# Subjective data arrangement using clustering techniques for training expert systems

Isaac Martín de Diego [a],*, Oscar S. Siordia [b], Alberto Fernández-Isabel [a], Cristina Conde [a], Enrique Cabello [a]

[a] *Face Recognition and Artificial Vision Group, Data Science Laboratory, Rey Juan Carlos University, c/ Tulipán, s/n, 28933, Móstoles, Spain*
[b] *Centro de Investigación en Ciencias de Información Geoespacial (CentroGeo), Laboratorio Nacional de Inteligencia (GeoInt), Parque Científico Tecnológico Yucatán (PCTY), México*

## ABSTRACT

The evaluation of subjective data is a very demanding task. The classification of the information gathered from human evaluators and the possible high noise levels introduced are ones of the most difficult issues to deal with. This situation leads to adopt individuals who can be considered as experts in the specific application domain. Thus, the development of Expert Systems (ES) that consider the opinion of these individuals have been appeared to mitigate the problem. In this work an original methodology for the selection of subjective sequential data for the training of ES is presented. The system is based on the arrangement of knowledge acquired from a group of human experts. An original similarity measure between the subjective evaluations is proposed. Homogeneous groups of experts are produced using this similarity through a clustering algorithm. The methodology was applied to a practical case of the Intelligent Transportation Systems (ITS) domain for the training of ES for driving risk prediction. The results confirm the relevance of selecting homogeneous information (grouping similar opinions) when generating a ground truth (a reliable signal) for the training of ES. Further, the results show the need of considering subjective sequential data when working with phenomena where a set of rules could not be easily learned from human experts, such as risk assessment.

© 2018 Elsevier Ltd. All rights reserved.

## 1. Introduction

The practice of Knowledge Engineering (Van Do, Le Thi, & Nguyen, 2018) has become a very useful approach to solve complex problems that require a high level of human expertise. This discipline involves integrating knowledge into computer systems which emulates the decision-making ability of a human expert in a specific domain. The systems in charge of achieving these tasks are the Expert Systems (ES) (Agarwal & Goel, 2014).

The building, maintaining and development of ES (Djamal et al., 2017) are mainly based on the interaction between the knowledge engineer and the domain expert (Yau & Sattar, 1994). The development of a reliable ES requires a deep understanding and a good representation of the knowledge of the domain expert.

In most of the cases, the knowledge representation is based on a set of rules (a production system) that ease the explanation of the decision-making made by the inference engine (Wick & Slagle, 1989). These rules are build from the knowledge acquired from human experts with the application of Machine Learning techniques (such as Neural Networks (Lin & Zhang, 2012), Deep Learning (Wei, He, Chen, Zhou, & Tang, 2017), Decision Trees (Sriram & Yuan, 2012), Fuzzy Logic (Wang, Lee, & Ho, 2007), Bayesian methods (WenBin, XiaoLing, YiJun, & Yu, 2010), Genetic Algorithms (Daza et al., 2011), among others).

Knowledge acquisition is a process which aims to extract knowledge, experience and problem-solving procedures from one or more domain experts. Several techniques have been proposed for a correct knowledge acquisition (see Hua, 2008 for a complete review). Nevertheless, there are several problems that must be considered when acquiring knowledge from human experts (Gaines, 1987):

- Experts may not be able to express their knowledge in a structured way.

- Experts may not be aware of the significance of the knowledge they have used.
- The expressed knowledge may be irrelevant, incomplete or not understandable.

In some cases, depending on the field of application, it may be easier to extract the knowledge from human experts through a continuous scale. This is the case of the risk assessment, where the knowledge could be acquired in a predefined scale (e.g. from 0, no risk, to 100, maximum risk). Here, the knowledge of the experts is gathered in form of subjective sequential data (Prelec, 2004) and could be treated as time series for its study and integration (see, for instance, de Diego, Crespo, Siordia, Conde, & Cabello, 2011; de Diego, Siordia, Conde, & Cabello, 2011; Siordia, de Diego, Conde, & Cabello, 2011a).

However, the integration of several opinions into a unique ground truth (i.e. a reliable signal) is a hard-to-achieve task (Liou & Nunamaker, 1990). Two different scenarios appear. The consideration of knowledge from too few experts could provide a ground truth with insufficient information. In contrast, the consideration of knowledge from too many experts could generate a noisy ground truth due to the appearance of possible contradictions between their evaluations (Turban, 1991). Different statistical approaches have been proposed in the past (see, for instance, meta-analysis methods in Brockwell & Gordon (2001)).

