

# Explaining classification systems using sparse dictionaries

A. Apicella, F. Isgrò, R. Prevete, A. Sorrentino, G. Tamburrini \*

Dipartimento di Ingegneria Elettrica e delle Teconologie dell'Informazione  
Università degli Studi di Napoli Federico II, Italy

**Abstract.** A pressing research topic is to find ways to explain the decisions of machine learning systems to end users, data officers, and other stakeholders. These explanations must be understandable to human beings. Much work in this field focuses on image classification, as the required explanations can rely on images, therefore making communication relatively easy, and may take into account the image as a whole. Here, we propose to exploit the representational power of sparse dictionaries to determine image local properties that can be used as crucial ingredients of humanly understandable explanations of classification decisions.

## 1 Introduction

Machine Learning (ML) techniques enable one to develop algorithmic systems that learn from observations. Many ML techniques (e.g., Support Vector Machines (SVM) and Deep Neural Networks (DNN)) give rise to systems whose behavior is often hard to interpret [1]. Although some ML techniques come with reasonably interpretable mechanisms and Input/Output (I/O) relationships (e.g., decision trees), this is not the case for a wide variety of ML systems, whose processing and I/O relationships are often difficult to understand [2]. Various senses of interpretability for learning systems have been distinguished and analyzed [3], and various approaches to overcoming their opaqueness are now being pursued [4, 5]. For example, in [6] a series of techniques for the interpretation of DNN are discussed, and in [7] a wide variety of motivations underlying interpretability needs are examined, thereby refining the notion of interpretability in ML systems. In the context of this multifaceted interpretability problem [8, 9], we focus on the issue of what it is to explain the behavior of ML classification systems for which only I/O relationships are accessible, i.e., the learning system is seen as a black-box. In literature, this type of approach is known as *model agnostic* [10].

Various model agnostic approaches have been developed to give *global* explanations exhibiting a class prototype which the input data can be associated to [4, 5, 8, 6]. These explanations are given in response to explanation requests that are usually expressed as why-questions: “Why were input data  $x$  associated to class  $C$ ?”. Specific why-questions which may arise in connection with actual learning systems are : “Why was this loan application rejected?” and “Why was this image classified as a fox?”. However, prototypes often make

---

\*The research presented in this paper was partially supported by the national project Perception, Performativity and Cognitive Sciences (PRIN Bando 2015, cod. 2015TM24JS-009).

rather poor explanations available. For instance, if an image  $x$  is classified as “fox”, the explanation provided by means of a fox-prototype is nothing more than a “because it looks like this” explanation: one would not be put in the position to understand what features (parts) of the prototype are associated to what characteristics (parts) of  $x$ . In order to go beyond this impoverished level of understanding, instead of merely giving the user a global explanation, one might attempt to provide a *local* explanation, which highlights salient parts of the input [10]

In this paper, we propose a model agnostic framework that returns local explanations based on *dictionaries* of local and humanly interpretable elements of the input. This framework can be functionally described in terms of a three-entity model, composed of an *Oracle* (an ML system, e.g. a classifier), an *Interrogator* raising explanations requests about the Oracle’s responses, and a *Mediator* helping the Interrogator to understand the answer given by the Oracle. The Mediator plays a crucial explanatory role here, by advancing hypotheses on what humanly interpretable elements are likely to have influenced the Oracle output. More specifically, elements are computed which represent humanly interpretable features of the input data, with the constraint that both prototypes and input can be reconstructed as linear combinations of these elements. Thus, one can establish meaningful associations between key features of the prototype and key features of the input. To this end, we exploit the representational power of sparse dictionaries learned from the data, where atoms of the dictionary selectively play the role of humanly interpretable elements, insofar as they afford a local representation of the data. Indeed, these techniques provide data representations that are often found to be accessible to human interpretation [11]. The dictionaries are obtained by a Non-negative Matrix Factorization (NMF) method [12, 2, 13], and the explanation is determined using an Activation-Maximization (AM) [4, 8] based technique, that we call *Explanation Maximization*.

The article is organized as follows: in Section 2 we present the overall architecture; experiments and results are discussed in Section 3, while Section 4 is devoted to concluding remarks and future developments.

## 2 Proposed Approach

Given an oracle  $\Omega$ , an input  $\vec{x}$  and an  $\Omega$ ’s answer  $\hat{c}$  (regardless of whether it is correct or not), we want to give an explanation of the answer provided by the model  $\Omega$  that is humanly interpretable.

