

A multi-agent architecture for labeling data and generating prediction models in the field of social services

Emilio Serrano, Pedro del Pozo-Jiménez, Mari Carmen Suárez-Figueroa,
Jacinto González-Pachón, Javier Bajo, Asunción Gómez-Pérez
emilioserra@fi.upm.es

Ontology Engineering Group, Universidad Politécnica de Madrid, Spain

Abstract. Prediction models are widely used in insurance companies and health services. Even when 120 million people are at risk of suffering poverty or social exclusion in the EU, this kind of models are surprisingly unusual in the field of social services. A fundamental reason for this gap is the difficulty in labeling and annotating social services data. Conditions such as social exclusion require a case-by-case debate. This paper presents a multi-agent architecture that combines semantic web technologies, exploratory data analysis techniques, and supervised machine learning methods. The architecture offers a holistic view of the main challenges involved in labeling data and generating prediction models for social services. Moreover, the proposal discusses to what extent these tasks may be automated by intelligent agents.

Keywords: Multi-agent systems, Human-agent societies, Social services, Machine learning.

1 Introduction

Building artificial intelligence and machine learning based systems has never been easier than today thanks to: (1) open-source tools such as TensorFlow or Spark; and, (2) massive amounts of computation power through cloud providers such as Amazon Web Services and Google Cloud [4]. Machine learning prediction models are widely employed, among others, in health services. These systems have a deep impact because there is strong evidence supporting that early detection of medical conditions results in less severe outcomes. For instance, the reader may calculate the risk of suffering a heart disease at different webs [3].

Social services, also called welfare services or social work, include publicly or privately provided services intended to aid disadvantaged, distressed, or vulnerable persons or groups. The economic crisis is undermining the sustainability of social protection systems in the EU [1]: 24% of all the EU population (over 120 million people) are at risk of poverty or social exclusion. The fight against poverty and social exclusion is at the heart of the Europe 2020 strategy for smart, sustainable and inclusive growth. Social services deal with a number of

undesirable conditions that affect not only the quality of life of individuals, but also the equity and cohesion of society as a whole [10].

Why not producing prediction models for the field of social services as in health services?. Machine learning could answer a number of questions such as: will this individual suffer chronic social exclusion?; will generational transmission of poverty occurs in this family?; how much economic aid is needed to integrate this person into society?; how long does it take aid to have an impact on a case?. Something that may go unnoticed by outsiders to the field of data science is that all these questions are forms of *supervised learning*. Therefore, these questions fall into two broad categories: (1) *classification* (“is this A or B?”, or “is this A or B, or C...?”); and (2) *regression* (questions answered with a number: “how much”, “how many”, “how long”).

Unsupervised learning and reinforcement learning achieve outstanding results when applicable, but they are not adequate to answer these questions. On the other hand, supervised learning is based on the premise that lots and lots of labeled and annotated data is available. There are a number of challenges in gathering this labeled data in social services. (1) There is not public and accepted datasets in the field, typically because of privacy reasons. (2) Even if there were such data, the labels to predict would not correspond to the needs of all social services because there is a strong coupling between the predictive tool and the data it is fed with. (3) Moreover, the conditions social services are concerned about depend on the society they deal with, not allowing prediction results to be extrapolated from a country or even a city to another one. (4) Finally, the complex and multi-dimensional nature of processes such as social exclusion may require a case-by-case debate and deciding a label is complicated even for social workers experts.

For the reasons explained above, the hardest part of building new artificial intelligence solutions for social services is not the machine learning algorithms, but the data collection and labeling. This paper copes with this problem by a multi-agent architecture that combines semantic web technologies, exploratory data analysis techniques, and supervised machine learning methods. The architecture is composed of a number of cooperating intelligent agents that assist social workers and data scientists in the labeling of sensitive data and in the subsequent generation of prediction services.

The rest of the paper is organized as follows: section 2 revises the related work. Section 3 presents the proposed architecture. Finally, the preliminary conclusions obtained are presented in section 4.

