



**HAL**  
open science

## Belief Temporal Analysis of Expert Users: case study Stack Overflow

Dorra Attiaoui, Arnaud Martin, Boutheina Ben Yaghlane

► **To cite this version:**

Dorra Attiaoui, Arnaud Martin, Boutheina Ben Yaghlane. Belief Temporal Analysis of Expert Users: case study Stack Overflow. Big Data Analytics and Knowledge Discovery DAWAK, Aug 2017, Lyon, France. hal-01576875

**HAL Id: hal-01576875**

**<https://hal.science/hal-01576875v1>**

Submitted on 24 Aug 2017

**HAL** is a multi-disciplinary open access archive for the deposit and dissemination of scientific research documents, whether they are published or not. The documents may come from teaching and research institutions in France or abroad, or from public or private research centers.

L'archive ouverte pluridisciplinaire **HAL**, est destinée au dépôt et à la diffusion de documents scientifiques de niveau recherche, publiés ou non, émanant des établissements d'enseignement et de recherche français ou étrangers, des laboratoires publics ou privés.

# Belief Temporal Analysis of Expert Users: case study Stack Overflow

Dorra Attiaoui<sup>1</sup>, Arnaud Martin<sup>2</sup>, and Boutheina Ben Yaghlane<sup>3</sup>

<sup>1</sup> DRIUD IRISA, University of Rennes 1, LARODEC, University of Tunis  
dorra.attiaoui@irisa.fr

<sup>2</sup> DRIUD IRISA, University of Rennes 1  
Arnaud.Martin@univ-rennes1.fr

<sup>3</sup> LARODEC, University of Carthage  
boutheina.yaghlane@ihec.rnu.tn

**Abstract.** Question Answering communities have known a large expansion over the last few years. Reliable people sharing their knowledge are not that numerous. Thus, detecting experts since their first contributions can be considered as a challenge. We are interested in studying the activity of these platforms' users during a defined period of time. As the data collected is not always reliable, imperfections can occur. In order to manage these imperfections, we choose to use the mathematical background offered by the theory of belief functions. People say that the more time they spend within a community, the more knowledge they acquire. We investigate this assumption in this paper by studying the behavior of users without taking into consideration the reputation system proposed by Stack Overflow. Experiments with real data from Stack Overflow demonstrate that this model can be applied to any expertise detection problem. Moreover, it allows to identify potential future experts. The analysis allows us to study the behavior of experts and non expert users over time spent in the community. We can see that some users keep on being reliable while others do gain knowledge and improve their expertise measure.

**Keywords:** Question answering community, theory of belief functions, expertise measure, classification

## 1 Introduction

With the emergence of Question Answering Communities (Q&A C), several platforms were developed aiming to help people. The main challenge of these websites is to provide helpful, quick and well organized answers for any posted question regarding any specific topic.

One of the most popular platforms is Stack Overflow (SO)<sup>4</sup>. It is the largest online community for programmers. Here, users can post questions, answers them, vote positively or negatively for both answers and questions in order to express their opinion on the quality of the posts.

---

<sup>4</sup> <http://stackoverflow.com>

Stack Overflow proposes a reputation system to reward active users. Actually, reputation<sup>5</sup> is the summery of users' activity in the web site. It is earned by convincing other users that he/she knows what he/she is talking about. Indeed, reputation reflects how involved a user is in the community and how other people see him/her. If this value is high, it means that a user is able to post fair questions or/and answers and how well he/she can communicate and interact with his/her peers. It also means that we can be in presence of a knowledgeable person. However, we assume this measurement as flawed. Reputation may support competitive gain of points rather than fair contributions in the community. Note that in this paper, we do not consider the reputation because it is a rough measurement of expertise according to only other people's opinion, and not founded on both their activities and opinions [19].

Detecting experts in online communities have been wildly investigated. We can distinguish two different methods: ranking based approaches and attribute based approaches. [15]. On one hand, the ranking based approaches intent to measure a score per user then select the top users [18, 19]. On the other hand, attributes based approaches aim to identify a number of features for the users and later apply machine learning techniques in order to classify them as experts and non-expert users [1, 12].

However, the literature suffers from few limitations like : 1) the dependence between training data and labels results on supervised machine learning, 2) high time consuming processes, and 3) the proposed approaches consider all the manipulated data as certain and perfect. Thus, this can not be taken into consideration, especially when we are dealing with real world applications. Several theories were proposed to manage uncertainty such as probability theory [14], possibility theory [6] and the theory of belief functions [2, 16]. The latter can be presented as a generalization of the other theories. Besides it offers a rich tool able to manage different types of data imperfections. When manipulating uncertainty, information fusion can be an interesting solution to obtain relevant information. Data fusion based on the theory of belief functions has been widely used in classification, image processing [8], clustering [4], etc. and more recently in social networks [10].

