

Toward a Unified Framework for Debugging Concept-based Models

Andrea Bontempelli,¹ Fausto Giunchiglia,^{1,2} Andrea Passerini,¹ Stefano Teso,¹

¹ University of Trento, Italy

² Jilin University, China
name.surname@unitn.it

Abstract

In this paper, we tackle interactive debugging of “gray-box” concept-based models (CBMs). These models learn task-relevant concepts appearing in the inputs and then compute a prediction by aggregating the concept activations. Our work stems from the observation that in CBMs *both* the concepts *and* the aggregation function can be affected by different kinds of bugs, and that fixing these bugs requires different kinds of corrective supervision. To this end, we introduce a simple schema for human supervisors to identify and prioritize bugs in both components, and discuss solution strategies and open problems. We also introduce a novel loss function for debugging the aggregation step that generalizes existing strategies for aligning black-box models to CBMs by making them robust to how the concepts change during training.

Introduction

A central tenet of eXplainable AI (XAI) is that explanations of a model’s predictions naturally uncover bugs and biases affecting the model (Lapuschkin et al. 2019; Schramowski et al. 2020). Post-hoc explanations of black-box models, however, can be unfaithful and ambiguous (Dombrowski et al. 2019; Teso 2019; Lakkaraju and Bastani 2020; Sixt, Granz, and Landgraf 2020), and the extraction process can be computationally challenging (Van den Broeck et al. 2021).

Concept-based models (CBMs) are designed to make the extraction step as straightforward as possible while retaining the performance of more opaque alternatives (Rudin 2019). To this end, CBMs learn a set of high-level, interpretable concepts capturing task-salient properties of the inputs, and then obtain predictions by aggregating the concept activations in a (typically) understandable manner (Alvarez-Melis and Jaakkola 2018; Losch, Fritz, and Schiele 2019; Koh et al. 2020; Chen et al. 2019; Hase et al. 2019; Rymarczyk et al. 2020; Nauta, van Bree, and Seifert 2021; Lage and Doshi-Velez 2020). A key feature of CBMs is that they explain their own predictions by supplying faithful concept-level attribution maps, which encompass both the concepts and the aggregation weights, thus facilitating the identification of bugs affecting the model.

Copyright © 2022, Association for the Advancement of Artificial Intelligence (www.aaai.org). All rights reserved.

Work on interactive troubleshooting of CBMs is, however, sparse and *ad hoc*: some approaches assume the concepts to be given and high-quality and focus on correcting the aggregation step (Teso 2019; Stammer, Schramowski, and Kersting 2021), others address issues with the learned concepts while ignoring how they are aggregated (Barnett et al. 2021; Lage and Doshi-Velez 2020). Fixing only one set of bugs is however insufficient.

In this paper, we outline a unified framework for debugging CBMs. Our framework stems from the simple observation that the quality of a CBM hinges on *both* the concepts vocabulary *and* on how the concepts are aggregated, and that both elements are conveniently captured by the CBM’s explanations. This immediately suggests a human-in-the-loop debugging strategy based on providing *supervision on the model’s explanations* that is composed of three different phases, namely (i) evaluating concept quality, (ii) correcting the aggregation weights, and (iii) correcting the concepts themselves. This novel yet intuitive setup allows us to identify limitations in existing works and highlight possible ways to overcome them.

As a first step toward implementing this framework, we introduce a new loss function on the aggregation weights that generalizes approaches for aligning local explanations (Ross, Hughes, and Doshi-Velez 2017; Lertvittayakumjorn, Specia, and Toni 2020) to be robust to changes in the underlying concepts during training. We then outline how the same strategy can be applied recursively to align the concepts themselves by correcting both their predictions and explanations.

Contributions. Summarizing, we:

1. Introduce a unified framework for debugging CBMs that explicitly distinguishes between and addresses bugs affecting how the concepts are defined and how they are aggregated.
2. Illustrate how to incorporate corrective feedback into the aggregation step in a manner that is invariant to changes to the learned concepts.
3. Discuss how to align the concepts by fixing the reasons behind their activations, opening the door to explanation-based debugging of the concepts themselves.

