



HAL
open science

PANDA: Human-in-the-Loop Anomaly Detection and Explanation

Grégory Smits, Marie-Jeanne Lesot, Véronne Yepmo Tchaghe, Olivier Pivert

► **To cite this version:**

Grégory Smits, Marie-Jeanne Lesot, Véronne Yepmo Tchaghe, Olivier Pivert. PANDA: Human-in-the-Loop Anomaly Detection and Explanation. IPMU 2022 - Information Processing and Management of Uncertainty in Knowledge-Based Systems, Jul 2022, Milan, Italy. hal-03696295

HAL Id: hal-03696295

<https://hal.inria.fr/hal-03696295>

Submitted on 15 Jun 2022

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.

PANDA: Human-in-the-Loop Anomaly Detection and Explanation

Grégory Smits¹,
Marie-Jeanne Lesot², Véronne Yepmo Tchaghe¹, and Olivier Pivert¹

¹University of Rennes – IRISA, UMR 6074, Lannion, France
{gregory.smits,veronne.yepmo,olivier.pivert}@irisa.fr

²Sorbonne Université, CNRS, LIP6, F-75005 Paris, France
marie-jeanne.lesot@lip6.fr

Abstract. The paper addresses the tasks of anomaly detection and explanation simultaneously, in the human-in-the-loop paradigm integrating the end-user expertise: it first proposes to exploit two complementary data representations to identify anomalies, namely the description induced by the raw features and the description induced by a user-defined vocabulary. These representations respectively lead to identify so-called data-driven and knowledge-driven anomalies. The paper then proposes to confront these two sets of instances so as to improve the detection step and to dispose of tools towards anomaly explanations. It distinguishes and discusses three cases, underlining how the two description spaces can benefit from one another, in terms of accuracy and interpretability.

Keywords: outlier detection, outlier explanation, XAI, human-in-the-loop, fuzzy vocabulary, linguistic description

1 Introduction

A common approach to provide users with eXplainable Artificial Intelligence (XAI) tools is to implement the human-in-the-loop paradigm, i.e. to offer the user a crucial role in the mining process itself. This paper considers the case of the anomaly detection task and proposes to take into account user knowledge expressed in the form of a fuzzy vocabulary to describe linguistically the data.

Informally, anomaly or outlier detection aims at identifying, in a data set, instances that are conspicuous and, as put in the commonly accepted definition, “deviate so much from other observations so as to arouse suspicions that they were generated by a different mechanism” [10]. There exist numerous methods to perform this task, as well as multiple surveys and taxonomies, see e.g. [4, 15]. However, most of them address the issue as a machine learning task, without taking into account the user who analyses the data. Recently, many methods have been proposed to provide *a posteriori* explanations about the identified outliers within the XAI framework, see e.g. [20] for a survey.

This paper proposes to take the user into account very early in the outlier detection process, leading to a knowledge-driven method that offers as additional

feature an integrated linguistic description of the identified points. It thus opens the way to their interpretation and understanding by the user, i.e. to an outlier explanation method.

In order to do so, the proposed approach called PANDA, that stands for Personalised ANomaly Detection and Analysis, takes as input, in addition to the data set to be processed, a fuzzy vocabulary defined by the user: this vocabulary allows building linguistic descriptions of the data and constitutes precious user knowledge. For instance, the vocabulary defines indistinguishable areas in the data, i.e. values that should be considered as equivalent although they numerically differ: it can lead to distance functions more relevant from the user point of view than the classical Euclidean distance [8].

The PANDA method proposes to exploit such a vocabulary to dispose of a second data representation, complementary to the description induced by the basic data features: it is built as the vector concatenating the membership degrees to all modalities of all features. It thus defines a knowledge-driven representation of the data. In addition to providing a formalization of a subjective interpretation of the data, this vector also provides a normalization (in the unit interval) and a unification of non commensurable values, easing the combination of numerical and categorical attributes within a data mining task.

PANDA then proposes to apply an outlier detection method in these two description spaces. This principle bears similarity with the method proposed in [11] that applies a clustering algorithm in the two data representation spaces: the initial data definition space and the symbolic space induced by the vocabulary. However the aim in [11] is to quantify the adequacy between the vocabulary and the data inner structure. The crucial analysis step of PANDA confronts the two sets of anomalies, identified in the two spaces, so as both to improve the detection step and to dispose of tools towards anomaly explanations: PANDA makes it possible to extend any anomaly detection method with a vocabulary-based personalization of the data and a cross analysis of the outliers detected in the two spaces. The isolation forest method [12] to anomaly detection is used as an illustration in this paper.