2 Related works

Predictions models are widely used in insurance companies to allow customers to estimate their policies cost. Manulife Philippines [2] offers a number of online tools to calculate the likelihood of disability, critical illness, or death before the age of 65; based on age, gender, and smoking status. Health is another application field where risk estimations are undertaken for preventive purposes.

More specifically, the risk of heart disease can be estimated at different websites such as at the Mayo clinic web [3]. The labeling of these cases is relatively simple a posteriori: roughly speaking there is no doubt when someone has suffered one of these conditions. There are also a few online tools that social services may use for early detection. In this manner, Rank and Hirschl [14] give an online calculator that evaluates the probability of experiencing poverty in the next 5, 10 or 15 years. Labeling poverty cases is something automatic when defined as falling below a certain annual income¹.

The multi-dimensional nature of conditions such as social exclusion makes considerably more challenging to analyze, detect, treat, and predict it than poverty. There are a number of data analysis works in social exclusion that are detailed enough to learn from their labeling methods for the presented work. Ramos and Valera [13] use the *logistic regression* (LR) model to study social exclusion in 384 cases labeled by social workers through a manual heuristic procedure. According to this procedure, an individual is considered at a consolidated phase of exclusion if: (1) he or she is living for at least 3 years in unstable accommodation; (2) has very weak links, or none at all, with family or friends; (3) is almost permanently unoccupied; and, (4) presents a substantial or total loss of working habits, self-care or motivation for inclusion. Similar conditions are defined for the initial phase of exclusion. This example of rule of thumb used by the social workers illustrates the complexity and ambiguity of deciding if someone is suffering social exclusion. Moreover, the heuristic has to be defined before starting gathering data so the social workers can use it. Finally, the fully manual approach only allows a very limited number of cases: less than 400. Lafuente-Lechuga and Faura-Martínez [9] undertake an analysis of 31 predictors based on segmentation methods and LR. The authors consider the aggregation of scores in different fields related to social exclusion to decide if a person is under this condition. After a cluster analysis, this score is used to rank and analyze the most important variables to decide whether there is vulnerability to social exclusion. In a similar style, Haron [6] studies the social exclusion in Israel labeling data by various indicators that are aggregated in a single weighted average score. The author proposes the *linear regression* as a better alternative to the LR. The problem with this approach is that, besides the difficulty in defining these aggregations functions and weights, the machine learning techniques will tend to calculate precisely the aggregation formula since it is defined based exclusively on the training data. Suh et al. [16] analyze over 35K cases of 34 European countries using LR. The particular objective of this work is a subjective study and not an objective measure of the social exclusion, for which the researchers use LR over responses to a survey of direct questions about whether people feel excluded from society. Therefore, as the authors point out, there is a subjectivity aspect that is the responsibility of the interviewee instead of the social worker

¹ In this vein, the adult dataset [8] is a well known public labeled dataset that allows predicting whether an adult income exceeds \$50K a year based on a 1994 census database. It can be used to train prediction models as a proof of concept before collecting and labeling the own proprietary data.

expert. These inspiring works support the hypothesis that machine learning can greatly benefit social services. Also, that the most interesting questions to assist social services belong to the supervised learning. However, there are no general methods and tools proposed for labeling the data before building a predictive model.

A number of approaches study interesting synergies between agent theory and machine learning [15]. Ponni and Shunmuganathan [12] propose multi-agent system for classification in multi-relational databases with, among others, Support Vector Machines (SVM). Kiselev and Alhajj [7] describe an efficient adaptive multi-agent approach to continuous online clustering of streaming data in complex uncertain environments. Giannella et al. [5] propose an implementation of distributed clustering algorithms with multi-agent systems. Park and Oh [11] introduce a multi-agent system to filter data that automatically selects and tunes a clustering or dimensionality reduction method. These significant contributions improve machine learning paradigms in a number of aspects by rethinking them from the perspective of multi-agent systems. More importantly for the work presented here, several of these references deal with exploratory data analysis techniques such as clustering and dimensionality reduction. These are natural solutions to summarize, simplify, condense, and distill a collection of data before labeling it. However, the revised multi-agent systems do not offer specific guidelines to go from the clusters or the principal components to the wanted labels. Clusters are in the eye of the beholder, and the architecture presented in this paper instead of focusing on implementing faster clustering techniques, is meant to adapt to the beholder and to recommend actions to label data intuitively.