Concept-based Models

We are concerned with learning a high-quality classifier $f : \mathbf{x} \mapsto y$ that maps instances $\mathbf{x} \in \mathbb{R}^d$ into labels $y \in [v] := \{1, \dots, v\}$. In particular, we focus on *concept-based models* (CBMs) that fit the following two-level structure.

At the lower level, the model extracts a vector of *concept activations* $\mathbf{c}(\mathbf{x}) = (c_1(\mathbf{x}), \dots, c_k(\mathbf{x})) \in \mathbb{R}^k$ from the raw input \mathbf{x} . The concepts $\{c_j\}$ are usually learned from data so to provide strong indication of specific classes (Chen et al. 2019). For instance, in order to discriminate between images of cars and plants, the model might learn concepts that identify wheels and leaves. The concepts themselves are completely black-box.

At the upper level, the model *aggregates* the concept activations into per-class scores, typically in a simulatable (Lipton 2018) manner. This is most often implemented by taking a linear combination of the activations:

$$s_y(\mathbf{x}) := \langle \mathbf{w}^{(y)}(\mathbf{x}), \mathbf{c}(\mathbf{x}) \rangle = \sum_j w_j^{(y)}(\mathbf{x}) \cdot c_j(\mathbf{x}) \quad (1)$$

where $\mathbf{w}^{(y)}(\mathbf{x}) \in \mathbb{R}^k$ is the weight vector associated to class y . Class probabilities are then obtained by passing the scores through a softmax activation, that is, $P(y | \mathbf{x}) := \text{softmax}(\mathbf{s}(\mathbf{x}))_y$.

CBMs are learned by minimizing an empirical loss $1/|D| \cdot \sum_i \ell(f, (\mathbf{x}_i, y_i))$, where i runs over the training set D , as customary. The loss function $\ell(f, (\mathbf{x}, y))$ itself combines a standard loss for classification, like the cross-entropy loss, together with additional regularization terms that encourage the learned concepts to be understandable to (sufficiently expert) human stakeholders. Commonly used criteria include similarity to concrete examples (Alvarez-Melis and Jaakkola 2018; Chen et al. 2019), disentanglement (Alvarez-Melis and Jaakkola 2018), and boundedness (Koh et al. 2020).

With CBMs, it is straightforward to extract local explanations that capture how different concepts contribute to a decision (\mathbf{x}, y) . These explanations take the form:

$$\mathcal{E}(\mathbf{x}, y) := \{(w_j^{(y)}(\mathbf{x}), c_j(\mathbf{x})) : j \in [k]\} \quad (2)$$

Notice that *both the concepts and the aggregation weights are integral to the explanation*: the concepts $\{c_j\}$ establish a vocabulary that enables communication with stakeholders, while the weights $\{w_j(\mathbf{x})\}$ convey the relative importance of different concepts. Crucially, the score of class y is independent from \mathbf{x} given the explanation $\mathcal{E}(\mathbf{x}, y)$, ensuring that the latter is faithful to the model’s decision process.

Implementations

Next, we introduce some well-known CBMs that match the above template. A summary can be found in Table 1. Additional models are briefly discussed in the Related Work.

Self-Explainable Neural Networks (SENNs), introduced by Alvarez-Melis and Jaakkola (2018), generalize linear models to representation learning. SENNs instantiate the above two-level template as follows. The weight functions $\mathbf{w}^{(y)}(\mathbf{x})$, $y \in [v]$, are implemented using a neural network

for multivariate regression. Importantly, the weights are regularized so to vary slowly across inputs. This is achieved by penalizing the model for deviating from its first-order (i.e., linear) Taylor decomposition. The concepts $\mathbf{c}(\mathbf{x})$ are either given and fixed or learned using a sparsity-regularized auto-encoder, and once learned they are conveyed to users by presenting a handful of concrete training examples on which the concepts maximally activate. All components of SENNs are training jointly in an end-to-end fashion.