The paper is structured as follows: Section 2 summarises related works, both on anomaly detection and explanation, Section 3 describes the proposed PANDA method, illustrating it with synthetic data and Section 4 presents a case study on real data describing car ads. Section 5 concludes the paper.

2 Related Works

This section briefly presents the two tasks to which the proposed PANDA method relates, considering anomaly detection and explanation in turn.

2.1 Anomaly Detection

There exist numerous methods to detect anomalies, i.e. points that deviate from so-called regular phenomena, as well as multiple surveys, see e.g. [4, 7, 20, 15].

Beyond a distinction between supervised and unsupervised approaches, there exists no consensus about a taxonomy or the categories, nor the number of categories, further structuring the domain. It has for instance been proposed to distinguish between approaches based on nearest neighbours, clustering, statistics, subspaces and classifiers [7], or between approaches based on density, distance and models [12], or between approaches based on distance, model and neural networks [20]. To name four, some classical examples of anomaly detection algorithm include LOF (Local Outlying Factor [3] and its numerous variants), Isolation Forests [12], One-class SVM [1] and Auto-Encoder Ensemble [5].

For the implementation of the generic PANDA method described in this paper, the isolation forest (IF) approach [12] is considered. It constitutes an unsupervised ensemble-based method that combines multiple isolation trees. Such an isolation tree recursively draws random features and values to partition the data, until a predefined tree depth is reached or each leaf contains only individual (or indistinguishable) data points. Based on the fact that, by definition, outliers are distant from dense regions, they are likely to be isolated early by the recursively defined node partition. They thus appear in leaves close to the tree root: an isolation score of any data point is defined as the length of the path to the leaf it is assigned to. An isolation forest then combines most often hundreds of such randomly built trees and, for any data point, outputs an anomaly score based on its average isolation score [12].

2.2 Anomaly Explanation

Given a set of identified outliers, a natural question from the user is to ask for the reason why they are considered as such, i.e. what makes them abnormal: this calls for anomaly explanation methods, at the cross-roads of anomaly detection and XAI. According to the recent survey [20], four categories of such methods can be distinguished, depending on the type of provided explanations. The first one, also the most represented one, groups feature importance approaches, that either compute a score for each individual feature, as [14] for instance, or determine relevant subspaces, as e.g. [13]. These approaches can also be distinguished depending on whether they apply locally to single outlier points or globally to sets of outliers, or whether they are detector specific or agnostic, within the so-called outlier aspect mining task [6].

A second category of anomaly explanation methods groups approaches that additionally associate the responsible features with the values they take, as for instance [2]. The latter can for instance be identified by rules expressed as conjunction of predicates, where explanations take a disjunctive normal form. A third category groups approaches based on point comparisons, that underline the difference between an outlying point and regular points, for instance in looking for counterfactual examples [9]. The fourth category focuses on analysing the structure of the data, identifying the relations between subsets, i.e. clusters, of regular points and individual anomalies or sets of anomalies, as e.g. [17]

To the best of our knowledge, none of these methods take into account user knowledge, so as to provide personalised and more understandable explanations.

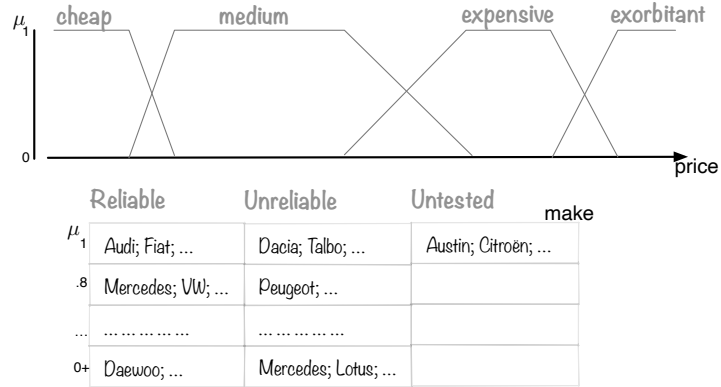


Fig. 1. Examples of fuzzy partitions describing car prices (top) and makes (bottom)

3 Proposed PANDA Approach

After presenting the notations used in the paper, this section describes the two steps of the proposed PANDA approach, respectively corresponding to the machine learning process applied to the considered data set and to the contrastive analysis of the identified anomalies.

