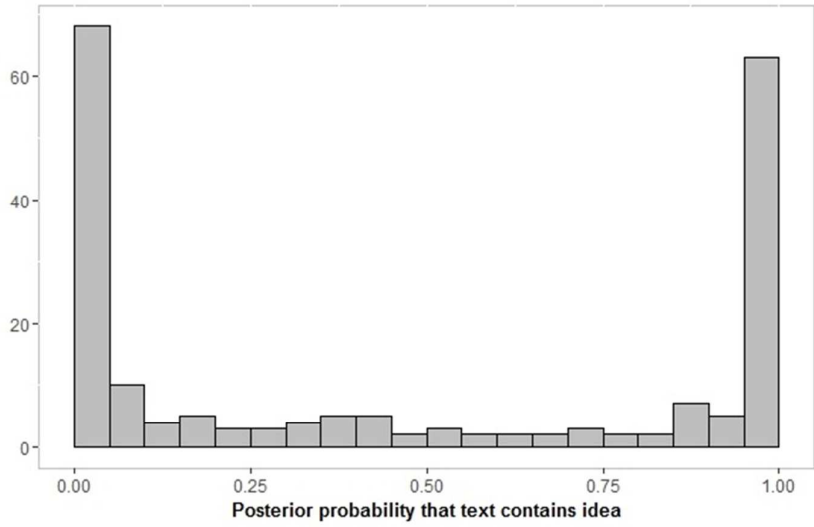


Figure 1 - Histogram of the posterior probability scores generated by the SVM-based automatic idea detection system for the 200 texts used in the present study

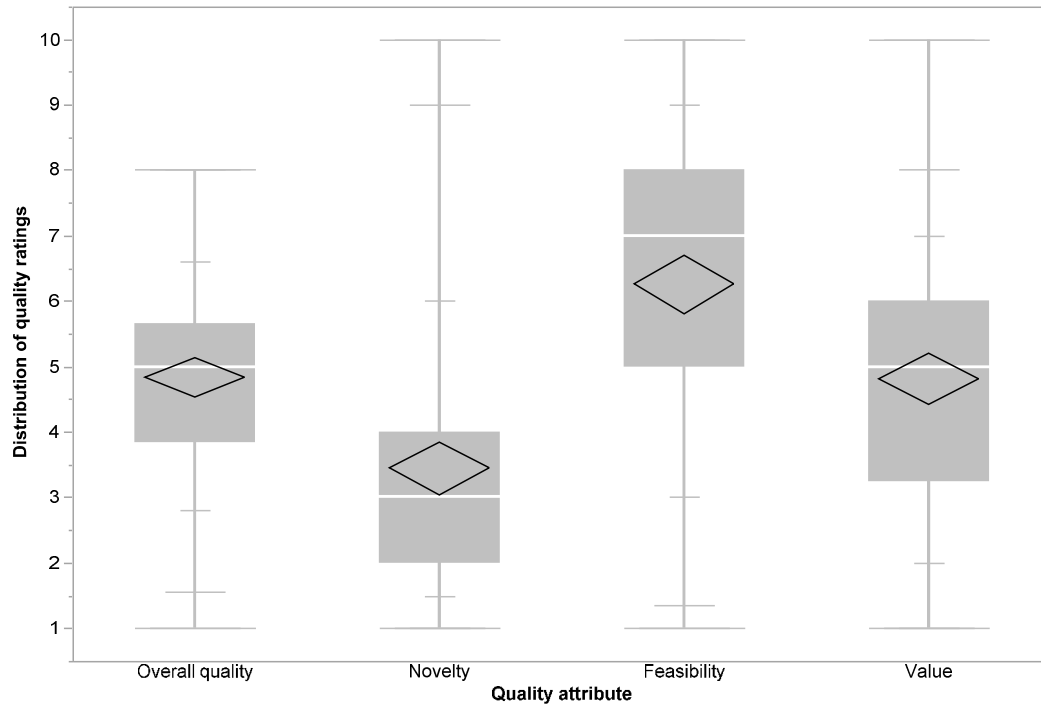


Figure 2 - Box plots of the distribution of quality ratings (overall quality = unweighted average of novelty, feasibility and value; diamonds represent 95% confidence intervals around distribution means)