3 Multi-agent architecture segmentation and prediction

This section describes the proposed architecture, see figure 1.

In the lower layer of the architecture, there is an interface that allows accessing the databases that are used for the different applications and records in social services and linking them with the rest of the architecture. The agents of this layer, besides controlling the protocols of access to the databases, will ensure that the information that the upper layers obtain is anonymous. Users with special privileges may require this layer, through services in upper layers and with the purpose of labeling a case, to link an anonymous identifier with an identity.

The data handled by the interface are also accessed by a layer of persistence transverse to the architecture. This persistence layer has capacities of semantic technologies as dealing with ontologies in languages such as RDFS and OWL. To favor the tagging service, the architecture offers functionality for the formal representation of the knowledge treated by social services through a network of ontologies. In addition to these ontologies, the layer stores intermediate data such as: data tables obtained from pre-processing the databases accessed by the interface, unsupervised learning models, supervised learning models, and users' preferences and history of decisions for recommendation and decision support.

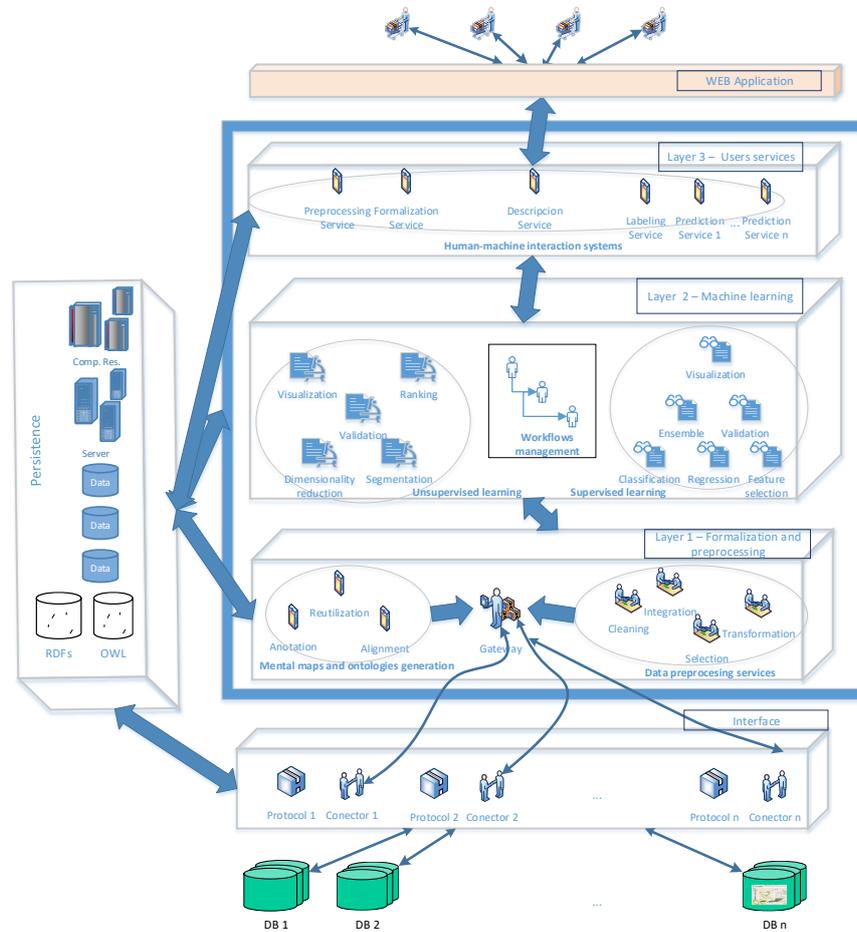


Fig. 1. A multi-agent architecture for labeling data and generating prediction models.