3.1 Notations and Illustrative Data Set

$\mathcal{D} = \{x_1, x_2, \dots, x_n\}$ denotes a set of n data points described by m attributes, A_1 to A_m , with respective domains D_1 to D_m . These attributes can be numerical or categorical.

$\mathcal{V} = \{P_1, \dots, P_m\}$ denotes a vocabulary defined as a set of linguistic variables: for $i = 1..m$, P_i is a triple $\langle A_i, \{\mu_i\}, \{l_i\} \rangle$ with q_i modalities. The μ_{ij} , $j = 1..q_i$ are the respective membership functions of the modalities defined on universe D_i and the l_{ij} their respective linguistic labels.

Figure 1 depicts two examples of fuzzy partitions: the top part applies to a numerical attribute describing second hand car prices, for which $q = 4$ and with labels ‘cheap’, ‘medium’, ‘expensive’ and ‘exorbitant’. The bottom part applies to a categorical attribute describing the car make: it shows a subjective interpretation of their reliability, with $q = 3$ and labels ‘reliable’, ‘unreliable’ and ‘untested’. The membership functions are defined through their α -cuts: for each term, each row shows the makes whose membership degrees equal the value given on the left of the table.

It is assumed that each P_i defines a strong partition [16], i.e. $\forall y \in D_i$, $\sum_{j=1}^{q_j} \mu_{ij}(y) = 1$. In addition, it is assumed that the partition is such that any value y can partially satisfy up to two modalities only. In the case of features with numerical domains, these two modalities are adjacent.

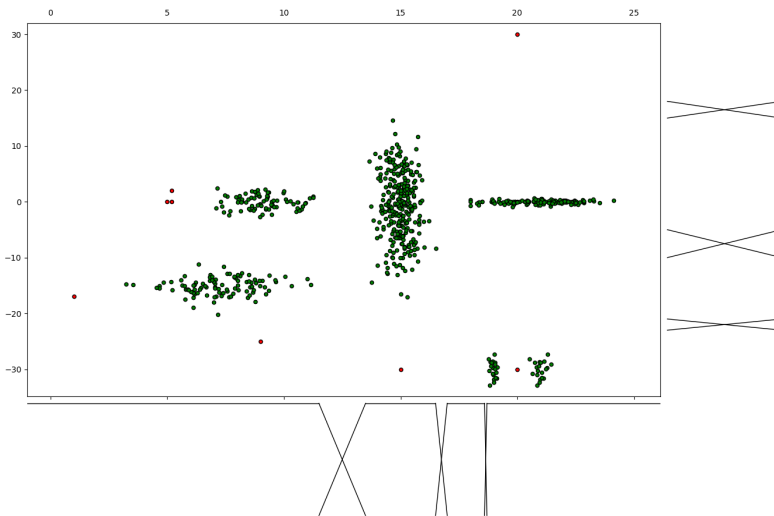


Fig. 2. Considered 2D illustrative data set

Throughout the section, we consider as example the data shown on Figure 2 with $n = 651$, $m = 2$ with A_1 and A_2 numerical attributes whose respective domains are $D_1 = [0, 25]$ and $D_2 = [-40, 30]$. P_1 contains $q_1 = 4$ modalities whose membership functions are shown below the graph and P_2 contains $q_2 = 4$ modalities as well, whose membership functions are shown on the right side of the graph. The data set contains several dense regions as well as some outliers.

3.2 Data Processing

Data Rewriting with the Fuzzy Vocabulary Each data point is first rewritten by computing its membership degrees to all modalities of all attributes and concatenating them: $x = \langle x^1, \dots, x^m \rangle$ is represented as the vector of $Q = \sum_{j=1}^m q_j$ components:

$$\langle \mu_{v_{11}}(x^1), \dots, \mu_{v_{1q_1}}(x^1), \dots, \mu_{v_{m1}}(x^m), \dots, \mu_{v_{mq_m}}(x^m) \rangle.$$

This vector is sparse, having at most $2m$ non-zero components due to the hypotheses on the partitions described in the previous section.

