

Explain your clusters with words. The role of metadata in interactive clustering

Maciej Mozolewski¹, Samaneh Jamshidi², Szymon Bobek¹ and Grzegorz J. Nalepa¹

¹Jagiellonian Human-Centered Artificial Intelligence Laboratory (JAHCAI) and Institute of Applied Computer Science, Jagiellonian University, 31-007 Kraków, Poland

²Center for Applied Intelligent Systems Research (CAISR), Halmstad University, Halmstad, Sweden
m.mozolewski@doctoral.uj.edu.pl, samaneh.jamshidi@hh.se, {szymon.bobek, grzegorz.j.nalepa}@uj.edu.pl,

Abstract

In this preliminary work, we present an approach for augmentation of clustering with natural language explanations. In clustering there are 2 main challenges: a) choice of a proper, reasonable number of clusters and b) cluster analysis and profiling. There is a plethora of technics for a) but not so much for b), which is in general a laborious task of explaining obtained clusters. In this work, we propose a method that aids experts in cluster analysis by providing iterative, human-in-the-loop methodology of generating cluster explanations. In a convincing example, we show how the process of clustering on a set of *objective variables* could be facilitated with textual *metadata*. In our case, images of products from online fashion store are used for clustering. Then, product descriptions are used for profiling clusters.

1 Introduction

Data analysts and data scientists are often faced with the task of describing phenomena in a way that is understandable to the audience. On the one hand, they have some data that they can describe statistically, categorize, or predict events based on it, etc. On the other hand, they want to deliver their observations in some kind of narrative that they can "sell" to decision-makers. There is even a phenomenon called *data storytelling*. As the authors of [Matei and Hunter, 2021] state, this is a type of *narrative* in which science provides explanations about cause-and-effect relationships. Science lends itself well to storytelling because new discoveries can be surprising, and therefore interesting, to the public. In our work we want to show that this intuitive approach is reflected in the work of data scientists, which manifests itself in the way they use different types of data.

Assigning labels to groups of similar objects is one of the ways how humans describe the world [Rosch and Lloyd, 1978]. It begins with the notion that some phenomena or entities differ from each other and that they could be divided into distinct classes. Clarification of the differences between groups gets better and better along with the knowledge gained about the instances that form different groups. Finally, one is giving names to those different categories of entities. In

essence, clustering in machine learning is no different process.

Clustering is an intrinsically subjective task and requires human assessment [Bae *et al.*, 2020]. It is a purely statistical method which finds homogeneous groups of entities. It belongs to the family of unsupervised learning algorithms in contrast to classification or regression, which are supervised. At every step of this process, the user makes decisions based on her/his domain knowledge. First, the user needs to select features (variables) used by the algorithm. Secondly, the user selects the type of algorithm, similarity measures, number of clusters or size of the smallest one. Finally, she or he checks clusters by describing objects belonging to subsequent groups. It also follows that the process is iterative.

From our expertise in e-commerce and Industry 4.0, we often see distinctions between two types of data. There are *objective data* and the *subjective data* or *metadata*. For instance, in e-commerce popular approach for recommendations is based on finding users similar to each other in terms of interactions with products. Thus, *objective data* is composed of the behaviors of shoppers. The categories, titles and descriptions of the products form *metadata*, which is usually the result of the joint work of many employees of the e-store. For rolling steel factories, predictive maintenance models are derived mainly from *objective sensory data*, such as temperature, force, etc. Factory accounting data are *metadata*.

Objective data can be viewed as story or the totality of facts that have occurred. They can be difficult for humans to understand but lend themselves well to clustering. *Metadata*, in contrast, is more a type of narrative, or how the algorithm's outcome is presented by data scientists to their audience. *Metadata* is more subjective, making it more suitable for justification and formulation of conclusions and explanations.