In this paper, we propose a new approach of measuring expertise and analyzing the behavior of experts and potential experts during a period of time based on uncertainty theories. The remainder of this paper is organized as follows. In Section 2, an overview on experts detection in social networks and in Stack Overflow. In Section 3, we present the basic necessary background related to the theory of belief functions. Section 4 details the approach for the representation of the proposed expertise measure and experts detection. Then in section 5, we present the results from experiments on Stack Overflow's data.

## 2 Related work

Most of the users of Stack Overflow aim to win as much reputation points as possible in order to obtain privileges like creating tags, moderating the forum etc. The reputation is defined according to the system presented in Table 1.

Every posted question or answer can be submitted either to positive or negative votes. A positive vote is a reward for the author, while the negative one penalizes him. Each

<sup>5</sup> <http://stackoverflow.com/help/whats-reputation>

**Table 1.** Gratification system of Stack Overflow

<b>Action</b>	<b>Reputation</b>
Answer voted up	+10
Question voted up	+5
Accepted answer	+15 (+2 to question asker )
Question voted down	-2
Answer voted down	-2 (-1 to voter)
Spammed answers	-100

person who posts a question is allowed to choose the best answer that seems to be the most helpful allowing his/her owner to gain reputation points.

Several researches focused on detecting experts in Stack Overflow. For [12], experts are known to provide the best answers in a very short time. They are more reactive and their answers are more useful than usual users. For [9], the authors proposed an analysis of Stack Overflow's reputation system. They focused on the contributors participation model. They consider the reputation as measurement of expertise. Any user with a reputation greater than 2400 points is an expert. However, their approach seems to be strict because it is only based on the value of the reputation gathered during users' activity in the platform. Another approach is proposed in [19] that is not founded on the reputation measure. They propose a metric called "Mean Expertise Contribution" that takes into account two indices: the debate generated by a question and the utility of the provided answers. The first index is related to the number of answers proposed for a given question. The second one is calculated according to the rating of an answer among all the answers provided.

Some other researches were interested in identifying experts and potential future experts in Stack Overflow using temporal analysis. For [11], authors modeled users' behavior based on their early participation in the community and showed that they could use classification as well as ranking algorithms to identify potential experts. They proposed that experts can be effectively identified from their early behavior. In [13], authors considered that expertise is present from the beginning and does not increase with the time spent in the community. Recently, [5] defined early expertise based on the number of best answers given by a user. Besides, they proposed an approach based on the combination of large number of textual, behavioral and time-aware features for detecting early expertise.

In [7], the authors identified three levels of uncertainty in question answering communities. The first level is related to the extraction and integration of the data. The second one deals with information sources, meaning the users of these platforms. The third level covers the uncertainty of the information itself. In the considered case, we are more interested in the evaluation of the sources and the part of uncertainty related to them. The main issue in these communities is that we are dealing with users that we do not usually have an *a priori* knowledge about them. We ignore everything about the sources' reliability, or expertise. In order to deal with this uncertainty, we will use the mathematical background provided by the theory of belief functions. This will help us to

consider the problem of early identification of potential experts with an uncertain point of view.

### 3 Theory of belief functions: an overview

This section recalls the necessary background notions related to the theory of belief functions. This theory has been developed by Dempster in his work on upper and lower probabilities [2]. Afterwards, it was formalized in a mathematical framework by Shafer in [16]. This theory is able to deal and represent imperfect (uncertain, imprecise and /or incomplete) information.

Let us consider a variable taking values in a finite set  $\Omega = \{\omega_1, \dots, \omega_n\}$  called *the frame of discernment*.

A *basic belief assignment (bba)* is defined on the set of all subsets of  $\Omega$ , named power set and noted  $2^\Omega$ . It affects a real value from  $[0, 1]$  to every subset of  $2^\Omega$  reflecting sources amount of belief on this subset. A bba  $m$  verifies:

$$\sum_{X \subseteq \Omega} m(X) = 1. \quad (1)$$

#### 3.1 Particular belief functions

Mass function is the common representation of evidential knowledge. Basic belief masses are degrees of support justified by available evidences. This section recalls some particular mass functions.