The whole dataset \mathcal{D} may thus be rewritten according to a vocabulary \mathcal{V} in linear time wrt. $|\mathcal{D}|$ but this process may easily be distributed to handle massive data [19]. The rewritten data $\mathcal{D}^\mathcal{V}$ are thus described as vectors of $[0, 1]^Q$.

Double Anomaly Detection To leverage the expert knowledge about the data embedded in his/her vocabulary, a same anomaly detection method is applied on

both \mathcal{D} and $\mathcal{D}^{\mathcal{V}}$. In this paper, the Isolation Forest (IF) [12] method is applied in the two spaces using recommended parameters (100 trees in the forest, anomaly score with threshold 0.5 and a subset minimum size 256).

The resulting sets of identified anomalies are denoted by \mathcal{A} and $\mathcal{A}^{\mathcal{V}}$ respectively. The former, identified in the initial feature space, are interpreted as data-driven anomalies; the latter, identified in the description space induced by the user vocabulary, are interpreted as knowledge-driven anomalies.

3.3 Anomaly Analysis: Cross Comparison of Detected Anomalies

The anomaly analysis step then consists in comparing the two sets \mathcal{A} and $\mathcal{A}^{\mathcal{V}}$, considering their intersection and differences, commented in turn in this section. The goal of this comparison is to help users better understand both the data and the vocabulary. It is shown that it makes it possible to refine the anomaly detection, turning data-driven anomalies into contextual regularities and conversely points looking regular in \mathcal{D} into contextual anomalies. Tools are thus provided towards the explanation of the identified outliers, as discussed below.

Figure 3 shows the result of an IF anomaly detection on \mathcal{D} (top part) and $\mathcal{D}^{\mathcal{V}}$ (bottom part). The blueish zones indicate the anomaly scores for each point of the domain, white zones corresponding to high anomaly scores. Black lines in the top part of Figure 3 are the separation lines of one isolation tree randomly drawn from the forest.

It can first be observed that, as expected, the general profiles of the anomaly score landscapes differ between the two graphs. In particular, in the rewritten case (bottom part), the regions homogeneous in terms of scores are parallel to the axes: the modalities of the fuzzy variables define indistinguishability zones within which all points have the same representation and are thus treated the same way. As a consequence, the data density is aggregated within each region, with fuzzy boundaries between the Cartesian product of the fuzzy set cores. The anomaly score landscape in the case of the initial representation space obviously follows the observed data density more closely.

Linguistic Description of Anomalies: $\mathcal{A} \cap \mathcal{A}^{\mathcal{V}}$ A first category of anomalies contains the points that are identified as such in both description spaces, i.e. the intersection of the two anomaly sets. These points can be considered as confirmed anomalies, for which in addition a linguistic description is available.

Indeed, a point $x \in \mathcal{A} \cap \mathcal{A}^{\mathcal{V}}$ is a data-driven outlier whose description in the vocabulary-induced space is considered as anomalous as well. Furthermore, this vocabulary-induced description characterises x , as it allows identifying it as an anomaly, and provides a linguistic description.

This case is illustrated with points x_1 with coordinates (1, -17), x_2 (20, 30), x_3 (9, -25) and x_4 (15, -30), on Figure 3: they are indeed outliers from the data density or separability point of view and from the vocabulary point of view. They illustrate two distinct cases: x_1 and x_2 possess extreme feature values, that make them outliers, whereas x_3 and x_4 possess anomalous features combinations as

