

Does Personalization Help? Predicting How Social Situations Affect Personal Values

Ilir KOLA^{a,1} Ralvi ISUFAJ^{a,b} and Catholijn M. JONKER^{a,c}

^a*Interactive Intelligence Group, Delft University of Technology, The Netherlands*

^b*School of Engineering, Autonomous University of Barcelona, Spain*

^c*Leiden Institute of Advanced Computer Science, The Netherlands*

Abstract. Personal values represent what people find important in their lives, and are key drivers of human behavior. For this reason, support agents should provide help that is aligned with the personal values of the users. To do this, the support agent not only should know the value preferences of the user, but also how different situations in the user's life affect these personal values. We represent situations using their psychological characteristics, and we build predictive models that given the psychological characteristics of a situation, predict whether the situation promotes, demotes or does not affect a personal value. In this work, we focus on predictions for the value 'enjoyment of life', and use different machine learning classifiers, all of them performing better than chance when training on data from multiple people. The best predictive model is a multi-layer perceptron classifier, which achieves an accuracy of 72%. Further, we hypothesize that the accuracy of such models would drop when tested on individual data sets. The data supports our hypothesis, and the accuracy of the best performing model drops by at least 11% when tested on individual data. To tackle this, we propose an active learning procedure to build personalized prediction models having the user in the loop. Results show that this approach outperforms the previously built model while using only 30% of the training data. Our findings suggest that how situations affect personal values can have subjective interpretations, but we can account for those subjective interpretations by involving the user when building a prediction model.

Keywords. Personal Values, Predictive Models, Active Learning, Support Agents

1. Introduction

Support agents that help users in their daily lives, such as personal assistants, health coaches etc., are becoming increasingly prevalent. It is a desideratum for such systems to offer support that is aligned with the personal values of the users [1]. Personal values represent what a person or group of people consider important in life [2], and this concept is seen as central for robust and beneficial AI [3].

¹Corresponding Author: Ilir Kola, Delft University of Technology, Van Mourik Broekmanweg 6, 2628 XE Delft, The Netherlands; E-mail: i.kola@tudelft.nl.

Existing work gives insight on how this can be achieved in practice. Generally, the support agent should be able to, on one hand, elicit which values are important to the user, and on the other hand, determine how different daily life activities affect those values. By using this information, the support agent can help the user conduct a lifestyle that is aligned with their personal values. For instance, Kayal et al. [4] demonstrate how combining information regarding which values are important to the user and which values are promoted by different social commitments can be used to automatically solve data sharing conflicts for users of location sharing platforms. Cranefield et al. [5] show how value preferences of the user can be included in the planning framework of a support agent in order to provide value-aligned support. However, these works do not focus on how to obtain this information. The topic of value elicitation has been widely researched, and different techniques have been proposed (e.g., [6,7]). Our focus in this work is on the other type of information: how the user's daily activities affect their personal values.

Personal values are an abstract concept, therefore building a model that given an activity determines how different values are affected is no trivial task. Different lists of values have been proposed (e.g. [8,9]) and formalized for use in AI systems (e.g. [10]), and that work serves as a starting point for us. Another challenge lies in how to conceptualize and formalize daily activities of the users. In this work, we view activities as *situations* in which the user finds themselves in their daily life. A way to formalize situations is through a textual description. That information can then serve as input to Natural Language Processing and machine learning techniques which predict the affected values, as done in [11]. However, the presence of values in language is often implicit, which makes this approach challenging. Furthermore, such an approach would require a large amount of data. Another option is to describe situations in terms of their psychological characteristics. For instance, Rauthmann et al. [12] posit that every situation can be described using characteristics such as duty, intellect etc., and their results show that different people would assign very similar psychological characteristics to a situation. In previous work [13], we show how psychological characteristics of situations can be used as a basis for social situation comprehension in support agents. Furthermore, in [14] we show that grouping situations based on their psychological characteristics gives insight into the personal values that are affected in these situations. Drawing on these insights, first of all we assess:

Research Question 1 (RQ1) - To what extent we can use machine learning techniques to predict how a situation affects personal values using as input the psychological characteristics of situations?

In practice, given a situation described in terms of its psychological characteristics and a personal value, we want to predict whether the situation promotes the value (i.e., helps the user achieve it), demotes the value (i.e., prevents the user from achieving it), or does not affect it. This classification is based on existing work in AI and social sciences [10,15]. Particularly, in this work we focus on the personal value *enjoyment of life*, and explore how social situations affect that value.