**Categorical mass functions** A categorical mass function is a normalized mass function which has a unique focal element  $X^*$ . This mass function is noted  $m(X)$  and defined as follows:

$$m_{X^*}(X) = \begin{cases} 1 & \text{if } X = X^* \subseteq \Omega \\ 0 & \forall X \subseteq \Omega \text{ and } X \neq X^* \end{cases} \quad (2)$$

We distinguish two particular cases of categorical mass functions: the vacuous mass functions when  $X^* = \Omega$  and the contradictory mass functions if  $X^* = \emptyset$ .

**Vacuous mass functions** A vacuous mass function is a particular categorical mass function focused on  $\Omega$ . It means that a vacuous mass function is normalized and has a unique focal element which is  $\Omega$ . This type of mass functions is defined as follows:

$$m_\Omega(X) = \begin{cases} 1 & \text{if } X = \Omega \\ 0 & \text{otherwise} \end{cases} \quad (3)$$

Vacuous mass function emphasizes the case of total ignorance.

**Simple support mass functions** Simple support mass functions are a special type that allow us to model both of the uncertainty and imprecision according the following equation:

$$\begin{cases} m(X) = 1 - \omega, X \subset \Omega \\ m(\Omega) = \omega \\ m(Y) = 0, Y \neq X \subset \Omega \end{cases} \quad (4)$$

where the mass on  $m(\Omega)$  represents the ignorance.

In the theory of belief function, Dempster in [2] proposed the first combination rule. It is defined for two bbas  $m_1, m_2, \forall X \in 2^\Omega$  with  $X \neq \emptyset$  by:

$$m_{DS}(X) = \frac{1}{1-k} \sum_{A \cap B = X} m_1(A)m_2(B), \quad (5)$$

where  $k$  is generally called the inconstancy of the combination, defined by  $k = \sum_{A \cap B = \emptyset} m_1(A)m_2(B)$  and  $1 - k$  is a normalization constant.

### 3.2 Discounting

Sometimes, it is possible to quantify the reliability of the body of evidence assessing degrees of support. The reliability of information sources reflects both its degrees of expertise and trust. When handling a mass function, we have to take into account the degree of reliability of its source. the degree of reliability of a source is taken into account by integrating it into all its mass functions. Using discounting operation in belief functions was first introduced in [16].

Discounting a mass function  $m$  consists in weighting every mass  $m(X)$  by a coefficient  $\alpha \in [0, 1]$  called reliability;  $\alpha$  is the discount rate. The *bba* is discounted as follows:

$$\begin{cases} \alpha m(X) = \alpha m(X) & \forall, X \subset 2^\Omega \setminus \Omega \\ \alpha m(\Omega) = 1 - (\alpha(1 - m(\Omega))) \end{cases} \quad (6)$$

### 3.3 Decision making

In order to make decision within the mathematical background of the theory of belief functions, [17] proposed to transform mass functions into probabilities (called BetP) using the pignistic probability transformation.

To do so, it transforms a *bba*  $m$  into a probability measure for all  $X \in 2^\Omega$  :

$$BetP(X) = \sum_{Y \neq \emptyset} \frac{|X \cap Y|}{|Y|} \frac{m(Y)}{1 - m(\emptyset)}. \quad (7)$$

where  $|Y|$  is the cardinality of  $Y$ .

## 4 Belief model of users in Stack Overflow

In this section, we detail the proposed approach that follows three main phases. In the first, we define the hypothesis describing each category of user and how do they behave in online communities. Then, parameters are estimated and mass functions are constructed. Finally, these latter are integrated in the general combination process. We explain in more details these two big steps in what follows.

First, we present some features related to very user  $i$ :

- **Number of votes related to answers** ( $AV_i$ ): the sum of positive votes collected by posted questions and answers.
- **Number of votes related to questions** ( $QV_i$ ): the sum of negative votes collected by posted questions and answers.
- **Time activity**: time of the activity of users from their registration to their last connection.
- **Number of posted questions** ( $NbQu_i$ ): number of questions posted in the dataset during the time activity of a user.
- **Number of posted answers** ( $NbAn_i$ ): number of answers provided in the dataset during the time activity of a user.
- **Number of posted answers** ( $NbAccAn_i$ ): number the answers chosen as the best.

### 4.1 Hypothesis

We applied an ascending hierarchical classification on the dataset allowing us to distinguish between three types of users in online communities. Here, we present these classes and the hypothesis proposed in order to identify each one of them.

- **Occasionals (O)**: these users represent the major part of members on the platform. They do not have a lot of knowledge. They occur occasionally only when they need an answer to a specific question that have not been treated before.
- **Apprentices (A)**: these users may have some expertise in a given topic. They aim to increase their reputation. To do so, they post a lot of answers that are not always very useful. The quality of their posts is not guaranteed and their answers can be down-voted.
- **Experts (E)**: these users are very reliable and recognized by the community. They provide a considerable number of useful answers that are chosen as the best ones. They are very active in the platform and guarantee a high quality content.