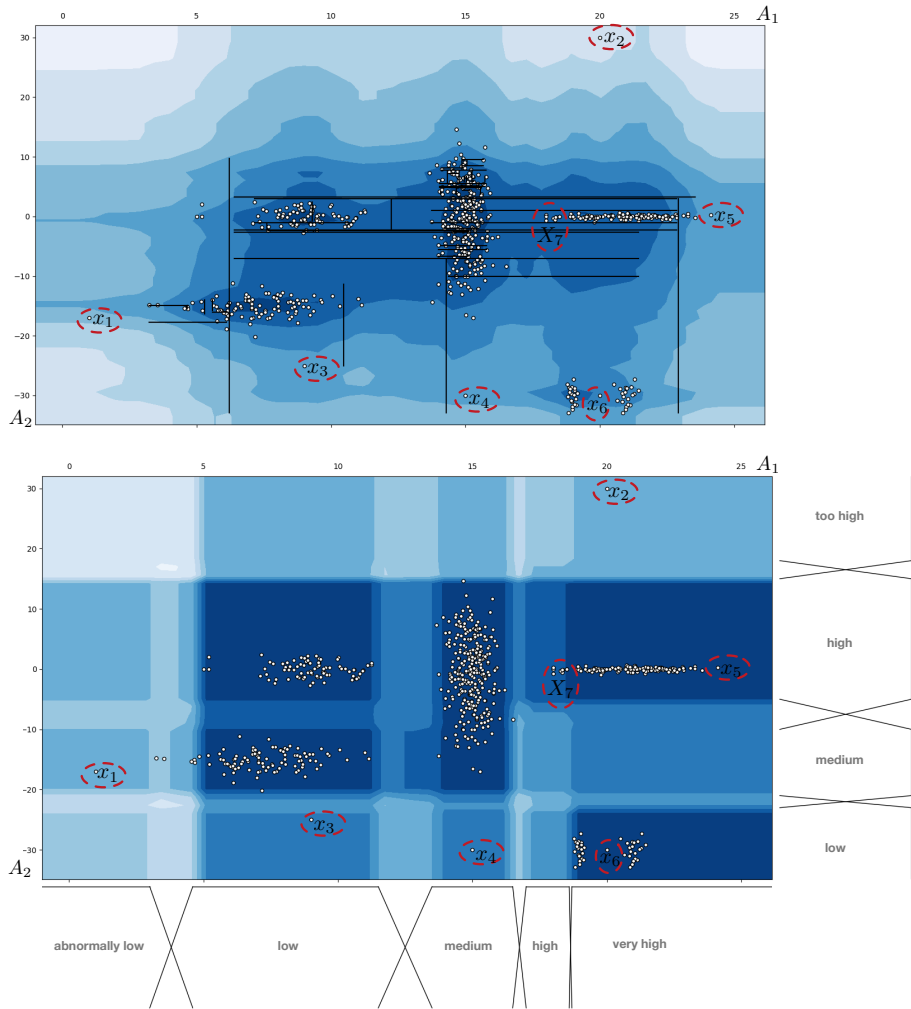


Fig. 3. Profile of the obtained anomaly scores: the lighter the colour, the higher the score. (Top) for \mathcal{D} , (bottom) for \mathcal{D}^V .

compared to the other data points. In addition, the vocabulary allows describing linguistically these outliers: x_1 can e.g. be described as “*value on attribute A_1 is abnormally low and value on attribute A_2 is medium*” whereas x_3 as “*value on attribute A_1 is low and on attribute A_2 is low*”.

The obtained linguistic description can be considered as a step towards an intelligible explanation of the anomalies. However, two caveats require for caution. First, the provided description involves all features, i.e. is of size m . As such it may be the case that it is actually not legible when the number of features is high. Enriching the proposed approach with outlier aspect mining (see Section 2.2) is a considered extension to address this issue. Second, the description does not explain the reason why the considered point is an anomaly: for instance this description takes the same form for the two above-mentioned illustrative points, whereas they correspond to different cases. This corresponds to a classic challenge of outlier explanation generation.

Unexpected Anomalies: $\mathcal{A} \setminus \mathcal{A}^{\mathcal{V}}$ A second category of anomalies contains data points that are identified as outliers in the initial description space, but not in the vocabulary-induced one. This case is illustrated with points x_5 (24, 0) and x_6 (20, -30) on Figure 3: in $\mathcal{D}^{\mathcal{V}}$, they are associated with the minimal anomaly score, i.e. they are considered as regular points. Indeed, they are described with terms that make them unanomalous whereas they are isolated from the data density point of view.

Such points can be described as “unexpected anomalies” insofar as the user does not dispose of a vocabulary that allows describing them and thus seems not to expect them. In an interactive information extraction process, it is highly relevant to draw his/her attention to these anomalies, so that the reason why they are not identified as such in the knowledge-driven approach can be explored. Several cases can indeed be distinguished, calling for different treatments.

A first possibility is that the user vocabulary is actually not adequate, i.e. does not correspond to the data distribution and content: it is useful to underline the existence of such specific cases the user may not have envisioned, suggesting to add new linguistic modalities to describe such subspaces specifically. This corresponds to a case of vocabulary data adequacy that is not captured by previous works on this topic, as e.g. [11]. For the considered illustrative data set, x_5 is an example of this case: it may call for splitting the *very high* modality so as to dispose of a term for this specific value of attribute A_1 .