As we want to obtain humanly interpretable elements which, combined together, can provide an acceptable explanation for the choice made by  $\Omega$ , we search for an explanation having the following qualitative properties: 1) the explanation must be expressed in terms of a *dictionary*  $V$  whose elements (atoms) are easily understandable by an interrogator; 2) the elements of the dictionary  $V$  (have to represent “local properties” of the input  $\vec{x}$ ; 3) the explanation must be composed by few dictionary elements.

We claim that considering as elements atoms of a sparse coding from a sparse

dictionary, and using sparse coding methods together with an AM-like algorithm we obtain explanations satisfying the properties described above.

## 2.1 Sparse Dictionary learning

The first step of the proposed approach consists in finding a “good” dictionary  $V$  that can represent data in terms of humanly interpretable atoms.

Let us assume that we have a set  $D = \{(\vec{x}^{(1)}, c^{(1)}), (\vec{x}^{(2)}, c^{(2)}), \dots, (\vec{x}^{(n)}, c^{(n)})\}$  where each  $\vec{x}^{(i)} \in \mathbb{R}^d$  is a column vector representing a data point, and  $c^{(i)} \in C$  its class. We arrange all  $\vec{x}^{(i)}$  s.t.  $c^{(i)} = c$  in a matrix  $X^{(c)} = (\vec{x}^{(i)})$ . The dictionary can be constructed following two different strategies: 1) *one-for-all*, a single dictionary for all classes; 2) *one-for-one*, a different dictionary for each class. The former strategy imposes to learn a single dictionary  $V \in \mathbb{R}^{d \times k}$  of  $k$  atoms across multiple classes and an encoding  $H \in \mathbb{R}^{k \times n}$  s.t.  $X = VH + \epsilon$  where  $X = (X^{(1)} | X^{(2)} | \dots | X^{(|C|)})$  and  $\epsilon$  is the error introduced by the coding. Every column  $\vec{x}_i$  in  $X$  can be expressed as  $\vec{x}_i = V\vec{h}_i$  with  $h_i$   $i$ -th column of  $H$ . The latter strategy imposes to learn for each class  $c \in C$  a different dictionary  $V^{(c)} \in \mathbb{R}^{d \times k^{(c)}}$  of  $k^{(c)}$  atoms and an encoding  $H^{(c)} \in \mathbb{R}^{k^{(c)} \times n}$  s.t.  $X^{(c)} = V^{(c)}H^{(c)} + \epsilon$  where  $\epsilon$  is the error introduced by the coding.

The dictionary forms (or the dictionaries form) the basis of our explanation framework for an ML system. Using the one-for-all or the one-for-one approach have different pros and cons. Intuitively, one-for-all gives a dictionary representing all data in an unsupervised manner. However, when data are too complex, this may not be the best choice. The one-for-one approach may generate dictionaries that are narrowly restricted to a single classes, but this is not a problem for our approach to explanation. In experiments we used both strategies, selecting the more suitable strategy for the data set complexity.

We selected as dictionary learning algorithm an NMF scheme [2] with the additional sparseness constraint proposed by [13]; this choice is motivated by the fact that it respects our requirements described above, giving a “local” representation of data, and *non-negativity*, that ensures only additive operations in data representations, giving a better human understanding with respect to other techniques. The sparsity level can be set using two parameters  $\gamma_1$  and  $\gamma_2$  which control the sparsity on the dictionary and the encoding, respectively.

## 2.2 Explanation Maximization

Unlike traditional dictionary-based coding approaches, our main goal is not to get an “accurate” representation of the input data, but to get a representation that helps humans to understand the decision taken by a trained model. To this aim, we modify the AM algorithm so that, instead of looking for the input that just maximizes the answer of the model, it searches for the dictionary-based encoding  $\vec{h}$  that maximizes the answer and, at the same time, is sparse enough but without being “too far” from the original input  $\vec{x}$ . More formally, indicating with  $\Pr(\hat{c}|\vec{x})$  the probability given by a learned model that input  $\vec{x}$  belongs to class  $\hat{c} \in C$ ,  $V$  the chosen dictionary,  $S(\cdot)$  a sparsity measure, the objective