According to the previous presentation of the classes of users, we can define the following hypothesis:

**Hypothesis 1** *If a user has a high score of answers this might mean that this person is an expert rewarded for the answers provided.*

**Hypothesis 2** *If a user has a high score of questions score this might mean that this person is an apprentice seeking for information, and rewarded for posting well asked and interesting questions*

**Hypothesis 3** *If a user has a high number of answers posted this can be justified by two facts. First, this person is an expert, providing high quality content. Second, it can be an apprentice trying to become an expert by proving to the community that he/she can be as reliable as an expert.*

**Hypothesis 4** *If a user has a high number of questions posted this can represent either an expert or an apprentice. Both of them ask a lot of questions.*

**Hypothesis 5** *If a user has a high number of accepted answers this can only represent experts. Experts are frequently chosen as the most helpful answers providers.*

## 4.2 Definition of mass functions

In this section we detail the mathematical model that defines the hypothesis presented bellow. For each hypothesis, we determine how to define the mass functions in order to represent the data relative to each user.

Each user  $u$  is characterized by the following features:

- According to the hypothesis 1, a high score on answers is represented by a mass function on the focal element "**Expert**" (E) and the remainder is given to the ignorance, for a user  $i$ :

$$\begin{aligned} m_1^i(E) &= \alpha_1(1 - e^{-\gamma_1 AV_i}) \\ m_1^i(\Omega) &= \alpha_1 e^{-\gamma_1 AV_i} \end{aligned} \quad (8)$$

- According to the hypothesis 2, a high score on questions is represented by a mass function on the focal element "**Apprentice**" (A) and the remainder is given to the ignorance, for a user  $i$ :

$$\begin{aligned} m_2^i(A) &= \alpha_2(1 - e^{-\gamma_2 QV_i}) \\ m_2^i(\Omega) &= \alpha_2 e^{-\gamma_2 QV_i} \end{aligned} \quad (9)$$

- According to the hypothesis 3, a high number of posted questions is represented by a mass on the union of two classes "**Apprentice**  $\cup$  **Expert**". Otherwise, when this value is low it is affected to the "**Occasional**" (O) and the reminder to the ignorance. When a mass is on the union, this means that we can not decide which one of these classes is concerned by the mass. For a user  $i$ :

$$\begin{aligned} m_3^i(E \cup A) &= \alpha_3 \left(1 - e^{-\gamma_3 NbQu_i}\right) \\ m_3^i(O) &= \alpha_3 e^{-\gamma_3 NbQu_i} \\ m_3^i(\Omega) &= 1 - \alpha_3 \end{aligned} \quad (10)$$



- According to the hypothesis 4, a high number of answers is represented by a mass on the union of "**Apprentice**  $\cup$  **Expert**" while on the opposite situation the mass is transferred to the "**Occasional**" and the reminder to the ignorance. For a user  $i$ :

$$\begin{aligned} m_4^i(E \cup A) &= \alpha_4(1 - e^{-\gamma_4 NbAn_i}) \\ m_4^i(O) &= \alpha_4 e^{-\gamma_4 NbAn_i} \\ m_4^i(\Omega) &= 1 - \alpha_4 \end{aligned} \quad (11)$$

- According to the hypothesis 5, a high number of accepted answers is represented by a mass on the focal element "**Expert**" and the reminder to the ignorance, for a user  $i$ :

$$\begin{aligned} m_5^i(E) &= \alpha_5(1 - e^{-\gamma_5 NbAccAn_i}) \\ m_5^i(\Omega) &= \alpha_5 e^{-\gamma_5 NbAccAn_i} \end{aligned} \quad (12)$$

In the previous equations, we fix  $\alpha_1, \alpha_5 = 0.9$ ,  $\alpha_2 = 1$ ,  $\alpha_3 = 0.8$  and  $\alpha_4 = 0.5$ . The values are fixed after several experimentation's in order to have the best representation of each class of users. These values are used to represent the ignorance in every mass function as described in [3]. As the apprentices are modeled only one time as focal element in equation (10) unlike experts and occasionnals, we choose to affect the value of 1 to  $\alpha_2$ . For  $\gamma$  after several experimentation, we decide to keep it as the maximum value of any attribute divided by 100.

**Example 1:** Let us consider a question posted by a user  $u_1$  in the online community. Two other users  $u_2$  and  $u_3$  will read the question and will try to identify to which class can the asker belong: occasional, apprentice A, Expert E, where  $\Omega = \{O, A, E\}$ .