Part Prototype Networks (ProtoPNets) (Chen et al. 2019) ground the two-level template to image classification as follows. The weights $\mathbf{w}^{(y)}$ are constant with respect to the input \mathbf{x} , while the concepts $\mathbf{c}(\mathbf{x})$ indicate the presence of “part-prototypes”, that is, prototypes that capture specific parts of images appearing in the training set. Specifically, the part-prototypes are implemented as (parts of) points in the embedding space of a (pre-trained) convolutional neural network. Each part-prototype \mathbf{p} is encoded as a $1 \times 1 \times q$ parameter vector and activates based on the contents of a corresponding rectangular receptive field in input space. Each class $y \in [v]$ is associated to its own $\lfloor \frac{k}{v} \rfloor$ part-prototypes, which are learned so as to maximally activate on (parts of) training images of the associated class and not to activate on those of the other classes. The concepts and the weights are learned sequentially.

IAIA-BL (Barnett et al. 2021) specializes ProtoPNets to image-based medical diagnosis. In particular, in addition to label annotations, IAIA-BL accepts per-example attribute relevance information (e.g., annotations of symptomatic regions in X-ray images) and penalizes part-prototypes that activate outside the relevant areas.

Concept Bottleneck Models (CBNMs) (Koh et al. 2020; Losch, Fritz, and Schiele 2019) are regular feed-forward neural networks in which the neurons in a given layer are trained to align with a known vocabulary of concepts $\mathbf{c}(\mathbf{x})$. This is achieved by supplying the CBNM with concept-level annotations. The aggregation step is not particularly restrained: the concept neurons are either using a single dense layer, in which case the weights $\mathbf{w}^{(y)}$ are constant with respect to \mathbf{x} , or through a sequence of layers, in which case the weights $\mathbf{w}^{(y)}(\mathbf{x})$ are *not* constant and must be inferred using post-hoc techniques (e.g., input gradients).

Existing Strategies for Debugging CBMs

Existing literature on debugging CBMs can be split into two groups. Approaches in the first group are concerned with bugs in the aggregation weights $\{\mathbf{w}^{(y)}(\mathbf{x})\}_y$. These can occur when the data fools the model into assigning non-zero weight to concepts that correlate with – but are not causal for – the label. A prototypical example are class-specific watermarks in image classification (Lapuschkin et al. 2019). Teso (2019) addresses this issue by leveraging explanatory interactive learning (XIL) (Teso and Kersting 2019; Schramowski et al. 2020). As in active learning, in XIL the machine obtains labels by querying a human annotator, but in addition it presents predictions for its queries and local

METHOD	CONCEPTS $\mathbf{c}(\mathbf{x})$	AGGREGATOR $f(\mathbf{c})$	EXTRA ANNOT.	TRAINING
SENN (Alvarez-Melis and Jaakkola 2018)	Autoencoder	Linear Comb.	–	End-to-end
ProtoPNet (Chen et al. 2019)	Conv. Filters	Linear Comb.	–	Multistep
IAIA-BL (Barnett et al. 2021)	Conv. Filters	Linear Comb.	Concept Attr.	Multistep
CBM (Koh et al. 2020)	Arbitrary	Arbitrary	Concept Labels	End-to-end

Table 1: Comparison between concept-based CBMs considered in this work.

explanations for its predictions. The user then provides corrective feedback on the explanations, for instance by indicating those parts of an image that are irrelevant to the class label but that the machine relies on for its prediction. The model is then penalized whenever its weights do not align with the user’s corrections. This setup assumes that the concepts \mathbf{c} are fixed rather than learned from data. Stammer, Schramowski, and Kersting (2021) extend XIL to neuro-symbolic models and structured attention, but also assume the concepts to be fixed.

Conversely, approaches in the second group are only concerned with how the concepts are defined: concepts learned from data, even if discriminative and interpretable, may be misaligned with (the stakeholders’s understanding of) the prediction task and may thus fail to generalize properly. CB-NMs work around this issue by leveraging concept-level label supervision (Koh et al. 2020), however this does not ensure that the concepts themselves are “right for the right reasons” (Ross, Hughes, and Doshi-Velez 2017). As a matter of fact, just like all other neural nets (Szegedy et al. 2013), part-prototypes learned by ProtoPNets can pick up and exploit uninterpretable features of the input (Hoffmann et al. 2021). Hoffmann et al. (2021) and Nauta et al. (2020) further argue that it may be difficult for users to understand what a learned prototype (concept) represents, unless the model explains *why* the concept activates. To this end, Nauta et al. (2020) propose a perturbation-based technique – analogous to LIME (Ribeiro, Singh, and Guestrin 2016) – for explaining why a particular prototype activates on a certain region of an image, which however does nothing to *fix* the bugs that it highlights. Finally, Barnett et al. (2021) introduce a loss term for ProtoPNets that penalizes concepts that activate on regions annotated as irrelevant by a domain expert.