In this paper, it is presented a novel methodology for the selection of subjective sequential data for the training of ES. This methodology upgrades the previous approaches in the domain (Siordia, de Diego, Conde, & Cabello, 2014) focusing on the inclusion of more experts. This increment of sources of information leads to produce heterogeneous and noisy evaluations that have to be arranged. A novel definition of similarity between experts' evaluations will be firstly presented here. In addition, in the previous method, the agreement between two or more evaluations was enough to define a unique ground truth. However, in the present paper, all the homogeneous evaluations will be used.

Delving into the main idea behind, the methodology consists of the arrangement of a set of evaluations acquired from human experts through a hierarchical clustering technique. In this way, similarities between the evaluations of experts could be identified and grouped together, filtering the contradictions. The resulting groups (clusters) could be analyzed in order to select the most appropriate ground truth labels (Healey, 2011) for the training of the ES.

The proposed methodology is a general purpose approach. Thus, it can be used in several domains where different human opinions should be managed. In this paper, the methodology is applied to a practical case on the Intelligent Transportation Systems (ITS) domain (Alam, Ferreira, & Fonseca, 2016). It is focused on the characterization of risky or safe situations for the driving task.

Regarding the experiments, three different have been considered to illustrate the performance of the approach. First, an experiment has been developed using synthetic data for demonstrative purposes. The other experiments are based on the practical case presented above. Thus, they have been achieved using real driving risk evaluations made by experts from urban and interurban scenarios respectively.

The paper is organized as follows: Section 2 situates the approach in the domain. Section 3 introduces the proposed methodology, explaining in detail the similarity measures to evaluate subjective sequential data. Section 4 describes the practical case where the approach has been applied. Section 5 presents the achieved experiments and their most relevant results. Finally, Section 6 concludes and provides future lines of work.

## 2. Related work

The ES have been widely used for multiple purposes (Wagner, 2017). They are systems that are able to exhibit features associated with human intelligence (e.g. problem solving or reasoning) (Hodson, 2018). They have a common architecture based on two main modules: a domain dependent knowledge database and the inference mechanism. Examples of them are Attwell, Leask, Meyer, Rokkas, and Ward (2017) or Meza-Palacios et al. (2017).

The architecture of the ES presented here comprehends both modules. The knowledge base is acquired from traffic experts that evaluate the behavior of drivers, while the inference mechanism is built applying similarity measures and unsupervised learning techniques.

Delving into these unsupervised learning techniques, clustering (see, for instance, Aggarwal, 2015) is an initial and fundamental step in data analysis. It has as a main goal to reveal a natural partition of data into a number of meaningful subclasses or clusters. Clustering of sequential data differs from clustering of static feature data mainly in how to compute the similarity between two data objects.

In the presented approach, *Agnes* clustering algorithm has been selected. It is an agglomerative hierarchical clustering technique that provides real-time updating (see Kaufman & Rousseeuw, 2009 for a complete description).

Regarding the characteristics of subjective sequential data (where sudden changes occur and where the key information is given by its trend), it is appropriated a piecewise representation of the data. Thus, a variety of algorithms to obtain a proper linear representation of sequential data have been proposed in the literature (see, for instance, Keogh, Chu, Hart, & Pazzani, 2004; Lachaud, Vialard, & De Vieilleville, 2005; Zhu, Wu, & Li, 2007)

Focusing on driving risk situations, there are multiple examples of their characterization through the analysis of data collected on driving sessions. These approaches are usually focused on the study of the drivers behavior and how their acts affect to the driving task. For instance, Cheng, Park, and Trivedi (2007) introduces an approach based on multi-perspective (several cameras recording the driver) in order to analyze the different body movements (mainly head and hands). In the case of Malta, Miyajima, Kitaoka, and Takeda (2011), it is oriented to identify the frustration and the different emotions of the driver and how these emotions affect to the driving task. These systems are related to the approach presented in this paper. Both examples use cameras to identify the movements of the driver, though in our case the face expressions are not considered.

Other studies have their key topic in learning from specific risk situations identifying patterns. For example, Wang, Zhu, and Gong (2010) has as a main purpose to infer the safe or dangerous actions achieved by drivers using time series and unsupervised learning. In this case, the presented approach could be considered as one of this type of systems.

There are similar approaches that evaluate specific tasks of the driver and not only the hands or the facial expressions. The pressures exerted on the break and throttle pedals are also interesting parameters to evaluate. Examples of these are Sathyanarayana, Boyraz, and Hansen (2008), that is oriented to route paths recognition and Rakha, El-Shawarby, and Setti (2007), which addresses the behavior of driver in intersections.

Delving into the behavior of drivers, multiple theoretical models have been developed. They can be classified into: taxonomic models and functional models. The firsts usually produce descriptive classifications of certain elements of traffic based on a context. They can be decomposed into features-based models (Bone & Mowen, 2006) and task-analysis models (Fastenmeier & Gstalter, 2007). The second ones can be organized into mechanical