The corresponding power set  $2^\Omega = \{\emptyset, O, A, O \cup A, E, O \cup E, A \cup E, \Omega\}$ .

To express their beliefs on the question asker,  $u_2$  will say that this person is an expert at 80% and 20% ignorance ( $u_2$  does not know). User  $u_3$  would say this person could be an expert or an apprentice with a belief of 70% and 20% as an occasional and 10% of ignorance. We obtain the following mass functions:

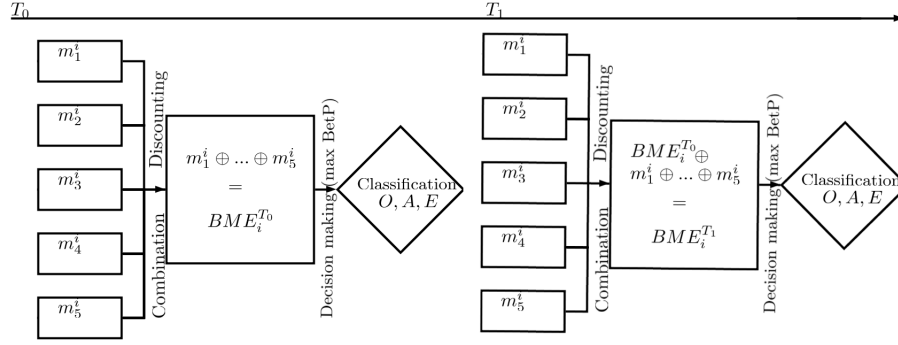
$$m_{u_2}(E) = 0.8, \quad m_{u_2}(\Omega) = 0.2 \quad (13)$$

$$m_{u_3}(E \cup A) = 0.7, \quad m_{u_3}(O) = 0.2, \quad m_{u_3}(\Omega) = 0.1 \quad (14)$$

Based on these beliefs, we will explain later how do we proceed to obtain to which class does user  $u_1$  belong.

### 4.3 Data aggregation and decision making

Coming to the combination of the belief functions for each feature, we adopt the Dempster's combination rule described in equation (5) for every time bucket. Finally, the



**Fig. 1.** Flow chart of the Belief Measure of Expertise and decision making

decision process is assured by the pignistic probability (BetP) described in equation (7). The final estimated classification label is the one having the higher pignistic probability.

The bucketing system provides us an overview of users' activity in the for a given period of time. For a user, we can then calculate the number of questions, answers and accepted answers posted by a given user during each time snap.

For every 30 days in the data set, we calculate for each user the number of questions asked, answers posted, the scores generated and the number of accepted answers. Each value is transformed into mass functions using Equations (8) to (12). Thus for every period, we obtain for every user 5 features: 5 mass functions for the number of questions, number of answers, score of questions, score of answers and a mass function for the accepted answers.

At  $t_0$  we combine these mass functions using the Dempster's combination rule presented in equation (5). Next we apply the pignistic probability and classify the user into Expert, Apprentice or Occasional for this specific time bucket. In  $t_1$ , we use the results of the previous period and combine them with the mass functions of this actual period. After, we define the class of belonging. We maintain this combination and classification process for the entire dataset.

The combination process allows us to estimate the actual belief expertise (noted BME) for each user during a period is expressed by the following equation:

$$BME_{t_1}(u_i) = \alpha^T BME_{t_0}^i \oplus m_1^i \oplus m_2^i \oplus \dots \oplus m_5^i \quad (15)$$

where  $\alpha^T$  is the discounting coefficient related to the time activity of a user. The value  $\alpha^T$  is the inverse of the number of days since the user first connected to the platform. The symbol  $\oplus$  represents the operator of combination.

BME will be in the interval  $[0, 1]$ . This process of combination and classification for every time bucket allows to follow the progress of users monthly during a defined period of time. Furthermore, based on that, we can distinguish clearly the evolution of each user during their time activity within the community. Thus, we can also detect potential experts on the onset of their participation.

**Example 2** We keep the same belief functions described in Example 1. We want to determine to which class does user  $u_1$  belong to. Let's assume that the user has been active on the platform for only 30 days. After defining the mass functions previously, during this step, we will first discount the masses  $m_{u_2}$  and  $m_{u_3}$  based on his time of activity by using  $\alpha_{u_1}^T = (1/30)$  and then, combine them using the Dempster's combination rule. We obtain the following results:

$$\begin{aligned} m_{\oplus}(O) &= 0, & m_{\oplus}(E) &= 0.0292, \\ m_{\oplus}(A \cup E) &= 0.0292, & m_{\oplus}(\Omega) &= 0.9441 \end{aligned} \quad (16)$$

After applying the pignistic transformation, we obtain the following probabilities:

$$Bet(O) = 0.3147, \quad Bet(A) = 0.3293, \quad Bet(E) = 0.3560. \quad (17)$$

We choose the highest probability, thus the user is defined as an Expert.