Collecting large amounts of data from the same user is difficult in practice because it is time consuming. For this reason, it is common to create data sets that combine observations from different users (e.g., [16]). While this is not an issue in domains where ground truths are less prone to subjective interpretations (e.g., image recognition), there can be issues in domains where labels are subjective. Results reported in Kayal et al. [4]

suggest that people are consistent when it comes to how they think an activity affects a personal value. However, their work is based on a narrow domain of activities that were designed by the researchers. Research in social sciences [15] introduces the concept of subjective value fulfillment to represent the extent to which people feel they can attain what they desire (i.e., fulfill their value preferences). As the name suggests, this process is considered to be subjective. For instance, some people feel that the activity ‘playing cards with friends’ promotes the value *enjoyment of life*, while others might feel like it demotes it, while everyone could assign similar psychological characteristics to the situation. We hypothesize that users are subjective when it comes to assessing how a situation affects a personal value, therefore:

Research Hypothesis (RH) - The prediction accuracy of the predictive model from RQ1 (trained on data from multiple users) declines when tested on an individual user.

Ideally, we want personal agents to make predictions that are calibrated to how situations affect values from the perspective of their user. As recognized by existing work, an intelligent agent must be designed to learn and act according to the preferences of its operators [3]. For this reason, we turn to active learning techniques. In such techniques, the algorithm iteratively identifies data points that it considers to be more informative, and asks an oracle (in this case, the user) to label these data points. The new data is incorporated in the learning procedure. This way, the agent can directly integrate feedback from the user. We explore the following research question:

Research Question 2 (RQ2) - Can we use active learning in order to build a personalized value prediction model that is more accurate on individual predictions than the model from RQ1?

Such an active learning approach would enable support agents to make personalized predictions, and provide user-centered support that is aligned with the personal values of the user. This contributes towards collaborative and adaptive AI in support agents.

The rest of this article is structured as follows. Section 2 provides the necessary background information. In Section 3 we introduce our proposed approach. In Section 4 we present how the data was collected, and in Section 5 we use that data to build and evaluate prediction models. Section 6 concludes the article.

2. Background

In this section we discuss concepts on which our work is built on.

2.1. Personal Values in Support Agents

Values represent key drivers of human decision making, and different lists of prominent human values have been proposed throughout the years [8,9]. According to Schwartz [17], in order for values to influence action not only should they be important to the actor, but they should also be salient in the situation where action will take place. Based on these insights, researchers have proposed that values should be taken into account when designing technology [2]. This has also had an impact on research on support agents, and different work highlights the importance of offering value-aligned support to users and

gives insight on how this can be achieved in practice. Cranefield et al. [5] include values in their planning system such as each action can help the user achieve specific values. By knowing which values are important to the user, the support agent selects actions that will help the user achieve the values that are important to them. Drawing on these insight, Tielman et al. [18] propose an approach that also includes how context influences the way in which values are affected by actions. Both these works focus on the computational frameworks needed to reason about integrating values in support agents. However, they assume the support agent already knows how actions affect personal values, and do not focus on how that information is acquired.

Kayal et al. [4] go one step further, not only they propose a reasoning framework for value-aligned support, but they also conduct a user study in which they elicit from the users information about which personal values are important to them, and how different actions affect those personal values in the domain of location sharing. Their findings suggest that personal values can inform support agents on how to resolve social commitment conflicts, thus showing in practice the importance of value-aligned support.

Recent work has focused on assessing which values are relevant in different domains of activities. Moonen and Tielman [19] show that it is possible to use expert knowledge to determine which value categories are more relevant for different daily life activities. Liscio et al. [20] go a step further, they propose a hybrid (human and AI) methodology for context-specific value identification. Specifically, the task of value identification is framed as a guided value annotation process of human annotators supported by Natural Language Processing techniques. Both these approaches can be used as a starting point for our work, since they allow the support agent to identify which values are relevant in a situation. For instance, they could inform a support agent that for the activity *going to a party* the values *creativity* or *national security* are not relevant while the value *enjoyment of life* is relevant. However, they do not provide information on whether those values are promoted or demoted in those situations from the point of view of a specific user. Our research aims to fill this gap.

2.2. Psychological Characteristics of Situations

Research in social psychology has explored ways in which situations can be systematically described. Among others, psychological characteristics of situations have been proposed as a way to formalize how people view situations. Different taxonomies have been developed (e.g., [12,21,22,23,24]), and they each propose a set of dimensions which can be used to represent situations. In this work, we refer to the DIAMONDS taxonomy [12], since it is meant to cover everyday situations and it offers a validated questionnaire that can be used in user studies. The taxonomy consists of the following dimensions:

- **Duty** - situations where a job has to be done, minor details are important, and rational thinking is called for;
- **Intellect** - situations that afford an opportunity to demonstrate intellectual capacity;
- **Adversity** - situations where you or someone else are (potentially) being criticized, blamed, or under threat;
- **Mating** - situations where potential romantic partners are present, and physical attractiveness is relevant;
- **Positivity** - playful, simple, clear-cut and potentially enjoyable situations;

- **Negativity** - stressful, frustrating, and anxiety-inducing situations;
- **Deception** - situations where someone might be deceitful. These situations may cause feelings of hostility;
- **Sociality** - situations where social interaction is possible, and close personal relationships are present or have the potential to develop.