A Unified Framework for Debugging CBMs

We propose a new unified framework – based on the right for the right reasons (RRR) principle (Ross, Hughes, and Doshi-Velez 2017) – that, in contrast with existing solutions, is designed for the more realistic case of CBMs affected by multiple, different bugs.

The RRR principle is straightforward: the goal is to ensure that the model outputs accurate *predictions* that are justified by high-quality *explanations*. In order to make this concrete, we need to define what we mean for an explanation $\mathcal{E}(\mathbf{x}, y)$ to be “high-quality”. Intuitively, for CBMs explanation quality depends on the quality of the learned concepts \mathbf{c} and aggregation weights $\mathbf{w}(\mathbf{x})$. We unpack this intuition as follows:

Definition 1. *A set of concepts \mathbf{c} is high-quality for a deci-*

sion (\mathbf{x}, y) if:

- C1.** *The set of concepts is sufficient to fully determine the ground-truth label y^* from \mathbf{x} .¹*
- C2.** *The various concepts are (approximately) independent from each other.*
- C3.** *Each concept is semantically meaningful.*
- C4.** *Each concept is easy to interpret (e.g., simple enough).*

Given high-quality concepts \mathbf{c} , a set of weights $\mathbf{w}^{(y)}(\mathbf{x})$ is high quality for a decision (\mathbf{x}, y) if:

- W1.** *It ranks all concepts in $\mathbf{c}(\mathbf{x})$ compatibly with their relevance for the prediction task.*
- W2.** *It associates (near-)zero weight to concepts in $\mathbf{c}(\mathbf{x})$ that are task-irrelevant.*

An explanation $\mathcal{E}(\mathbf{x}, y)$ is high-quality for a decision (\mathbf{x}, y) if both \mathbf{c} and $\mathbf{w}^{(y)}(\mathbf{x})$ are.

It is worth discussing the various points in detail. Requirement C1 ensures that the concepts are task-relevant and jointly sufficient to solve the prediction task. This requirement is imposed by the learning problem itself. Notice that, since in all CBMs the prediction $f(\mathbf{x})$ is independent from the input \mathbf{x} given the explanation $\mathcal{E}(\mathbf{x}, f(\mathbf{x}))$, C1 is necessary for the model to achieve good generalization. Although some works stress the need for the concept vocabulary to be complete (Yeh et al. 2020; Bahadori and Heckerman 2021), we argue that – when it comes to local explanations – sufficiency is more relevant, although it is clear that if the concepts \mathbf{c} are sufficient for all instances \mathbf{x} , then they do form a complete vocabulary.

Requirements C2, C3, and C4, on the other hand, are concerned with whether the concept set can be used for *communicating* with human stakeholders, and lie at the core of CBMs and human-in-the-loop debugging. Notice also that requirements C1–C4 are mutually independent: not all task-relevant concepts are understandable (indeed, this is often not the case), not all semantically meaningful concepts are easy to interpret, *etc.*

Debugging CBMs in Three Steps

Assume to be given a decision (\mathbf{x}, y) and an explanation $\mathcal{E}(\mathbf{x}, y)$ that is *not* high quality. What bugs should the user prioritize? We propose a simple three-step procedure:

- Step 1** *Evaluating concept quality:* Determine if $\mathcal{E}(\mathbf{x})$ contains a high-quality subset $\mathbf{c}' \subseteq \mathbf{c}$ that is *sufficient* to produce a correct prediction y^* for the target instance \mathbf{x} .

¹In case the ground-truth label is ill-defined, it should be replaced with the Bayes optimal label.

Step 2 Correcting the aggregation weights: If so, then it is enough to fix how the model combines the available concepts by supplying corrective supervision for the aggregation weights \mathbf{w} .

Step 3 Correcting the learned concepts: Otherwise, it is necessary to create a high-quality subset \mathbf{c}' by supplying appropriate supervision on the concepts \mathbf{c} themselves.