## 5 Experimental evaluation and analysis

The first step in this analysis of users is to build the temporal series of number of questions, answers and accepted answers given by users during a period of time. To do this, we divide the periods of the dataset into monthly and bi-weekly buckets. The begging of the first bucket is be the time of the earliest question in the dataset, noted  $t_0$ , and the end of the first bucket would be  $t_0 + 30$  days. We work on data covering 15 months allowing us to have 15 time snaps for monthly buckets.

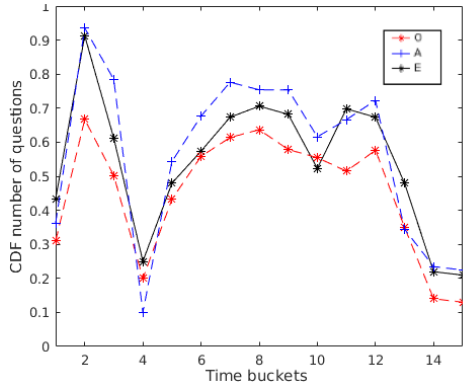
### 5.1 Time analysis of the data set

Figure 2 shows the cumulative distribution functions (CDF) related to the mean of the number of questions posted by contributors over a period of time of several months. We notice that apprentices ask more questions than the other users. This is due to the fact that they are seeking for information, and that they lack of knowledge.

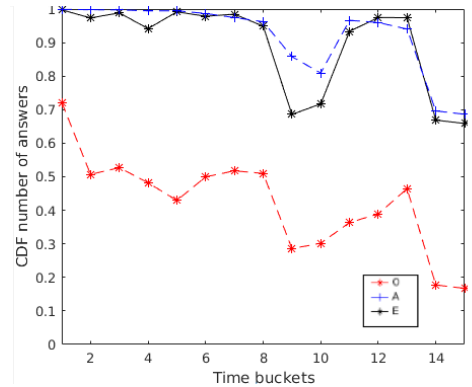
We also witness that experts do also questions. This can be justified by the idea that experts can not know everything about anything: they are knowledgeable on some specific topics only. Moreover, they are known to post difficult questions that only other experts can answer.

The CDF related to the mean of the number of answers is represented in Figure 3. We notice the same phenomenon described for the CDF of questions. At the beginning both experts and apprentices have almost similar values. However, over time, experts are less and less present within the community. They do not post as much answers as the apprentices. Though, the latter users try to provide a lot of contributions because they are motivated by gaining reputation points in Stack Overflow, sometimes without taking care of the quality of their posts. The fact they anyone posts answers may discourage experts to sharing their knowledge on the platform. This can cause the decrease of their interest on posting helpful answers.

The number of accepted answers is a very important indicator on how to evaluate the expertise of a user in Stack Overflow. Over time buckets, the CDF of the number of best

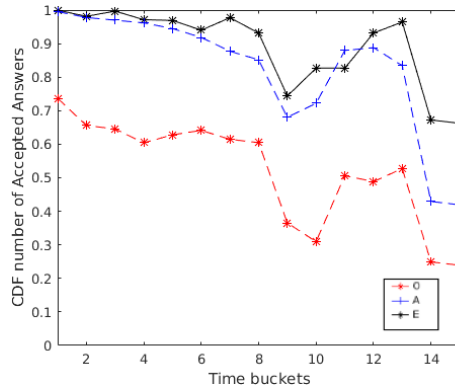


**Fig. 2.** CDF number of Questions



**Fig. 3.** CDF number of Answers

answers provided by each class of users is presented in 4. We notice that apprentices becoming future experts post a lot of answers that are considered as best. The more time they spend on the community the more expertise they have. However, both experts and apprentices lose interest in the community which is reflected by the decrease of their contributions over time.



**Fig. 4.** CDF number of Accepted Answers

**5.2 Analysis of users’ behavior over time**

In this section we provide an analysis of the activity of the users during the 15 months of the dataset. As described before, we classify users according to the belief expertise measure presented in equation (15) for every time bucket. We randomly choose  $n$  users from the big dataset and we obtain the results presented in Figure 5.