Rauthmann et al. [12] conducted studies in which they presented users with situation descriptions which had to be rated in terms of the above mentioned dimensions, and they showed moderately strong correlation between the answers of different users, suggesting a high level of agreement on how people interpret the psychological characteristics of a situation.

2.3. Active Learning

Active Learning refers to a subset of machine learning techniques based on the hypothesis that if the learning algorithm is allowed to choose the data from which it learns, it will perform better with less training [25]. This is particularly relevant in domains where labeled data is not freely available. For instance, active learning has been used in activity recognition [26], recommender systems [27], and hearing aids personalization [28]. All these tasks have in common the fact that subjective interpretations are relevant, like we hypothesized is the case with personal values. In the pool-based active learning setting, the principle is that given a set of unlabeled data and a learning algorithm, active learning uses an uncertainty sampling query technique which selects the situations for which the learning algorithm is least certain what label to assign and asks an oracle (e.g., human annotator) to label those instances, which are added to the training set (right part of Figure 3). By selecting the most informative data points, active learning allows for better performance while requiring only a subset of the data to be labeled.

3. Proposed Approach

The objective of this work is to enable support agents to automatically assess how personal values are affected by different activities from the user's point of view. In this section, we present our proposed approach, its merits and its limitations.

3.1. Design Choices and Data Collection

To model activities, we conceptualize them as situations. Specifically, in this work we focus on social situations, since that accounts for a wide number of our life situations. Each situation is represented through its psychological characteristics, a concept that has been shown to be a good predictor both for human behavior and values that are afforded in situations [12,14]. When it comes to values, we take them from the Schwartz value survey [8]. Following insights from AI (e.g., [10]) and social sciences (e.g., [15]), we assume that a situation can either promote a personal value, demote it, or not affect it. Therefore, at this point our goal is to build a predictive affect model that takes as input the psychological characteristics of a situation, and predicts whether that situation promotes, demotes, or does not affect a specific personal value, from the point of view of the user.

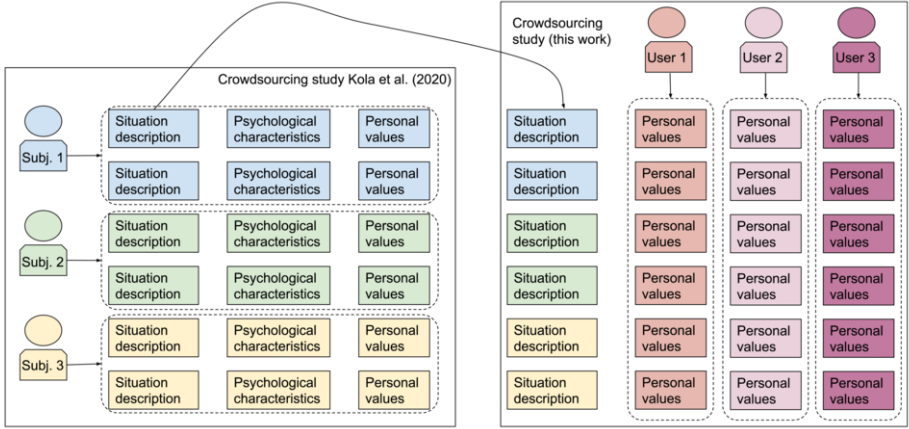


Figure 1. Experimental setup for collecting the data. Participants of our crowdsourcing study were presented with the situation descriptions collected in the crowdsourcing study of Kola et al. [14].

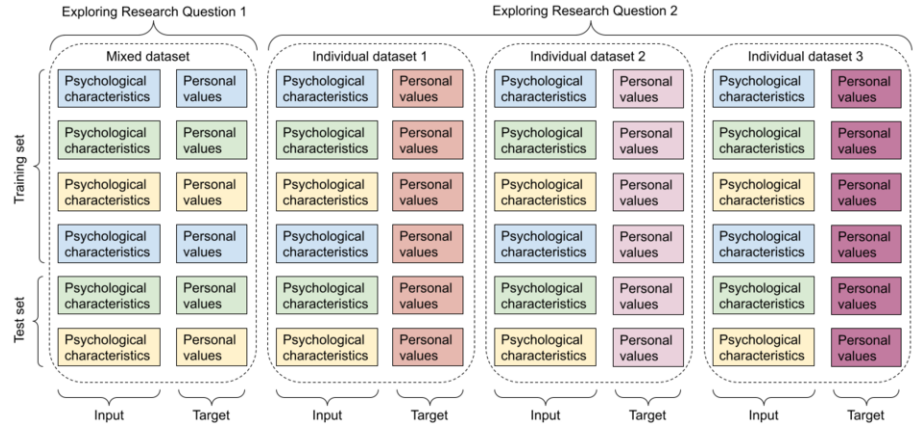


Figure 2. Different data sets used to train and evaluate the models. The mixed data set consists of the data from the participants of the crowdsourcing study of Kola et al. [14], whereas the individual data sets consists of the data from the participants of our crowdsourcing study.