The more *objective* the data is, the more it is suited for modeling the phenomena, be it physical, business, sociological or psychological in nature. *Metadata* is more suitable for explaining the model to the user, convincing her or him, and prompting to make decisions and actions based on this knowledge. It is more prone to error because of its conventionality and subjectivity, but they speak to humans.

In this work we propose a method that allows for clustering dataset with *objective data*, and explain differences between clusters with *metadata*. We use XAI methods to ex-

plain differences between clusters using *metadata* which is perfectly understandable by humans, but may not be of sufficient quality to perform valid clustering. The selection of the most interpretable *metadata* is iterative and human-guided. In our example, we show how image-based clustering can be enhanced with textual description of clusters. We argue that such an approach can lead to better utilisation of *metadata* for cluster analysis purposes, which results in better understanding of clusters which is the final goal of every clustering task. Furthermore, it allows for checking the consistency between two or more possible instance representations (image and text) which might be crucial in domains that rely on both (e.g. e-commerce).

The rest of the paper is organized as follows. In Section 2 we present current research in the area of interactive clustering and human-guided cluster analysis. The description of our method along with use-case studies is given in Section 3. Finally, we conclude our work and show perspectives of its further development is presented in Section 4.

2 Related works

In the survey on interactive clustering [Bae *et al.*, 2020] authors have distinguished 2 groups of approaches in terms of how interaction occurs. In general, users can interact indirectly with the tool, by changing the parameters of the algorithm, or directly by giving feedback to the result of clustering. The parameters adjusted most frequently are the number of clusters and the similarity threshold [Andrienko and Andrienko, 2015; Arin *et al.*, 2018]. For topic modeling, users are given the option to select keywords and set their relative importance [El-Assady *et al.*, 2018]. On the other hand, direct feedback might be realized by highlighting incorrect instances of splitting or merging the resulting clusters [Yang *et al.*, 2017]. For textual data, users can provide the tool with blacklists of incorrect topic labels or set similar restrictions [Chang *et al.*, 2016]. Based on this signal, the clustering tool learns user preferences and tries to incorporate the knowledge in the next iteration.

Explainable AI approaches have become particularly important, and although most work is generally focused on supervised learning, some works have been done to explain clusters. One of the most common methods for understanding clustering methods is visualization. By using low-dimensionality embedding and displaying them in two- or three-dimensional dimensions, one can get an overview of the clusters and their data. However, these visualizations are not always understandable and explainable.

The decision tree is one of the inherently interpretable algorithms. So one common way to explain models is to use decision trees. Nevertheless, the critical point for explaining the decision tree is its depth because decision trees with high depth no longer are interpreted, so we must pay attention to the depth of the tree produced. Using a small decision tree to divide a dataset into k clusters provides explainable clusters, but this approach has a trade-off between being explainable and accuracy. IMM algorithm [Dasgupta *et al.*, 2020] approximates k -means and k -median clustering by a threshold tree with k leaves. ExKMC [Frost *et al.*, 2020] uses a threshold

tree to provide an explainable k -mean clustering in which the number of tree leaves is greater than the number of clusters.

Besides, visualization or providing some conditions on features, using text data is reasonable to generate explanations to users. In [Hendricks *et al.*, 2016] authors use captions of the images along with the images to create a more discriminative classification. In addition, they use this *metadata* to provide language explanation and generate a text description for each class. However, by blending textual and image modalities into one dataset, authors limit the possibility of checking consistency between these two types of data and implicitly assume the correctness of possibly wrong image descriptions.

Similarly, in many other methods that aim at explaining differences between discovered clusters, the clustering task is transformed to classification one, and the classifier is then explained with available XAI methods such as LIME [Ribeiro *et al.*, 2016], Anchor [Ribeiro *et al.*, 2018], LUX [Bobek and Nalepa, 2021], etc. One of the most recent implementations of such approach can be found in [Lötsch and Malkusch, 2021].