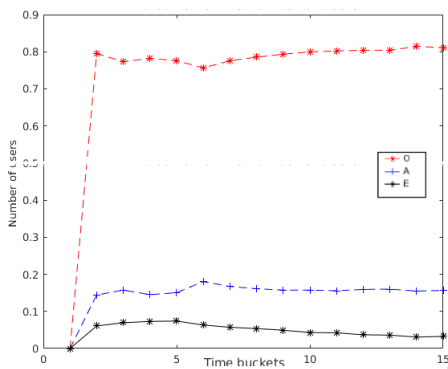


Fig. 5. Evolution of the percentage of Occasionals, Apprentices Experts per time bucket

First of all, we notice that the number of Occasionals is always the highest class of users present in the platform. After that, proportionally to the number of newbies, apprentices are not that numerous. However, we witness that their number changes over the months. Finally, for the experts, we find that their number fluctuates for the period of time described in the dataset. For the last time buckets, they become more and more scarce. The community may risk high-potential users leaving because of the lack of recognition regarding their efforts by other contributors.

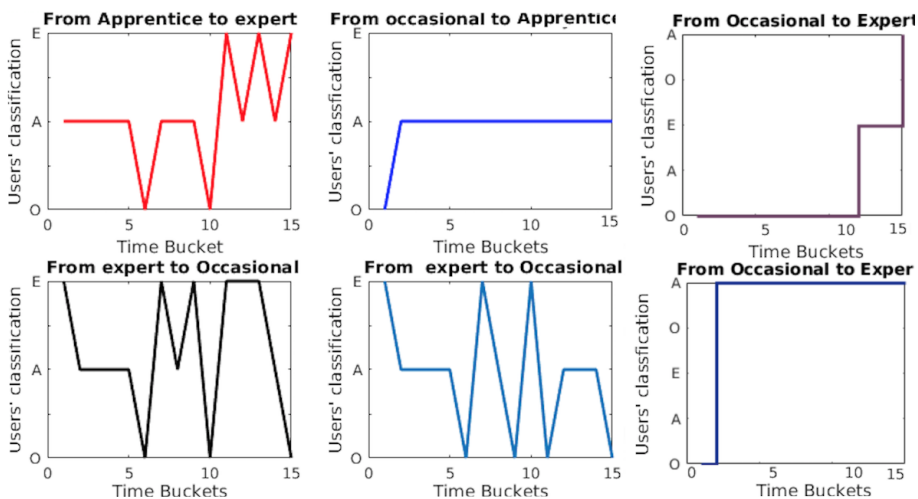


Fig. 6. Evolution of users over time

In Figure 6, we present the evolution of classification of some users during the hole time buckets. Some of them stay always as experts and some stay always as occasionals. However, we can witness the evolution of persons over the time spent within the community. we can see that contributors may evolve during their time activity in the community from occasional to apprentice to expert. Thus, we may notice that some of interrupt their contribution for some months and then restart posting. For some other users. Therefore, we can find users that can be experts for a period and then start posting less and less until leaving the community becoming occasionals.

In Figure 7, we present some values of the Belief Measure of Expertise for different classes of users. The BME is the mass affected to the experts. We notice that the value of this expertise measure for experts is high and always close to 1. However, for Occasionals it is very low with a  $BME = 0.1$  and decreases over time to 0 if this user does not contribute anymore. Therefore, for the Experts who Apprentices then Occasionals, the value of BME fluctuates over time until reaching 0.

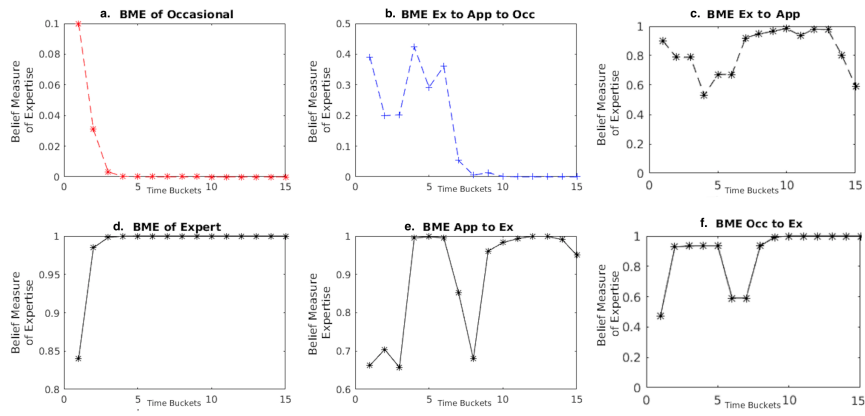


Fig. 7. Belief Measure of Expertise

With the value of BME of new users during the first months of their activity, we are able to detect future experts based on their posts and the time spent within the community. The BME being the mass allocated to the experts, potential experts users are detected early based on their behavior. Some of them are identified since the 2 or 3 times buckets like presented in Figure7. f where their BME increases from 0.45 to 0.9 in 2 months. However, some of them need a lot more time to acquire knowledge.