Ideally, to train such a model, we would need data from a large number of situations from the user’s life, annotated by the user in terms of psychological characteristics of situations as well as personal values. However, that is difficult in practice because it requires continuous interaction for a long amount of time. To overcome this issue, in previous work [14] we conducted a crowd-sourcing study in which different users described two social situations from their past few days (left part of Figure 1). Each situation was described textually², as well as through its psychological characteristics and personal

²An example of a situation description: ‘I had a talk with my boss about clear work expectations. There have been many grey areas lately and I was trying to clear them up. After a slightly confrontational conversation, we were on the same page and have more clarity going forward.’

values. This can be used to explore RQ1, however, we would not have information about the subjectivity of how people view values and therefore we cannot test RH and RQ2.

For this reason, we conduct a new crowdsourcing study in which each participant is presented with the textual situation descriptions from Kola et al. [14], and for each situation annotates whether a personal value is promoted, demoted or not affected (right part of Figure 1). These annotations are *ex-situ*, since participants are not annotating situations from their lives, but rather situations which they have not participated in. The motivation for this choice is presented in Section 3.3. With the data from the crowdsourcing study we create individual data sets (Figure 2). The annotated values from these data sets can be used to test RH and RQ2.

3.2. Prediction Models

The collected data can be used to address the research questions and hypothesis. By using the mixed data set, we can build a predictive model (like in the left part of Figure 3) that is trained on data coming from multiple people (Figure 2, mixed data set - training set) and explore RQ1 by testing it on the mixed test set and RH by testing it on the individual test sets. With the individual data sets we then also explore RQ2, as illustrated in the right part of Figure 3. The model starts with an initial amount of situations from the mixed data set. All other situations from the training set are considered to be unlabeled. To build a personalized prediction model, in each iteration, the active learning module selects the unlabeled situation that deems to be more informative, and poses a query to the oracle (in this case, the individual data set of the user for which we are personalizing the model). The now labeled situations are added to the training set, and the model is trained again. After a fixed amount of iterations, the personalized prediction model can be evaluated in order to answer RQ2.

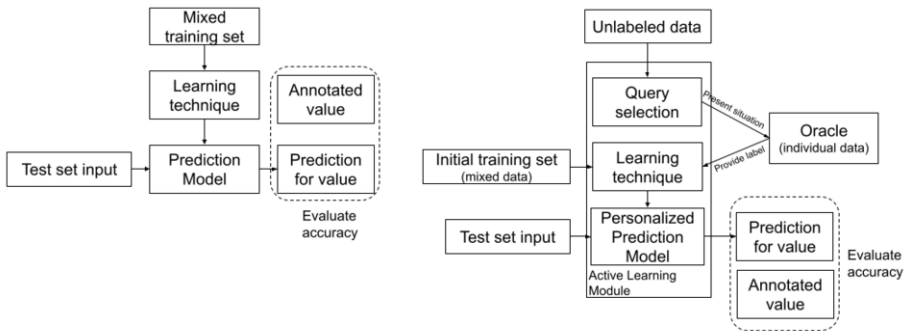


Figure 3. Regular learning and evaluation procedure used to build the general prediction model (left), and active learning procedure used to build the personalized prediction models (right). Different learning techniques were used.

3.3. Analyzing the assumptions

Our work has an exploratory nature, therefore some of the design choices are based on assumptions that can add noise to the predictive models. In this section we analyze those assumptions.

First of all, in order to build the individual data sets we ask participants to annotate the personal values of relatively short textual descriptions of situations provided by other people. This process is often referred to as *ex-situ* rating on thin slices of information [12], and is common in social psychology research. However, to truly assess the subjectivity of how people view that situations affect values, in future work it will be important to collect a large corpus of situations from the same users interacting with an agent for an extended amount of time.

Furthermore, participants in our crowdsourcing study are presented with the textual description of situations, however the model takes as input the psychological characteristics of a situation instead. This choice was made because on one hand, using as model input the textual descriptions of situations would require a high amount of data, and on the other hand, presenting participants with the psychological characteristics of the situations would create a high cognitive load since they are not concepts that we think about in our daily lives. This can add noise to the model, since there can be some level of variance and subjectivity in how people interpret the psychological characteristics of a situation. In this work, we consider that when reading the description of a situation different users would implicitly assign similar psychological characteristics to the situation, thus the model input would have been very similar across users. This is based on the findings of Rauthmann et al. [12]. In their work, in a setting which is similar to ours, users are presented with verbal descriptions of situations provided by other people, and rate them in terms of their psychological characteristics. Results show moderately strong correlation between the ratings of different users (0.64 on average) across the psychological characteristics, which suggests that different people rate the psychological characteristics of situations similarly.