Layer 1 of the architecture is devoted to the formalization and preprocessing of social services data, on the one hand, to have a powerful query model based on semantic technologies and, on the other hand, to allow the machine learning methods to learn from this data. These two differentiated services are connected through a gateway: data processing services and ontology construction. The first service includes agents specialized in: (1) data selection; (2) its integration from various sources accessed by the interface; (3) cleaning the data by detecting noise and inconsistencies; (4) and, transforming the data into forms suitable for mining. The agents in charge of assisting the process of generating mental maps and ontologies are specialized in tasks such as: (1) data annotation; (2) reuse of already built ontologies; and, (3) ontology alignment, i.e. determining correspondences between concepts in ontologies.

Layer 2 of the architecture is the machine learning layer. This layer provides unsupervised learning services to obtain simplified representations of data for labeling. The agents of these services specialize in segmentation through different algorithms, dimensionality reduction, validation methods of unsupervised learning (as the silhouette method to determine an adequate number of clusters or the information loss in principal component analysis), the visualization of clusters and data in n-dimensions (with methods as star diagrams or Chernoff faces), and the calculation of rankings by similarity to a given case. The layer also offers supervised learning services so that, once the first labels are available, interpretable models of these data are built such as rule-based classifiers. These machine learning models can be used by social workers as heuristics to label new cases. In addition, learning paradigms of higher predictive power will also be generated requiring little or no parameter tuning by social services, such as random forest or AdaBoost with decision trees (considered one of the best out-of-the-box classifier). Quality metrics for these supervised models can be used as a stop criterion in the process of labeling cases. The agents of this service specialize in: classification, regression, feature selection, ensemble methods to combine the results of several models, validation (among others: leave-one-out cross-validation, fold cross-validation, and with a test set), and visualization of models and evaluation metrics. A workflow manager allows combining this layer processes in different workflows.

Layer 3 offers distinguished user services: (1) pre-processing services for data scientists; (2) formalization services for ontology engineers; and (3) concept description services and (4) labeling services for social workers. The agents these services are composed of are the only ones that interact with the agents of lower layers and the transverse layer of persistence. For each user service, there is a *decision support system* (DSS) that provide users with explicit decision suggestions. For the labeling service, an example of decision suggestion could come after the first input of labels. If an underfitting situation is detected with a few tagged examples, labeling new examples will not improve the future prediction model. In this case, some suggestions include: (1) revising the labels that might be inconsistent; (2) collecting more fields for the cases; (3) or, considering changing the purpose of the prediction model. On the other hand, if a high accuracy is achieved with the currently labeled data (or it does not improve in many iterations), stopping the labeling process could also be suggested. Finally, these services agents have to learn from the users' preferences and recommend actions and alternatives through techniques as *collaborative filtering*. Once again, clusters are in the eye and there are no inherently better cluster analysis methods than others.

In the top layer, user services are accessed through responsive web applications. In this way, users can use these services through a variety of devices: smartphones, tablets, laptops, etcetera. The ultimate goal is to provide social workers with intelligent prediction models in the palm of the hand, which allow them to anticipate events for the sake of their "social patients". But as discussed, the hardest part to get there is labeling the data.

4 Conclusion and future works

This paper presents a multi-agent architecture for labeling data and generating social services prediction models. The proposal responds to the enormous importance of social services in today's Europe and to the difficulty of generating predictive models that help social workers in their day-to-day work. To improve this situation, the architecture supports an iterative and incremental data labeling until a predictive model that attends a specific social service is obtained.

The core of the proposal is based on offering services and assistance to social workers in cluster analysis and dimensionality reduction as a means to summarize, simplify, condense, and distill a collection of data before labeling it. Intelligent agents not only automate this analysis as much as possible, but also learn from workers' preferences since there is a lack of objectivity in the generation of useful summaries and visualizations of unlabeled data. Furthermore, the architecture includes agents for supervised learning that, in each iteration in which new labels are added, contribute with: explanatory models of the data; selections of the most important predictors; and, stopping conditions to the labeling process automatically checked. Finally, it provides a service for the consultation of ontology networks that facilitates the unambiguous description of the concepts included in the cases to be tagged and their relations.

Although the architecture has not been implemented, many of the ideas and proposals have been put into practice for the generation of an online social exclusion prediction service in the Spanish region of Castilla y León (<http://webpact.oeg-upm.net/>).