Another approach is given in [Bobek *et al.*, 2021], where authors present a toolkit for conformance checking between expert knowledge and automatic clustering. The differences between expert-based clustering and automated clustering are justified with XAI methods and the process is iterative. However, the explanations are not human-guided, and the expert has no impact on the way they are generated. In particular, it is not possible to provide additional *metadata* for explanations, nor modify the set of concepts that are used for explanations.

In all the cases the process is not iterative, nor human-guided. Finally, to the best of authors' knowledge, neither of the approaches known in the literature divides data into *objective* part with a good quality for cluster algorithms, but poor explanation capabilities and *metadata* with possibly worse potential as clustering features, but better explanation capabilities and possible inconsistencies with *objective* data that should be fixed. Addressing these issues was the primary motivation of our work that will be described in more detail in the following sections.

3 Cluster analysis with *metadata*

In this section, we will show how our method could be applied to real-case scenarios. We choose an example from the e-commerce field because the authors have experience working in this industry. Specifically, we work with online stores to provide them, among others, with recommendations of products to their end-users (clients).

In real-life scenarios, data about products are stored in product catalogs in shop databases, and most often exchanged with so-called product feeds (XML documents). We used a public dataset from Kaggle¹. This dataset in terms of content resembles a product feed for an online store of a medium size product catalog. It consisted of 44000 products with category labels, titles, and images. For the code accompanying this

¹ See: <https://www.kaggle.com/datasets/paramaggarwal/fashion-product-images-small>

example see GitHub repository².

As has been said before, we treat images as *objective data*. We used embeddings of images obtained via MobileNetV2 [Sandler *et al.*, 2019] as an input to clustering pipeline. The fully-connected layer at the top of the network was disregarded because we were not interested in the classification done by the model. The output of the final layer of the model was of length 20480. We used Singular Value Decomposition (SVD) with normalization to reduce the dimensionality of embeddings, leaving at least 90 percent of the variance.

In this section, we will present tools dedicated to data scientists who would like to perform clustering. We propose a 2-step clustering loop, which consists of k-means clustering and textual explanations of clusters. Data preparation also could be performed more than once, if needed. For the sake of simplicity, we call it "step 0" in this work.

3.1 Data preparation

The method requires 2 types of data: *objective* and *metadata* as defined in the previous section. In "step 0" method provides users with helper functions to prepare both types of data. For *objective data* there is a function that performs a reduction of dimensionality via SVD followed by normalization. It works on any numerical data, which could be as well as one-hot variables and continuous real values (floats). User sets the percentage of explained variance left after SVD reduction. The optimal count of new dimensions could be determined automatically by our algorithm. This is done by probing different dimension counts with *scipy.optimize* package, so the user does not need to do this manually. Regarding *metadata* which is textual, there are wrappers built on top of the SpaCy³ and NLTK⁴ libraries. Users can contact text columns, lemmatize, remove stopwords and perform Tf-Idf vectorization. For numerical *metadata*, we found a way to incorporate them into textual explanations. For instance, the year could be re-coded as the label "year2022", which will be easily interpreted along the pipeline. Other numerical variables could be re-coded to low/medium/high bins, based on quartiles. Finally, the user constructs the "Pipeline" object and initializes it with 2 datasets: *objective* and *metadata*.

3.2 Assistance in clustering

The first step corresponds to running the unsupervised clustering algorithm. Typically, the person performing the analysis starts with the dilemma of choosing the number of clusters. It can be resolved with her/his background knowledge, intuition, practicality prerequisites, or just a trial and error approach. To give our users a hint in this regard, we use the T-SNE [Hinton and Roweis, 2003] 2-dimensional projection of the data. At the moment, this is a solely visual clue. It is depicted in Figure 1. If data have an underlying structure, points representing observations will cluster, which would be