4. Data collection

In this section we introduce the data sources that were used to build the models.

4.1. Existing Data

As a starting point, we use the data set collected in Kola et al. [14]. In that study, 150 participants described two social situations from their daily lives. For each situation, participants provided a textual description, completed the S8* questionnaire [29] to determine the psychological characteristics of the situations, as well as were presented with a list of personal values, and they were asked on a slider with a range from -10 (fully demote) to 10 (fully promote), how much is each value promoted or demoted in each of their two situation. After being read by two researchers and controlled for spelling errors, the final data set contains 283 situations.

We group the answers regarding personal values as follows: for each personal value, situations where the personal value got an answer from -10 to -4 are labeled as ‘Demotes personal value’, situations where the personal value got an answer from -3 to 3 are labeled as ‘Does not affect personal value’, and situations where the personal value got an answer from 4 to 10 are labeled as ‘Promotes personal value’. This grouping was done to facilitate the classification task when predicting personal values, since now there are 3 target labels from the original 21. This data is used to form the mixed data set (Figure 2).

4.2. Crowd-sourcing study

We conducted a crowd-sourcing study to collect the data needed for the individual data sets. The study design was approved by the university's ethics committee.

4.2.1. Material

We use the descriptions of the 283 situations from Kola et al. [14]. Each situation is written in first person from the point of view of the subject who described it. For the purpose of this work we focus on the personal value *enjoyment of life*, since in the mixed data set it is the value that affects (ie., promotes or demotes) more situations.

4.2.2. Participants

We recruited 8 subjects through sharing the survey with university employees and through social media. No demographic information was collected in order to ensure that subjects could not be identified, due to the small sample size.

4.2.3. Procedure

Subjects completed an online survey. After being briefed about the purpose of the study, subjects were explained the procedure. Then, subjects were presented with textual descriptions of the social situations taken from the study in Kola et al. [14]. For each situation, subjects were asked "From your perspective, how does each situation affect the value enjoyment of life for the user?". The answer options were 'Promotes enjoyment of life', 'Demotes enjoyment of life', and 'Does not affect enjoyment of life'. The procedure is described in the right side of Figure 1. The data from each subject is then used to form the individual data sets (Figure 2).

5. Models and Results

In this section we describe how the collected data was used to build and evaluate predictive models³.

5.1. Value Prediction Model

Our first goal is to assess RQ1, namely to what extent we can use machine learning techniques to predict how a situation affects personal values. To tackle this, we train and evaluate different models on the mixed data set (Figure 2). The followed procedure is illustrated in the left side of Figure 3. We start by dividing the data set into a training set that contains 75% of the data and a test set that contains 25% of the data. We use the training set to build four classification models: a Decision Tree classifier, a Random Forest classifier, the XGBoost classifier, and a Multi-Layer Perceptron (MLP) classifier. For this, we use the implementations from the Scikit-learn package in Python [30]. These classifiers were selected to account for some of the most common machine learning techniques used on tabular data sets: decision trees, ensemble techniques, and neural

³The code and data can be accessed in: <https://doi.org/10.4121/19292246>

networks. Our goal was not to find the best possible classifier, but rather to explore how different learning techniques perform in this task. Furthermore, we add as baseline a model that makes random predictions, as commonly done in new machine learning tasks with no predetermined benchmarks [31]. Each model takes as input the psychological characteristics of a situation, and has to predict whether the situation promotes, demotes or does not affect the value *enjoyment of life*. In Table 1 we report the performance of the models on the test set. As we can notice, the MLP classifier outperforms all the other prediction models, with a prediction accuracy of 72%.

Table 1. Performance of different prediction models when tested on the mixed data set. The calculation of precision, recall and F1-score is weighted to take into account possible class imbalances.

Model	Accuracy	Precision	Recall	F1-score
MLP Classifier	0.72	0.74	0.71	0.68
Random Forest Classifier	0.66	0.66	0.66	0.63
XGBoost Classifier	0.67	0.66	0.67	0.64
Decision Tree Classifier	0.57	0.55	0.58	0.56
Random Classifier	0.34	0.41	0.34	0.36

Our hypothesis (RH) was that the accuracy of the model trained on the mixed data set drops when tested on the individual data sets. To evaluate this hypothesis, we take the predictive models trained on the mixed data set (i.e., data from participants of the study in Kola et al. [14]) and we evaluate their predictions on each individual data set collected from our crowd-sourcing study. The accuracies are reported in Table 2. As we can see, the accuracy drops when predictions are tested on data from individual users, thus supporting our hypothesis. This does not hold for the baseline model, which is to be expected since there is no learning taking place there. We also notice that the drop in accuracy in the decision tree classifier is smaller, which can be explained by the fact that the accuracy of that model is overall low.