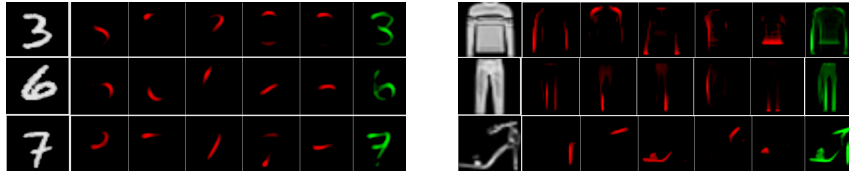


Fig. 1: Visual explanations obtained on three samples from the MNIST (left) and Fashion-Mnist (right) data sets correctly classified by the Oracle. In red are the meaningful parts determined by the systems producing explanations. In green are the encodings of the input image obtained from the sparse dictionary.

function that we optimise is

$$\max_{\vec{h} \geq 0} \log \Pr(\hat{c} | V\vec{h}) - \lambda_1 \|V\vec{h} - \vec{x}\|_2 + \lambda_2 S(\vec{h}) \quad (1)$$

where  $\lambda_1, \lambda_2$  are hyper-parameters regulating the input reconstruction and the encoding sparsity level, respectively. The first regularization term leads the algorithm to choose dictionary atoms that, with an appropriate encoding, form a good representation of the input, while the second regularization term ensures a certain sparsity degree, i.e., that only few atoms are used. The  $\vec{h} \geq 0$  constraint ensures that one has a purely additive encoding. Thus, each  $h_i, \forall i.1 \leq i \leq d$ , measures the “importance” of the  $i$ -th atom. Equation 1 is solved by using a standard gradient ascent technique, together with a projection operator given by [13] that ensures both sparsity and non-negativity.

### 3 Experimental Assessment

To test our framework, we chose as Oracle a convolutional neural network architecture, LeNet-5 [14], generally used for digit recognition as MNIST. We have trained the network from scratch using two different datasets: MNIST [14], obtaining an accuracy of 98.86% on the test set, and Fashion-MNIST [15], obtaining an accuracy of 91.43% on the test set. The training set is composed of 50000 images, while the test set is composed of 10000 images; the model is learned using the Adam algorithm [16].

NMF with sparseness constraints [13] is used to determine the dictionaries. We set the number of atoms to 200, relying on PCA analysis which showed that the first 100 principal components explain more than 95% of the data variance. We construct different dictionaries with different sparsity values in the range  $\gamma_1, \gamma_2 \in [0.6, 0.8]$  [13], then we choose the dictionaries having the best trade-off between sparsity level and reconstruction error.

In our experiments we learn a single dictionary for all the 10 classes of the MNIST data set, which is a relatively simple data set. For the Fashion-MNIST data set, which is more complex, we opted for learning one dictionary for each class. The dictionaries are determined by looking for a good trade-off between reconstruction error and sparsity level.

The atoms forming our explanations are selected by taking those with larger encoding values (i.e., those that are more “important” in the representation). In figure 1 (left) we show the atoms forming the explanation on three inputs on which the Oracle gave the correct answer. The chosen atoms seem to describe well the visual impact of the input numbers, by providing elements that appear to be discriminative, such as the curved elements that are present for “3” or the straight lines for “7”. To probe empirically the impact of sparsity on this representation, we performed the same experiment using a dictionary with a very low sparsity (0.1), obtaining encodings without any preponderant value, thereby making it difficult to select appropriate atoms for explanation.

For the Fashion-MNIST we found a specific dictionary for each class. Furthermore, in the dictionary learning procedure, we set the sparsity in an exclusive manner, i.e.  $\gamma_a \neq 0 \iff \gamma_b = 0, \forall a, b \in \{1, 2\}$ . This choice is motivated by the fact that having sparsity both on dictionary and encoding leads to poor atoms. So, for each class, we construct different dictionaries for different values on  $\gamma_1, \gamma_2$  and then we choose the one with a good trade off between reconstruction error and sparsity level.

In figure 1 (right) we show the more “important” atoms obtained on three input images, a *female sandal* and a *shirt*, all of them correctly classified by the Oracle. Selecting the atoms with higher encoding values seems to give rise to representative parts of the selected input, returning parts that can be easily interpreted by a human interrogator (e.g., the raised sole for the female sandal and the sleeves for the shirt). As for MNIST, we performed the same experiment using a dictionary with low sparsity, ending up with results that are difficult to interpret.