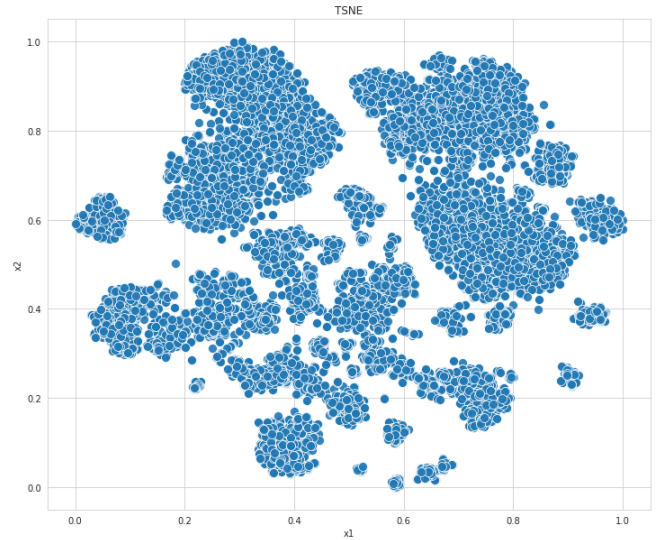


Figure 1: Preliminary visualisation of *objective data* in 2-D projection with T-SNE dimensionality reduction.

observed on the chart. As T-SNE on massive data could be resource intensive, the default is to run this process on random subsample and cache results. Additionally, users can apply textual labels to the T-SNE chart, plotted on a subsample of data, to avoid cluttering the chart. Labels could represent the most important pieces of *metadata*, such as the label, the observation id, and summary of description. The next clue is derived from the silhouette score on a plot in the Figure 2. The

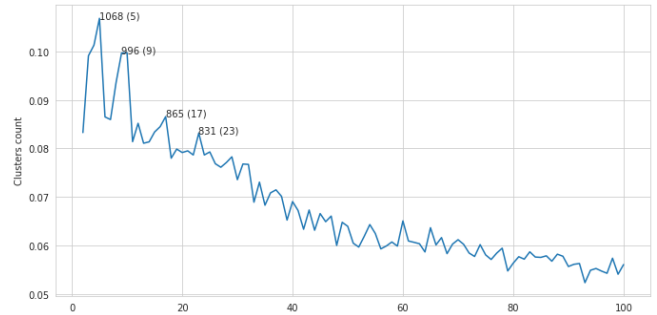


Figure 2: Silhouette score with cluster counts and values.

range of the number of clusters to be tested is provided in accordance with the previous clue. To speed-up computations, this plot could be obtained on a random subset, and results are cached for further reference. For now, the user interprets the plot on her/his own. Finally, clustering with k-means is performed on all observations. Visualization with T-SNE is presented, this time with clusters colored in different colors, which is depicted in Figure 3.

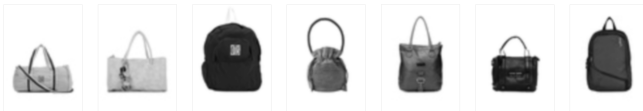
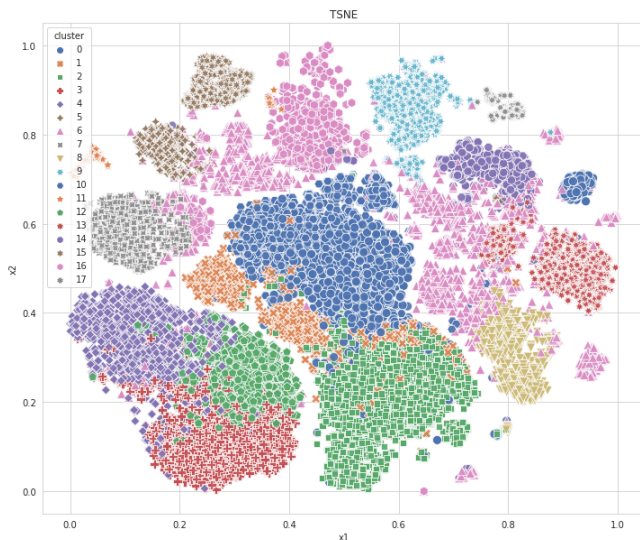
3.3 Interactive explanations