On the other hand, a second possibility is that these data point should indeed not be identified as anomalies: the knowledge provided through the vocabulary allow to diagnose issues in the data, suggesting the need to add information so that they are not considered as anomalies. It may for instance be the case that the processed data set is actually incomplete and misses data points, that would e.g. connect the candidate anomaly to a denser data region: the data set may be not representative of the underlying data distribution, about which the user vocabulary provides information. This can again be illustrated by data point x_5 , which may be connected to the cluster of regular data observed for lower values

of attribute A_1 . Similarly, the vocabulary can be considered as suggesting that data point x_6 should not be considered as isolated, further suggesting that there should be no distinction between the two clusters it is inbetween.

Inadequate Vocabulary: $\mathcal{A}^{\mathcal{V}} \setminus \mathcal{A}$ A third category of anomalies contains, as a reciprocal of the second category, data points that are identified as outliers in the vocabulary-induced description space, but not in the initial one. This case is illustrated with the set of points X_7 , around $(17, 0)$ on Figure 3: they build a minor group of points described with the *high* modality of attribute A_1 .

Such points can be characterised as special cases based on the vocabulary, whereas they are not in the raw feature space. They may indicate a type of vocabulary inadequacy, different from the one discussed above: the vocabulary can be interpreted as being too subtle and introducing fine distinctions that are not justified in a data-driven analysis. As illustrated with Figure 3, such cases e.g. occur when a modality splits a dense data area, here with the distinction between modalities *high* and *very high* of attribute A_1 . They may suggest the need for vocabulary revision, in the same manner as the one explored in [18].

On the other hand, the fuzzy vocabulary is a model of the knowledge an expert possesses about a specific applicative context, explaining how subsets of the different attribute domains have to be interpreted [8]. Thus, these points in $\mathcal{A}^{\mathcal{V}}$ not identified as anomalies in \mathcal{D} could also correspond to contextual false negative. However, the identification of such cases relies on additional contextual knowledge: the user may be interested in detecting the occurrence of such cases and the fine distinction may be required from an expert point of view. The vocabulary then offers the mean to identify them. As a concrete example of such a situation, let us consider the temperature monitoring a combustion engine whose ideal temperature is around 90°C . Whereas observing operating temperatures in the range $[60, 91]$ may not be problematic (it may e.g. be records during the warm-up phase), it may be crucial for the expert to know when the temperature reaches 92°C . A dedicated vocabulary is a solution to avoid contextual false negatives and false positives.

4 Use case: Secondhand Car Ads

This section presents preliminary experiments conducted on a real data set describing classified ads about secondhand cars and discusses the results obtained when applying the proposed PANDA method. Identifying anomalies then aims at detecting both possible description errors, e.g. typing errors that make the ads unrealistic, and very specific cars, e.g. vintage cars or rare models.

4.1 Experimental Protocol

The considered real data set contains 49,188 ads about secondhand cars described by six attributes *price*, *mileage*, *year*, *priceNew*, *make* and *model*. The *priceNew* attribute indicates the price of the car of the considered make and

Table 1. Terms of the vocabulary used to rewrite the car descriptions

Attr.	Linguistic values
Price	almostOffered, veryLow, low, medium, expensive, veryExpensive, exorbitant
Mileage	almostNull, veryLow, low, medium, high, veryHigh, huge
Year	vintage, old, acceptable, recent, almostNew
PriceNew	veryLow, low, medium, expensive, veryExpensive, exorbitant
Make	luxury, highClass, mediumClass, lowClass

model when sold new. Table 1 gives the labels of the linguistic variables defined for the five first attributes; the associated membership functions, omitted for size constraints, correspond to common sense definitions of the modelled properties.

The proposed PANDA method is applied on this data set \mathcal{D} and its rewritten form $\mathcal{D}^\mathcal{V}$. The Isolation Forest algorithm is run with the hyper-parameter values suggested in [12]: 100 trees are built on random subsamples of the data set each containing 256 data points, the anomaly threshold is set to 0.5.