## 4 Conclusions

We proposed a model-agnostic framework to explain the answers given by classification systems. To achieve this objective, we started by defining a general explanation framework based on three entities: an Oracle (providing the answers to explain), an Interrogator (posing explanation requests) and a Mediator (helping Interrogator to interpret the Oracle’s decisions). We propose a Mediator using known and established techniques of sparse dictionary learning, together with Interpretability ML techniques, to give a humanly interpretable explanation of a classification system outcomes. We tried our proposed approach by using an NMF-based scheme as sparse dictionary learning technique. However, we expect that any other technique that meets the requirements outlined in Section 2 may be successfully used to instantiate the proposed framework. The results of the experiments that we carried out are encouraging, insofar as the explanations provided seem to be qualitatively significant. Nevertheless, more experiments are necessary to probe the general interest of our approach to explanation. We plan to perform both a quantitative assessment, to evaluate explanations by techniques such as those proposed in [6], and a subjective quality assessment to test how do humans perceive and interpret explanations of this kind.

The proposed approach does not take so far into account factors such as the internal structure of the dictionary used. Accordingly, the present work can be extended by considering, for example, whether there are atoms that are sufficiently “similar” to each other or whether the presence in the dictionary of atoms which can be expressed as combinations of other atoms may affect the explanations that are arrived at. Another interesting direction of research concerns contrastive explanations, which enable one to answer “why not?” negative questions, by explaining why some given object was not given another classification, differing from the classification that the Oracle actually provided.

## References

- [1] B. Letham, C. Rudin, T. H. McCormick, D. Madigan, et al. Interpretable classifiers using rules and bayesian analysis: Building a better stroke prediction model. *The Annals of Applied Statistics*, 9(3):1350–1371, 2015.
- [2] D. D. Lee and H. S. Seung. Algorithms for non-negative matrix factorization. In *Advances in neural information processing systems*, pages 556–562, 2001.
- [3] D. Doran, S. Schulz, and T. R. Besold. What does explainable ai really mean? a new conceptualization of perspectives. *arXiv preprint arXiv:1710.00794*, 2017.
- [4] D. Erhan, Y. Bengio, . Courville, and P. Vincent. Visualizing higher-layer features of a deep network. *University of Montreal*, 1341(3):1, 2009.
- [5] A. Nguyen, A. Dosovitskiy, J. Yosinski, T. Brox, and J. Clune. Synthesizing the preferred inputs for neurons in neural networks via deep generator networks. In *Advances in Neural Information Processing Systems*, pages 3387–3395, 2016.
- [6] G. Montavon, W. Samek, and K.R. Müller. Methods for interpreting and understanding deep neural networks. *Digital Signal Processing*, 73:1–15, 2018.
- [7] Z. C. Lipton. The mythos of model interpretability. *Queue*, 16(3):30:31–30:57, June 2018.
- [8] K. Simonyan, A. Vedaldi, and A. Zisserman. Deep inside convolutional networks: Visualising image classification models and saliency maps. *arXiv preprint arXiv:1312.6034*, 2013.
- [9] I. Sturm, S. Lapuschkin, W. Samek, and K.R. Müller. Interpretable deep neural networks for single-trial eeg classification. *Journal of neuroscience methods*, 274:141–145, 2016.
- [10] M. T. Ribeiro, S. Singh, and C. 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*, pages 1135–1144. ACM, 2016.
- [11] A. Mensch, J. Mairal, B. Thirion, and G. Varoquaux. Dictionary learning for massive matrix factorization. In *Proceedings of The 33rd International Conference on Machine Learning*, pages 1737–1746, 2016.
- [12] C. Bao, H. Ji, Y. Quan, and Z. Shen. Dictionary learning for sparse coding: Algorithms and convergence analysis. *IEEE transactions on pattern analysis and machine intelligence*, 38(7):1356–1369, 2016.
- [13] P. O. Hoyer. Non-negative matrix factorization with sparseness constraints. *Journal of machine learning research*, 5(Nov):1457–1469, 2004.
- [14] Y. Lecun, L. Bottou, Y. Bengio, and P. Haffner. Gradient-based learning applied to document recognition. *Proceedings of the IEEE*, 86(11):2278–2324, 1998.
- [15] H. Xiao, K. Rasul, and R. Vollgraf. Fashion-mnist: a novel image dataset for benchmarking machine learning algorithms. *CoRR*, abs/1708.07747, 2017.
- [16] D. Kingma and J. Ba. Adam: A method for stochastic optimization. In *International Conference on Learning Representations*, 12 2014.