In the following we discuss how to implement these steps.

Step 1: Assessing concept quality

In this step, one must determine whether there exists a subset $\mathbf{c}' \subseteq \mathbf{c}$ that is high-quality. This necessarily involves conveying the learned concepts \mathbf{c} to a human expert in detail sufficient to determine whether they are “good enough” for computing the ground-truth label y^* from \mathbf{x} .

The most straightforward solution, adopted for instance by SENNs (Alvarez-Melis and Jaakkola 2018) and ProtoP-Nets (Chen et al. 2019), is to present a set of instances that are most *representative* of each concept. In particular, ProtoPNets identify parts of training instances \mathbf{x} that maximally activate each $c_j \in \mathbf{c}$. A more refined way to characterize a concept is to explain *why* it activates for certain instances, possibly the prototypes themselves (Nauta et al. 2020; Hoffmann et al. 2021). Such explanations can be obtained by extracting an attribution map $\text{attr}(c_j, \mathbf{x})$ that identifies those inputs – either raw inputs \mathbf{x} or higher-level features of \mathbf{x} like color, shape, or texture – that maximally contribute to the concept’s activations. Since concepts are black-box, this map must be acquired using post-hoc techniques, like input gradients (Baehrens et al. 2010; Sundararajan, Taly, and Yan 2017) or LIME (Ribeiro, Singh, and Guestrin 2016). Explanations are especially useful to prevent stakeholders from confusing one concept for another. For instance, in a shape recognition task, a part-prototype that activates on a yellow square might do so because of the shape, of the color, or both. Without a proper explanation, the user cannot tell these concepts apart. This is fundamental for preventing users from wrongly trusting ill-behaved concepts (Teso and Kersting 2019; Rudin 2019).

Although assessing concept quality does *not* require to communicate the full semantics of a learned concept to the human counterpart, acquiring feedback on the aggregation weights *does*. Doing so is non-trivial (Zhang et al. 2019), as concepts learned by CBMs are not always interpretable (Nauta et al. 2020; Hoffmann et al. 2021). However, we stress that uninterpretable concepts are not high-quality, and therefore they must be dealt with in Step 3, i.e., they must be improved by supplying appropriate concept-level supervision.

Step 2: Fixing how the concepts are used

In this second step, the goal is to fix up how the concepts are aggregated. This case is reminiscent of debugging black-box models. Existing approaches for this problem acquire a (possibly partial (Teso and Kersting 2019; Teso 2019)) ground-truth attribution map $\mathbf{m} \in \{0, 1\}^d$ that specifies what input attributes are (ir)relevant for the target prediction, and

then penalize the model for allocating non-zero relevance to them.

More specifically, let f be a *black-box* model and $\text{attr}(f, (\mathbf{x}, y)) \in \mathbb{R}^d$ be an attribution mechanism that assigns numerical responsibility for the decision (\mathbf{x}, y) to each input x_i , for $i = 1, \dots, d$, for instance integrated gradients (Sundararajan, Taly, and Yan 2017). The model’s explanations are corrected by introducing a loss of the form (Ross, Hughes, and Doshi-Velez 2017; Schramowski et al. 2020; Shao et al. 2021):²

$$\ell_{\text{attr}}(f, (\mathbf{x}, y), \mathbf{m}) := \sum_{i \in [d]} (1 - m_i) \cdot \text{attr}_i(f, (\mathbf{x}, y))^2 \quad (3)$$

Now, consider a CBM f and a decision (\mathbf{x}, y) obtained by aggregating *high-quality* concepts \mathbf{c} using *low-quality* weights $\mathbf{w}(\mathbf{x})$. It is easy to see how Eq. 3 could help with aligning the aggregation weights. However, simply replacing attr_i with the weights \mathbf{w} does not work. The reason is that in CBMs the concepts *change* during training, and so do the semantics of their associated weights. Hence, feedback of the form “don’t use the j -th concept” becomes obsolete (and misleading) whenever the j -th concept changes.