4.2 Result Analysis

Tables 2 and 3 show the ten instances, respectively in \mathcal{D} and $\mathcal{D}^\mathcal{V}$, with the highest anomaly scores. Deviating values and value combinations are shown in bold, based on manual analysis. It can be observed that the first PANDA category, $\mathcal{A} \cap \mathcal{A}^\mathcal{V}$, is empty, this section discusses the reason why and comments the two other categories in turn, comparing \mathcal{A} and $\mathcal{A}^\mathcal{V}$.

Regarding \mathcal{A} given in Table 2, it can first be observed the ads ranked 1, 4, 5 and 10 can legitimately be considered as anomalies due to their erroneous prices, that take values greater than one million. The analysis of the other ads in this list shows they can be interpreted as anomalies because they correspond to rare luxury sport cars, that despite not being new models are still very expensive even with a medium mileage (see e.g. the third ad).

Observing the results for the rewritten data in Table 3, one can first remark that the integration of expert knowledge using the fuzzy vocabulary leads to a very different list of anomalies. Indeed, due to the fact that luxury makes are now grouped within a dedicated modality, they do not appear anymore as anomalies: *luxury* make having a *very expensive* price despite an *acceptable* year is now a sufficiently frequent conjunction of properties describing a subset of the analyzed ads. As a consequence, the knowledge driven anomaly detection makes it possible to identify other outliers and its combination with the data driven approach to get a better understanding of their respective contents.

Anomalies $\mathcal{A}^\mathcal{V}$ can be interpreted as being of two types: typing errors leading to unrealistic values, as a mileage equal to 1 for a *vintage* car (e.g. ad 3), and suspicious combinations of properties. Ads 5, 8 and 9 are examples of the latter: they correspond to cars from a luxury make with an expensive or very expensive price and a very low sale price. Looking more in depth at the ad description reveals that these cars are sold with a broken engine.

Table 2. Top-10 anomalies found in the secondhand cars dataset, \mathcal{A}

	Price	Mileage	Year	PriceNew	Make	Model	Score
1	7,500,000	112,000	1993	98,754	mercedes	500 SL A	0.705
2	110,000	15,000	1984	80,570	ferrari	BB 512 5	0.696
3	62,000	50,000	1992	168,174	ferrari	F 512 4.9i	0.69
4	42,600,000	22,000	2010	44,020	mercedes	Classe C 350 CDI	0.688
5	17,490,000	202,000	2005	54,440	mercedes	Classe CLS 320 CDI	0.682
6	109,000	3,800	2007	104,719	porsche	911 3.6i	0.681
7	93,900	41,900	2007	168,372	ferrari	F430 Spider V8	0.68
8	112,000	21,750	2009	136,882	audi	R8 V10 5.2 FSI 525	0.68
9	115,000	22,154	2009	136,882	audi	R8 V10 5.2 FSI 525	0.68
10	12,500,000	334,000	2007	32,774	mercedes	Classe C 220 CDI	0.677

Table 3. Top-10 anomalies found in the rewritten secondhand cars data set, \mathcal{A}^V

	Price	Mileage	Year	PriceNew	Make	Model	Score
1	450	100	1988	8,232	renault	Super 5 Tiga	0.609
2	850	229,000	1983	8,345	bmw	315	0.604
3	25	1	2010	26,798	audi	A3 Sportback 2.0 TDI	0.602
4	2,350	4,801	2009	9,639	dacia	Sandero 1.5 dCi 70	0.599
5	1,000	450,000	1988	36,550	mercedes	300 TD	0.598
6	6,999	159	2004	30,387	jaguar	X	0.597
7	700	10	1994	28,178	bmw	525 TD	0.597
8	1,000	500,000	1985	36,416	bmw	628 CSi	0.596
9	2,990	290,000	1988	76,441	bmw	750 iL	0.596
10	500	320	1991	27,116	bmw	524 TD	0.594

5 Conclusion and Perspectives

Addressing the task of identifying and explaining outliers in a data set, the PANDA approach proposed in this paper makes it possible to integrate user expertise so as to detect and compare both data-driven and knowledge-driven anomalies. Analyses based on an illustrative toy data set and a real data set show how they enrich each other: the PANDA approach provides a personalised outlier detection method, drawing the user attention to different types of specific cases of interest. It thus constitutes a human-in-the-loop outlier detection and offers tools towards outlier explanation.