Table 2. Accuracy of classifiers trained on the mixed data set when evaluated on the individual data sets.

Model/Test set	Mixed data	User 1	User 2	User 3	User 4	User 5	User 6	User 7	User 8
MLP Classifier	0.72	0.61	0.55	0.48	0.58	0.53	0.59	0.46	0.42
Random Forest	0.66	0.54	0.51	0.48	0.62	0.54	0.62	0.44	0.45
XGBoost	0.67	0.59	0.54	0.54	0.64	0.52	0.57	0.46	0.48
Decision Tree	0.57	0.56	0.52	0.46	0.55	0.49	0.55	0.45	0.44
Random Classifier	0.34	0.31	0.35	0.37	0.32	0.34	0.31	0.34	0.36

5.2. Active Learning for Personalized Predictions

In this section we explore RQ2, namely evaluating whether we can use active learning in order to build personalized value prediction models that are more accurate on individual predictions than the prediction model trained on the mixed data set. We use the active learning framework *modal* in Python [32]. As a base classifier we use the MLP classifier since it was the prediction model that achieved the best performance in Table 1. We build personalized prediction models for each user, using the individual data sets (Figure 2).

For each personalized prediction model, we start with 10 situations from the mixed data set, and remove the labels for all other situations. An initial prediction model is built

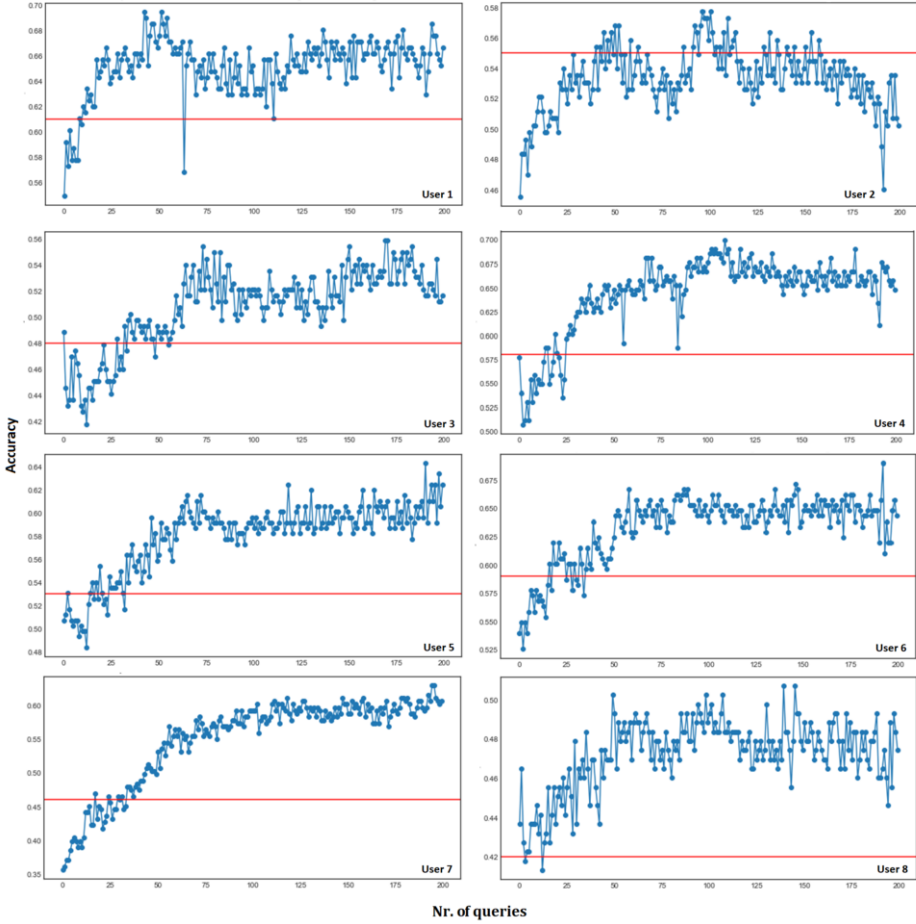


Figure 4. Accuracy of personalized models after each query. The red lines represent the accuracy of the prediction model from Section 5.1 on each individual data set (Table 2). The x axis represents the number of queries, and the y axis represents the accuracy of the models after each query.

using as input the 10 situations. Then, the active learner starts an interactive procedure in which it selects the unlabeled instance for which its current prediction is maximally uncertain, i.e. the instance for which the classification probability is the lowest. A query is posed for that instance, and the label is retrieved from the individual data set. The labeled instance is added to the training set, and the model is re-trained. We allow the model to pose 200 queries, since that was the amount of unlabeled data. The procedure is described in the right part of Figure 3.