Another major problems with Eq. 3 is that, by penalizing concepts by their index, it does not prevent the model from re-learning a forbidden concept under a different index. Consider a toy image classification task in which the positive class consists of images that contains a green circle (50% of the images) or a pink triangle (the other 50%), but the data is confounded such that 100% of the positive training images also contain a yellow square. A ProtoPNet with a budget of two prototypes per class is encouraged to rely on the confounder to achieve 100% accuracy on the training set. If the model is penalized with Eq. 3 for using the confounder as part-prototype 1, it still has plenty of room to learn the confounder using the *other* prototype of the positive class.

Roughly speaking, this shows that CBMs are too flexible for regular alignment losses and can therefore easily work around them. Our solution to this problem is presented in the next Section.

Step 3: Fixing how the concepts are defined

Finally, consider a decision (\mathbf{x}, y) that depends on a *low-quality* set of concepts \mathbf{c} . In this case, the goal is to ensure that at least some of those concepts quickly become useful by supplying additional supervision.

The order in which the various concepts should be aligned is left to the user. This is reasonable under the assumption that she is a domain expert and that those concepts that are worth fixing are easy to distinguish from those that are not. An alternative is to debug the concepts sequentially based on their overall impact on the model’s behavior, for instance sorting them by decreasing relevance $|w_j^{(y)}(\mathbf{x})|$. More refined strategies will be considered in future work.

²The loss can be adapted to also encourage the model to rely on concepts that are deemed relevant, for instance by penalizing the model whenever the concept’s weight is too close to zero.

Now, let c_j be a low-quality concept. Our key insight is that there is no substantial difference between the models’ output $f(\mathbf{x})$ and a specific concept $c_j(\mathbf{x})$ appearing in it.³ This means that work on understanding and correcting predictions of black-box models can be applied for understanding and correcting concepts in CBMs – with some differences, see below. For instance, the work of (Nauta et al. 2020) can be viewed as a concept-level version of LIME (Ribeiro, Singh, and Guestrin 2016). To the best of our knowledge, this useful analogy has never been pointed out before.

A direct consequence is that concepts, just like predictions, can be aligned by providing labels for them, as is done by CBMs (Koh et al. 2020). One issue with this strategy is that, if these annotations are biased, the learned concept may end up relying on confounders (Lapuschkin et al. 2019). A more robust alternative, that builds on approaches for correcting the model’s explanations (Ross, Hughes, and Doshi-Velez 2017; Teso and Kersting 2019), is to also align the concepts’ *explanations*. This involves extracting an explanation $\text{attr}(c_j, \mathbf{x})$ that uncovers the reasons behind the concept’s activations in terms of either inputs, as done in (Barnett et al. 2021), or higher-level features, and supplying corrective feedback for them. One caveat is that these strategies are only feasible when the semantics of c_j are already quite clear to the human annotator (e.g., the concept clearly captures “leaves” or “wheels”, rather than seemingly random but class-discriminative blobs, as can occur during training). If this is not the case, then the only option is to instruct the model to avoid using it by using ℓ_{aggr} or analogous penalties.

On-demand “conceptification”. Since concepts in CBMs are black-box, explanations for them must be extracted using post-hoc attribution techniques, which – as we argued in the introduction – is less than optimal. A better alternative is to model the concepts themselves as CBMs, making cheap and faithful explanations *for the concepts* readily available. We call the idea of replacing a black-box model with an equivalent CBM “*conceptification*”. Said CBM can be obtained using, for instance, model distillation (Gou et al. 2021), and it would persist over time. A nice benefit of conceptification is that all debugging techniques discussed so far would immediately become available for the newly introduced concepts too.

Conceptification makes the most sense in applications where concepts can be computed from simpler, lower-level concepts (like shape, color, and texture (Nauta et al. 2020)) or where concepts exhibit a hierarchical, parts-of-parts structure (Fidler and Leonardis 2007). Notably, conceptification needs not be applied to all concepts unconditionally. One interesting direction of research is to let the user decide what concepts should be conceptified – for instance those that need to be corrected more often, or for which supervision can be easily provided.

³A similar point was brought up in (Stammer, Schramowski, and Kersting 2021).