The second step is to explain the clusters so that the person performing the data analysis can assess the result. We would like to give users the freedom to refine explanations. Thus,

²See: <https://github.com/mozo64/xai-survey/blob/main/src/example1-clustering-products-fashion.ipynb>

³See: <https://spacy.io/>

⁴See: <https://www.nltk.org/>



we provide her or him with the possibility to influence explanations by extending stopwords with his own terms. On the other hand, we initialize the whitelist with keywords like "year2022", defined in "step 0". Then we use the Tf-Idf vectorizer, taking into account the aforementioned lists. Vectors are used to train decision tree classifiers. The size of the list of additional terms is under the control of a user. She or he can change it and interactively observe the result in a Figure 6, which gives an insight into what terms were relevant for classifier. Furthermore, for every cluster, example observations, word clouds and LIME explanation for description of random observation are presented. For instance for the cluster 13 in

4 Summary

In this work, we presented the method that allows for explaining clusters with concepts that could be more human-readable than the data which was used as an input to clustering algorithm. We based our method on the observation that different types of data are suitable in different degrees to clustering and explaining tasks. We demonstrated the feasibility of our approach on the e-commerce example, where images were treated as input for clustering and textual descriptions of images as basis for cluster descriptions.

In its current version, adaptability to the needs of a human expert is provided by the possibility of customization of *Metadata*. *Metadata*-based explanations can be refined in two ways. The user can influence the number of tokens used in the explanation or can directly influence the list of tokens by adding words to the whitelist or blacklist. After each such change, the expert can see how it affects the classifier used as an explainer, both globally and at the level of random observations for a cluster. One can modify the scope of the *metadata* and the way individual observations are presented. If the results, despite the changes, are not satisfactory, it may imply the need to go back to choosing the number of clusters.

We treat metadata as something given. From our experience in e-commerce, we know that product descriptions are the result of the work of many e-store employees, but we do not want to interfere with them. For example, a product recommendation system in which the objective data would be product images and the online shopping behavior would be switched on for a particular e-commerce after the explanation is accepted by a decision maker in the e-store. The explanations could be refined by technical support of the platform before being shown to the e-store employee. The tool is intended for use by a sole data scientist. However, in other situations, collaborative knowledge engineering approaches are also possible. Metadata could be created and enhanced in systems such as LOKI⁵ [Kutt, 2016].

In future work we would like to improve our method with several extensions. We will focus on automatically proposing number of clusters based on both embedding features with methods similar to T-SNE and metrics like silhouette score. We want to test clustering techniques other than k-means. For instance, hierarchical clustering could be more suitable in e-commerce, where taxonomies of products are multilayer. Word clouds could be replaced with topic analysis with Latent Dirichlet Allocation or techniques derived from Natu-