We report the accuracies of each model in Figure 4. For 7 out of the 8 individual data sets, the personalized prediction model outperforms the accuracy on those data sets of the prediction model trained on the mixed data set. Furthermore, we notice that the accuracy tends to converge after approximately 70 queries, suggesting that active learning can identify the most informative instances in this domain. This greatly helps practical applications since it reduces the labeling burden for the users. We also notice that the

accuracies of some personalized models are lower than others. This can be explained by the different levels of noise introduced by the fact that users were labeling other people's situations and they were presented with textual descriptions while the model had as input psychological characteristics of the situation (see Section 3.3).

6. Conclusion

6.1. Contributions

Our goal was to predict how social situations affect personal values of users. This information can be used by support agents to provide value-aligned support. In exploring RQ1, we show that different learning algorithms that take as input the psychological characteristics of a situation can predict how that situation affects the personal value *enjoyment of life* better than baseline predictors. The highest accuracy of 72% is achieved by the MLP classifier, outperforming the other used models. This result is achieved for a model trained on data coming from multiple people. We hypothesize that people are subjective when assessing how a situation affects personal values, and thus that the accuracy of the models drops when tested on individual data. To evaluate this, we collected data from eight participants of an online study. The accuracy of the MLP classifier dropped on each individual data set by at least 11%, thus supporting our hypothesis. The accuracy also dropped in the other predictive models. To achieve better individual performance, we used an active learning procedure that builds personalized prediction models (RQ2). Results suggest that this approach leads to an improved accuracy in 7 of the 8 individual data sets. Furthermore, the improved accuracy is reached by using approximately 70 labeled data points, thus reducing the annotation burden in possible applications. Overall, our results suggest that an active learning approach that includes the user in the loop is highly beneficial in order to achieve personalized predictions.

6.2. Future Work

Based on our findings, a larger study can be conducted in which more situations are collected from the same user over continuous interactions. This would allow for truly subjective and personalized prediction models. A larger amount of data would also contribute in building more robust and accurate models, which is crucial for implementing this procedure in support agents. Another interesting line of research would be to include in the study online active learning techniques, in which for each new situation the model decides whether to ask the user for labels or not. Such an approach would on one hand benefit from having the user in the loop, and on the other hand would reduce the amount of interruptions that the user gets from the support agent. Towards a practical application, more research is also needed on how to perceive user situations and how to automatically assess the psychological characteristics of situations [13,33]. Next, it would be important to replicate our results for other personal values. Lastly, it would be interesting to assess similarities in participants' subjective views, and whether it is possible to benefit from those similarities like it is done in collaborative filtering.

Acknowledgements

The first author is part of the research programme CoreSAEP, with project number 639.022.416, which is financed by the Netherlands Organisation for Scientific Research (NWO). The second author is funded by the SESAR Joint Undertaking under the EU H2020 Research and Innovation Programme under grant agreement No 783287. The third author is funded by the [Hybrid Intelligence Center](#), a 10-year programme funded the Dutch Ministry of Education, Culture and Science through NWO grant number 024.004.022 and by EU H2020 ICT48 project “Humane AI Net” under contract #952026.

References

- [1] Akata Z, Balliet D, De Rijke M, Dignum F, Dignum V, Eiben G, et al. A research agenda for hybrid intelligence: augmenting human intellect with collaborative, adaptive, responsible, and explainable artificial intelligence. *Computer*. 2020;53(08):18-28.
- [2] Friedman B, Kahn PH, Borning A. Value sensitive design and information systems. *The handbook of information and computer ethics*. 2008:69-101.
- [3] Soares N. The value learning problem. In: *Artificial Intelligence Safety and Security*. Chapman and Hall/CRC; 2018. p. 89-97.
- [4] Kayal A, Brinkman WP, Neerincx MA, Riemsdijk MBV. Automatic resolution of normative conflicts in supportive technology based on user values. *ACM Transactions on Internet Technology (TOIT)*. 2018;18(4):1-21.
- [5] Cranefield S, Winikoff M, Dignum V, Dignum F. No Pizza for You: Value-based Plan Selection in BDI Agents. In: *IJCAI*; 2017. p. 178-84.
- [6] Le Dantec CA, Poole ES, Wyche SP. Values as lived experience: evolving value sensitive design in support of value discovery. In: *Proceedings of the SIGCHI conference on human factors in computing systems*; 2009. p. 1141-50.
- [7] Manolios S, Hanjalic A, Liem CC. The influence of personal values on music taste: towards value-based music recommendations. In: *Proceedings of the 13th ACM Conference on Recommender Systems*; 2019. p. 501-5.
- [8] Schwartz SH. Universals in the content and structure of values: Theoretical advances and empirical tests in 20 countries. *Advances in experimental social psychology*. 1992;25(1):1-65.
- [9] Rokeach M. *The nature of human values*. Free press; 1973.
- [10] van der Weide TL, Dignum F, Meyer JJC, Prakken H, Vreeswijk GA. Practical reasoning using values. In: *International Workshop on Argumentation in Multi-Agent Systems*. Springer; 2009. p. 79-93.
- [11] Wilson S. *Natural language processing for personal values and human activities*. PhD Thesis, University of Michigan; 2019.
- [12] Rauthmann JF, Gallardo-Pujol D, Guillaume EM, Todd E, Nave CS, Sherman RA, et al. The Situational Eight DIAMONDS: A taxonomy of major dimensions of situation characteristics. *Journal of Personality and Social Psychology*. 2014;107(4):677.
- [13] Kola I, Jonker CM, van Riemsdijk MB. Using psychological characteristics of situations for social situation comprehension in support agents. *arXiv preprint arXiv:211009397*. 2021.
- [14] Kola I, Jonker CM, Tielman ML, van Riemsdijk MB. Grouping Situations Based on their Psychological Characteristics Gives Insight into Personal Values. In: *11th International Workshop Modelling and Reasoning in Context*; 2020. p. 17-26.
- [15] Oppenheim-Weller S, Roccas S, Kurman J. Subjective value fulfillment: A new way to study personal values and their consequences. *Journal of Research in Personality*. 2018;76:38-49.
- [16] Kola I, Tielman ML, Jonker CM, van Riemsdijk MB. Predicting the Priority of Social Situations for Personal Assistant Agents. In: *International Conference on Principles and Practice of Multi-Agent Systems*. Springer; 2020. .
- [17] Schwartz SH. An overview of the Schwartz theory of basic values. *Online readings in Psychology and Culture*. 2012;2(1):2307-0919.
- [18] Tielman ML, Jonker CM, van Riemsdijk MB. What should I do? Deriving norms from actions, values and context. In: *MRC@ IJCAI*; 2018. .