The Debugging Loop

Summarizing, interactive debugging of CBMs involves repeatedly choosing an instance (\mathbf{x}, y) and then proceeding as in steps 1–3, acquiring corrective supervision on weights and concepts along the way. The target decisions can be either selected by the machine, as in explanatory active learning (Teso and Kersting 2019), or selected by a human stakeholder, perhaps aided by the machine (Popordanoska, Kumar, and Teso 2020). Since in CBMs concepts and weights are learned jointly, the model is retrained to permit supervision to flow to both concepts and weights, further aligning them. This retraining stage proceeds until the model becomes stable enough in terms of, e.g., prediction or explanation accuracy on the training set or on a separate validation set. At this point, the annotator can either proceed with another debugging session or stop if the model looks good enough.

This form of debugging is quite powerful, as it allows stakeholders to identify and correct a variety of different bugs, but it is not universal: some bugs and biases that cannot be uncovered using local explanations (Popordanoska, Kumar, and Teso 2020). Other bugs – like those due to bad choice of hyper-parameters, insufficient model capacity, and failure to converge – are not fixable with any form of supervision. Dealing with these issues, however, is beyond the scope of this paper and left to future work.

A Robust Aggregation Loss

Simply discouraging the model from using incorrectly learned concepts does not prevent it from learning them again with minimal changes. To address this problem, we propose to penalize the model for associating non-zero weight to concepts that are *similar* to those concepts c indicated as irrelevant to the target decision during the debugging session. We assume these concepts to have been collected in a memory M . The resulting loss is:

$$\ell_{\text{aggr}}(f, (\mathbf{x}, y), M) := \sum_{\substack{c \in M(\mathbf{x}, y) \\ j \in [k]}} w_j^{(y)}(\mathbf{x})^2 \cdot \kappa(c, c_j) \quad (4)$$

The sum iterates over old concepts c irrelevant to (\mathbf{x}, y) , denoted $M(\mathbf{x}, y)$, and over the current concepts c , while $0 \leq \kappa(c, c') \leq 1$ measures the similarity between concepts.

Notice how Eq. 4, being independent of the concept’s order, prevents the model from re-learning forbidden concepts. Moreover, it automatically deals with redundant concepts, which are sometimes learned in practice (Chen et al. 2019), in which cases all copies of the same irrelevant concept are similarly penalized.

Measuring concept similarity. A principled way of doing so is to employ a kernel between functions. One natural approach is to employ a product kernel (Jebara, Kondor, and Howard 2004), for instance:

$$\kappa_{\text{act}}(c, c') := \int_{\mathbb{R}^d} (c(\mathbf{x}) \cdot c'(\mathbf{x}))^\rho p^*(\mathbf{x}) d\mathbf{x} \quad (5)$$

Here, $p^*(\mathbf{x})$ is the ground-truth distribution over instances and $\rho > 0$ controls the smoothness of the kernel. It is easy

to show that κ_{act} is a valid kernel by rewriting it as an inner product $\langle (q^*(\mathbf{x})c(\mathbf{x}))^\rho, (q^*(\mathbf{x})c'(\mathbf{x}))^\rho \rangle$, where we set $q^*(\mathbf{x}) := p^*(\mathbf{x})^{1/(2\rho)}$.

This kernel measures how often two concepts co-activate on different inputs, which is somewhat limiting. For instance, when discriminating between images of cars and plants, the concepts “wheel” and “license plate” are deemed similar because they co-activate on car images and not on plant images. This suggests using a more fine-grained kernel that considers *where* the concepts activate, for instance:

$$\kappa_{\text{attr}}(c, c') := \int_{\mathbb{R}^d} \langle \text{attr}(c, \mathbf{x}), \text{attr}(c', \mathbf{x}) \rangle^\rho p^*(\mathbf{x}) d\mathbf{x} \quad (6)$$

Since co-localization entails co-activation, this kernel specializes κ_{act} , in the sense that $\kappa_{\text{attr}}(c, c') \leq \kappa_{\text{act}}(c, c')$. Unfortunately, p^* is not known and the integrals in Eqs. 5 and 6 are generally intractable, so in practice we replace the latter with a Monte Carlo approximation on the training data. Notice that the activation and localization maps of irrelevant concepts in M can be cached for future use, thus speeding up the computation.