Future works will aim at including further developments regarding the outlier explanation component, in particular the generation of linguistic description of the identified anomalies, e.g. combining the proposed methodology with outlier aspect mining components. They will also address the question of integrating PANDA within relational data base management systems, as exploratory tool for a user to get a global view on the data content and global structure. Experiments with real data and real users will be conducted to measure the extent to which it contributes to the user understanding and satisfaction when interacting with massive data sets.

References

1. Amer, M., Goldstein, M., Abdennadher, S.: Enhancing one-class support vector machines for unsupervised anomaly detection. In: Proc. of the ACM SIGKDD Workshop on Outlier Detection and Description. pp. 8–15 (2013)
2. Barbado, A., Corcho, O., Benjamins, R.: Rule extraction in unsupervised anomaly detection for model explainability: Application to OneClass SVM. *Expert Systems with Applications* 189 (2022)
3. Breunig, M.M., Kriegel, H.P., Ng, R.T., Sander, J.: LOF: identifying density-based local outliers. *ACM sigmod record* 29(2), 94–104 (2000)
4. Chandola, V., Banerjee, A., Kumar, V.: Anomaly detection: A survey. *ACM computing surveys (CSUR)* 41(3), 1–58 (2009)
5. Chen, J., Sathe, S., Aggarwal, D., Turaga, D.: Outlier detection with autoencoder ensembles. In: Proc. of the SIAM Int. Conf. on Data Mining. pp. 90–98 (2017)
6. Duan, L., Tang, G., Pei, J., Bailey, J., Campbell, A., Tang, C.: Mining outlying aspects on numeric data. *Data Mining and Knowledge Discovery* 29, 116–1151 (2014)
7. Goldstein, M., Ushida, S.: A comparative evaluation of unsupervised anomaly detection algorithms for multivariate data. *PLoS One* 11(4) (2016)
8. Guillaume, S., Charnomordic, B., Loisel, P.: Fuzzy partitions: a way to integrate expert knowledge into distance calculations. *Information sciences* 245, 76–95 (2013)
9. Haldar, S., Johnand, P.G., Saha, D.: Reliable counterfactual explanations for autoencoder based anomalies. In: Proc. of the 8th ACM IKDD CODS and 26th COMAD Conf. pp. 83–91. ACM (2021)
10. Hawkins, D.M.: Identification of outliers, vol. 11. Springer (1980)
11. Lesot, M.J., Smits, G., Pivert, O.: Adequacy of a user-defined vocabulary to the data structure. In: Proc. of the IEEE Int. Conf. on Fuzzy Systems. IEEE (2013)
12. Liu, F.T., Ting, K.M., Zhou, Z.H.: Isolation-based anomaly detection. *ACM Transactions on Knowledge Discovery from Data (TKDD)* 6(1), 3 (2012)
13. Myrtakis, N., Tsamardinos, I., Christophides, V.: Proteus: Predictive explanation of anomalies. In: Proc. of the 37th IEEE Int. Conf. on Data Engineering (ICDE). pp. 1967–1972. IEEE (2021)
14. Pevný, T.: LODA: Lightweight on-line detector of anomaly. *Machine Learning* 102, 275–304 (2015)
15. Ruff, L., Kauffmann, J., Vandermeulen, R., Montavon, G., Samek, W., Kloft, M., Dietterich, T., Müller, K.R.: A unifying review of deep and shallow anomaly detection. *Proc. of the IEEE* 109(5), 756–795 (2021)
16. Ruspini, E.H.: A new approach to clustering. *Information and Control* 15(1), 22 – 32 (1969)
17. Shukla, A.K., Smits, G., Pivert, O., Lesot, M.J.: Explaining data regularities and anomalies. In: Proc. of the Int. Conf. on Fuzzy Systems. IEEE (2020)
18. Smits, G., Pivert, O., Lesot, M.J.: A vocabulary revision method based on modality splitting. In: Proc. of the Int. Conf. on Information Processing and Management of Uncertainty in Knowledge-Based Systems, IPMU. CCIS, vol. 442, pp. 376–385. Springer (2014)
19. Smits, G., Pivert, O., Yager, R.R., Nerzic, P.: A soft computing approach to big data summarization. *Fuzzy Sets and Systems* 348, 4–20 (2018)
20. Tchaghe, V.Y., Smits, G., Pivert, O.: Anomaly explanation: A review. *Data & Knowledge Engineering* 137 (2021)