- [19] Moonen DD, Tielman ML. Linking Actions to Value Categories-a First Step in Categorization for Easier Value Elicitation.; 2020. .
- [20] Liscio E, van der Meer M, Siebert LC, Jonker CM, Mouter N, Murukannaiah PK. Axes: Identifying and Evaluating Context-Specific Values. In: Proceedings of the 20th International Conference on Autonomous Agents and MultiAgent Systems; 2021. p. 799-808.
- [21] Parrigon S, Woo SE, Tay L, Wang T. CAPTION-ing the situation: A lexically-derived taxonomy of psychological situation characteristics. *Journal of personality and social psychology*. 2017;112(4):642.
- [22] Gerpott FH, Balliet D, Columbus S, Molho C, de Vries RE. How do people think about interdependence? A multidimensional model of subjective outcome interdependence. *Journal of Personality and Social Psychology*. 2018;115(4):716.
- [23] Ziegler M. Big Five Inventory of personality in occupational situations. Mödling, Austria: Schuhfried GmbH. 2014.
- [24] Brown NA, Neel R, Sherman RA. Measuring the evolutionarily important goals of situations: Situational affordances for adaptive problems. *Evolutionary Psychology*. 2015;13(3):1-15.
- [25] Settles B. Active learning literature survey. 2009.
- [26] Miu T, Missier P, Plötz T. Bootstrapping personalised human activity recognition models using online active learning. In: 2015 IEEE International Conference on Computer and Information Technology; Ubiquitous Computing and Communications; Dependable, Autonomic and Secure Computing; Pervasive Intelligence and Computing. IEEE; 2015. p. 1138-47.
- [27] Rubens N, Elahi M, Sugiyama M, Kaplan D. Active learning in recommender systems. In: *Recommender systems handbook*. Springer; 2015. p. 809-46.
- [28] Nielsen JBB, Nielsen J, Larsen J. Perception-based personalization of hearing aids using Gaussian processes and active learning. *IEEE/ACM Transactions on Audio, Speech, and Language Processing*. 2014;23(1):162-73.
- [29] Rauthmann JF, Sherman RA. Measuring the Situational Eight DIAMONDS characteristics of situations: An optimization of the RSQ-8 to the S8*. *European Journal of Psychological Assessment*. 2016;32(2):155.
- [30] Pedregosa F, Varoquaux G, Gramfort A, Michel V, Thirion B, Grisel O, et al. Scikit-learn: Machine Learning in Python. *Journal of Machine Learning Research*. 2011;12:2825-30.
- [31] Gu S, Kelly B, Xiu D. Empirical asset pricing via machine learning. *National Bureau of Economic Research*; 2018.
- [32] Danka T, Horvath P. modAL: A modular active learning framework for Python. Available on arXiv at <https://arxiv.org/abs/1805.00979>. Available from: <https://github.com/cosmic-cortex/modAL>.
- [33] Kola I, Murukannaiah PK, Jonker CM, Van Riemsdijk MB. Towards Social Situation Awareness in Support Agents. *IEEE Intelligent Systems*. 2022.