⁵See: <https://loki.re/wiki/docs:start#loki>

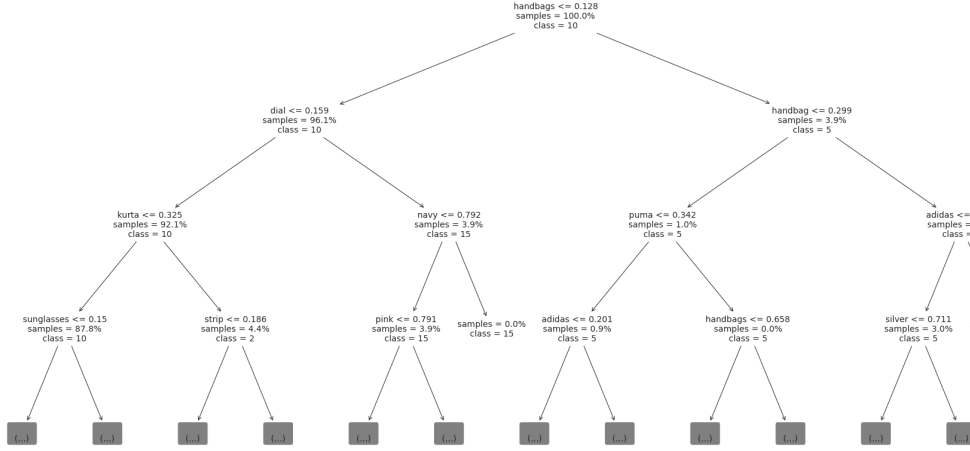


Figure 6: Decision tree classifier which explains how clusters differ in terms of *metadata*, here pruned to level 4.

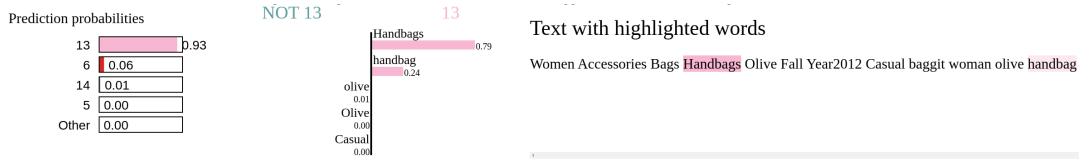


Figure 7: LIME explanation for a random observation drawn from the cluster 13.

ral Language Generation. Another interesting direction is to construct explanations with other modalities, like visual, by something more sophisticated than presenting example images. It could be done for instance with image captioning.

Acknowledgments

This paper is funded from the XPM (Explainable Predictive Maintenance) project funded by the National Science Center, Poland under CHIST-ERA programme Grant Agreement No. 857925 (NCN UMO-2020/02/Y/ST6/00070).

The work of Szymon Bobek has been additionally supported by a HuLCKA grant from the Priority Research Area (Digiworld) under the Strategic Programme Excellence Initiative at the Jagiellonian University (U1U/P06/NO/02.16).

The work of Samaneh Jamshidi was supported by CHIST-ERA grant CHIST-ERA-19-XAI-012 funded by Swedish Research Council.

The work of Maciej Mozolewski has been additionally supported by Edrone Sp. z o.o.⁶, which provided computer resources for machine learning.

References

[Andrienko and Andrienko, 2015] Gennady Andrienko and Natalia Andrienko. Visualization support to interactive cluster analysis. In Albert Bifet, Michael May, Bianca Zadrozny, Ricard Gavaldà, Dino Pedreschi, Francesco

Bonchi, Jaime Cardoso, and Myra Spiliopoulou, editors, *Machine Learning and Knowledge Discovery in Databases*, pages 337–340, Cham, 2015. Springer International Publishing.

[Arin *et al.*, 2018] Inanc Arin, Mert Erpam, and Yucel Saygin. I-tvec: Interactive clustering tool for twitter. *Expert Systems with Applications*, 96:1–13, 04 2018.

[Bae *et al.*, 2020] Juhee Bae, Tove Helldin, Maria Riveiro, Sławomir Nowaczyk, Mohamed-Rafik Bouguelia, and Göran Falkman. Interactive clustering: A comprehensive review. *ACM Comput. Surv.*, 53(1), feb 2020.

[Bobek and Nalepa, 2021] Szymon Bobek and Grzegorz J. Nalepa. Introducing uncertainty into explainable ai methods. In Maciej Paszynski, Dieter Kranzlmüller, Valeria V. Krzhizhanovskaya, Jack J. Dongarra, and Peter M. A. Sloot, editors, *Computational Science – ICCS 2021*, pages 444–457, Cham, 2021. Springer International Publishing.