## 6 Conclusion

This paper is focused on two major issues: first on identifying three classes of users on Stack Overflow: Occasionals, Apprentices and Experts. Then, detecting potential experts on their early time of activity. The strength of the proposed model is that it could be

applied to any topic in the platform. Based on a belief model of the users' behavior, we calculated the general degree of expertise called the BME. This measure takes into account the combination of all the masses that describe a user during a defined period of time of activity on the web site. Once the expertise measure calculated for each time bucket, it allows us to have an overview on the users' behavior. Potential experts can be detected since the early few months of their entrance to the community. In future works, we will search to study the expertise of users based on their topical interests.

## References

1. Bouguessa, M., Romdhane, L.B.: Identifying authorities in online communities. *ACM TIST* 6(3), 30 (2015)
2. Dempster, A.P.: Upper and lower probabilities induced by a multivalued mapping. In: *Classic Works of the Dempster-Shafer Theory of Belief Functions*, vol. 219 (1967)
3. Denoeux, T.: A k-nearest neighbor classification rule based on dempster-shafer theory. *IEEE Trans. Systems, Man, and Cybernetics* 25(5), 804–813 (1995)
4. Denoeux, T., Kanjanatarakul, O.: Evidential clustering: A review. In: *IUKM. Lecture Notes in Computer Science*, vol. 9978, pp. 24–35 (2016)
5. van Dijk, D., Tsagkias, M., de Rijke, M.: Early detection of topical expertise in community question answering. In: *SIGIR*. pp. 995–998. ACM (2015)
6. Dubois, D., Prade, H.: Possibility theory and its applications: Where do we stand? In: *Handbook of Computational Intelligence*, pp. 31–60. Springer (2015)
7. Kasneci, G., Gael, J.V., Stern, D.H., Graepel, T.: Cobayes: bayesian knowledge corroboration with assessors of unknown areas of expertise. In: *WSDM*. pp. 465–474. ACM (2011)
8. Khaleghi, B., Khamis, A.M., Karray, F., Razavi, S.N.: Multisensor data fusion: A review of the state-of-the-art. *Information Fusion* 14(1), 28–44 (2013)
9. Movshovitz-Attias, D., Movshovitz-Attias, Y., Steenkiste, P., Faloutsos, C.: Analysis of the reputation system and user contributions on a question answering website: Stackoverflow. In: *Rokne, J.G., Faloutsos, C. (eds.) ASONAM*. pp. 886–893. ACM (2013)
10. Nguyen, V.D., Huynh, V.N.: Integrating with social network to enhance recommender system based-on dempster-shafer theory. In: *CSoNet. Lecture Notes in Computer Science*, vol. 9795, pp. 170–181. Springer (2016)
11. Pal, A., Farzan, R., Konstan, J.A., Kraut, R.E.: Early detection of potential experts in question answering communities. In: *UMAP. Lecture Notes in Computer Science*, vol. 6787, pp. 231–242. Springer (2011)
12. Pal, A., Harper, F.M., Konstan, J.A.: Exploring question selection bias to identify experts and potential experts in community question answering. *ACM Trans. Inf. Syst.* 30(2) (2012)
13. Posnett, D., Warburg, E., Devanbu, P.T., Filkov, V.: Mining stack exchange: Expertise is evident from initial contributions. In: *SocialInformatics*. pp. 199–204 (2012)
14. Reyni, A.: *Probability Theory*. North-Holland (1962)
15. Sahu, T.P., Nagwani, N.K., Verma, S.: Multivariate beta mixture model for automatic identification of topical authoritative users in community question answering sites. *IEEE Access* 4, 5343–5355 (2016)
16. Shafer, G.: *A Mathematical Theory of Evidence* (1976)
17. Smets, P., Kennes, R.: The transferable belief model. *Artificial Intelligence* 66, 191–234 (1994)
18. Tang, X., Yang, C.C.: Ranking user influence in healthcare social media. *ACM TIST* 3(4), 73 (2012)
19. Yang, J., Tao, K., Bozzon, A., Houben, G.J.: Sparrows and owls: Characterisation of expert behaviour in stackoverflow. In: *UMAP*. vol. 8538, pp. 266–277. Springer (2014)