Future work includes: implementing the architecture in a multi-agent platform; extending the decision support system for the labeling service; and, a better exploitation of the ontologies and semantic resources to include forms of advanced learning such as case-based reasoning, transfer learning, and graph mining.

Acknowledgments

This publication would not have been possible without the inputs and collaboration of the Social Services of Castilla y León. This research work is supported by the "Junta de Castilla y León" under the public contract: "Servicios de elaboración de modelos matemáticos para realizar segmentación poblacional" (A2016/000271); by the EU Programme for Employment and Social Innovation (EaSI) under the project PACT ("people-oriented case management for social inclusion proactive model"); and, by the Spanish Ministry of Economy, Industry and Competitiveness under the R&D project Datos 4.0: Retos y soluciones (TIN2016-78011-C4-4-R, AEI/FEDER, UE).

References

1. European Commission's DG for Employment, Social Affairs & Inclusion. <http://ec.europa.eu/social/main.jsp?catId=751>. Accessed: February of 2017.
2. Manulife Philippines. Calculate your risk, your partner's risk or both. <http://www.insureright.ca/what-is-your-risk>. Accessed: February of 2017.
3. Mayo Clinic. Heart Disease Risk Calculator. <http://www.mayoclinic.org/diseases-conditions/heart-disease/in-depth/heart-disease-risk/itt-20084942>. Accessed: February of 2017.
4. L. de Oliveira. Fueling the Gold Rush: The Greatest Public Datasets for AI. <https://goo.gl/mJ08nf>. Accessed: February of 2017.
5. C. Giannella, R. Bhargava, and H. Kargupta. *Multi-agent Systems and Distributed Data Mining*, pages 1–15. Springer Berlin Heidelberg, Berlin, Heidelberg, 2004.
6. N. Haron. *On Social Exclusion and Income Poverty in Israel: Findings from the European Social Survey*, pages 247–269. Springer US, Boston, MA, 2013.
7. I. Kiselev and R. Alhajj. A self-organizing multi-agent system for adaptive continuous unsupervised learning in complex uncertain environments. In D. Fox and C. P. Gomes, editors, *Proceedings of the Twenty-Third AAAI Conference on Artificial Intelligence, AAAI 2008, Chicago, Illinois, USA, July 13-17, 2008*, pages 1808–1809. AAAI Press, 2008.
8. R. Kohavi and B. Becker. Adult Data Set. <https://archive.ics.uci.edu/ml/datasets/Adult>. Accessed: February of 2017.
9. M. Lafuente-Lechuga and U. Faura-Martínez. Análisis de los individuos vulnerables a la exclusión social en españa en 2009. *Anales de ASEPUMA*, (21), 2013.
10. R. Levitas, C. Pantazis, E. Fahmy, D. Gordon, E. Lloyd, and D. Patsios. *The multi-dimensional analysis of social exclusion*. London: Social Exclusion Task Force, Cabinet Office, 2007.
11. J.-E. Park and K.-W. Oh. Multi-agent systems for intelligent clustering. *International Journal of Computer, Electrical, Automation, Control and Information Engineering*, 1(11):275 – 280, 2007.
12. J. Ponni and K. L. Shunmuganathan. Multi-agent system for data classification from data mining using svm. In *2013 International Conference on Green Computing, Communication and Conservation of Energy (ICGCE)*, pages 828–832, Dec 2013.
13. J. Ramos and A. Varela. Beyond the margins: Analyzing social exclusion with a homeless client dataset. *Social Work & Society*, 14(2), 2016.
14. M. R. Rank and T. A. Hirschl. Calculate Your Economic Risk. *New york times*, 2016.
15. E. Serrano, M. Rovatsos, and J. Botia. A qualitative reputation system for multi-agent systems with protocol-based communication. In *Proceedings of the 11th International Conference on Autonomous Agents and Multiagent Systems - Volume 1, AAMAS '12*, pages 307–314, Richland, SC, 2012. International Foundation for Autonomous Agents and Multiagent Systems.
16. E. Suh, P. TiffanyVizard, and T. AsgharBurchardt. Quality of life in europe: Social inequalities. *3rd European Quality of Life Survey*, 2013.