[Bobek *et al.*, 2021] Szymon Bobek, Michal Kuk, Jakub Brzegowski, Edyta Brzychczy, and Grzegorz J. Nalepa. Knac: an approach for enhancing cluster analysis with background knowledge and explanations. *CoRR*, abs/2112.08759, 2021.

[Chang *et al.*, 2016] Shuo Chang, Peng Dai, Lichan Hong, Cheng Sheng, Tianjiao Zhang, and Ed H. Chi. Appgrouper: Knowledge-based interactive clustering tool for app search results. *Proceedings of the 21st International Conference on Intelligent User Interfaces*, 2016.

⁶See: <https://edrone.me/en/>

- [Dasgupta *et al.*, 2020] Sanjoy Dasgupta, Nave Frost, Michal Moshkovitz, and Cyrus Rashtchian. Explainable k -means and k -medians clustering. *arXiv preprint arXiv:2002.12538*, 2020.
- [El-Assady *et al.*, 2018] Mennatallah El-Assady, Rita Sevastjanova, Fabian Sperrle, Daniel A. Keim, and Christopher M. Collins. Progressive learning of topic modeling parameters: A visual analytics framework. *IEEE Transactions on Visualization and Computer Graphics*, 24:382–391, 2018.
- [Frost *et al.*, 2020] Nave Frost, Michal Moshkovitz, and Cyrus Rashtchian. Exkmc: Expanding explainable k -means clustering. *arXiv preprint arXiv:2006.02399*, 2020.
- [Hendricks *et al.*, 2016] Lisa Anne Hendricks, Zeynep Akata, Marcus Rohrbach, Jeff Donahue, Bernt Schiele, and Trevor Darrell. Generating visual explanations. In *European conference on computer vision*, pages 3–19. Springer, 2016.
- [Hinton and Roweis, 2003] Geoffrey Hinton and Sam Roweis. Stochastic neighbor embedding. *Advances in neural information processing systems*, 15:833–840, 2003.
- [Kutt, 2016] Krzysztof Kutt. Semantic wikis versioning with bifrost. In *Doctoral Consortium of the 15th Edition of AI*IA*, volume 1769 of *CEUR Workshop Proceedings*, pages 34–39. CEUR-WS.org, 2016.
- [Lötsch and Malkusch, 2021] Jörn Lötsch and Sebastian Malkusch. Interpretation of cluster structures in pain-related phenotype data using explainable artificial intelligence (xai). *European Journal of Pain*, 25(2):442–465, 2021.
- [Matei and Hunter, 2021] Sorin Adam Matei and Lucas Hunter. Data storytelling is not storytelling with data: A framework for storytelling in science communication and data journalism. *The Information Society*, 37(5):312–322, 2021.
- [Ribeiro *et al.*, 2016] Marco Tulio Ribeiro, Sameer Singh, and Carlos Guestrin. “Why should i trust you?”: Explaining the predictions of any classifier. In *Proceedings of the 22nd ACM SIGKDD International Conference on Knowledge Discovery and Data Mining, KDD ’16*, page 1135–1144, New York, NY, USA, 2016. Association for Computing Machinery.
- [Ribeiro *et al.*, 2018] Marco Tulio Ribeiro, Sameer Singh, and Carlos Guestrin. Anchors: High-precision model-agnostic explanations. In *AAAI*, 2018.
- [Rosch and Lloyd, 1978] Eleanor Rosch and Barbara Lloyd. *Cognition and Categorization*. Lawrence Elbaum Associates, 1978.
- [Sandler *et al.*, 2019] Mark Sandler, Andrew Howard, Menglong Zhu, Andrey Zhmoginov, and Liang-Chieh Chen. Mobilenetv2: Inverted residuals and linear bottlenecks, 2019.
- [Yang *et al.*, 2017] Lei Yang, Dai Yu, Zhang Bin, and Yang Yang. Interactive k -means clustering method based on user behavior for different analysis target in medicine. *Computational and Mathematical Methods in Medicine*, 2017:1–9, 10 2017.