Certain models may admit more efficient kernels between concepts. For instance, in ProtoPNets one can measure concept similarity by directly comparing the convolutional filter representations \mathbf{p} of the learned concepts, e.g., $\langle \mathbf{p}, \mathbf{p}' \rangle$. Using this kernel gives the following aggregation loss:

$$\sum_{\substack{\mathbf{p} \in M(\mathbf{x}, y) \\ j \in [k]}} \langle \mathbf{p}, \mathbf{p}_j \rangle \cdot w_j(\mathbf{x})^2 = \left\langle \sum_{\mathbf{p} \in M(\mathbf{x}, y)} \mathbf{p}, \sum_{j \in [k]} w_j(\mathbf{x})^2 \cdot \mathbf{p}_j \right\rangle \quad (7)$$

which dramatically simplifies the computation. However, the intuitive meaning of inner products in parameter space is not entirely obvious, rendering this kernel harder to control. A thorough analysis of this option is left to future work.

Beyond instance-level supervision. Eq. 4 essentially encourages the distribution encoded by the CBM not to depend on a known-irrelevant concept. In previous work, relevance annotations are provided at the level of individual examples and, as such, specify an invariance that only holds for few similar examples (Teso and Kersting 2019; Schramowski et al. 2020; Lertvittayakumjorn, Specia, and Toni 2020; Shao et al. 2021; Stammer, Schramowski, and Kersting 2021).

Users, however, may find it more natural to provide relevance information for entire classes or for individual concepts. For instance, although the concept of “snow” is not relevant for a specific class “wolf”, it may be useful for other classes, like “winter”. Similarly, some concepts – for instance, artifacts in X-ray images due to the scanning process – should never be used for making predictions for any class. These more general forms of supervision encode invariances that applies to entire *sets* of examples, and as such have radically different expressive power. A nice property of Eq. 4 is that it can easily encode such supervision by letting the outer summation range over (or at least sample) sets of instances, i.e., all instances that belong to a certain class or all instances altogether.

Related Work

Our work targets well-known concept-based CBMs, including SENNs (Alvarez-Melis and Jaakkola 2018), CBMs (Koh et al. 2020; Losch, Fritz, and Schiele 2019), and ProtoPNets (Li et al. 2018; Chen et al. 2019; Hase et al. 2019; Rymarczyk et al. 2020; Barnett et al. 2021) and related approaches (Nauta, van Bree, and Seifert 2021), but it could be easily extended to similar models and techniques, for instance concept whitening (Chen, Bei, and Rudin 2020). The latter is reminiscent of CBMs, except that the neurons in the bottleneck layer are normalized and decorrelated and feedback on individual concepts is optional. Of course, not all CBMs fit our template. For instance, (Al-Shedivat, Dubey, and Xing 2020) propose a Bayesian model that, although definitely gray-box, does not admit selecting a unique explanation for a given prediction. More work is needed to capture these additional cases.

Our debugging approach is rooted in explanatory interactive learning, a set of techniques that acquire corrective from a user by displaying explanations (either local (Teso and Kersting 2019; Selvaraju et al. 2019; Lertvittayakumjorn, Specia, and Toni 2020; Schramowski et al. 2020; Shao et al. 2021) or global (Popordanoska, Kumar, and Teso 2020)) of the model’s beliefs. Recently, (Stammer, Schramowski, and Kersting 2021) designed a debugging approach for the specific case of attention-based neuro-symbolic models. These approaches assume the concept vocabulary to be given rather than learned and indeed are not robust to changes in the concepts, as illustrated above. Our aggregation loss generalizes these ideas to the case of concept-based CBMs. Other work on XIL has focused on example-based explanations (Teso et al. 2021; Zylberajch, Lertvittayakumjorn, and Toni 2021), which we did not include in our framework, but that could provide an alternative device for controlling a CBM’s reliance on, e.g., noisy examples.

A causal approach for debiasing CBMs has been proposed by (Bahadori and Heckerman 2021). This work is orthogonal to our contribution, and could be indeed integrated with it. The strategy of (Lage and Doshi-Velez 2020) acquires concept-based *white* box models that are better aligned with the user’s mental model by interactively acquiring concept-attribute dependency information. The machine does so by asking questions like “does depression depend on lorazepam?”, acquiring more impactful dependencies first. This approach is tailored for specific white-box models, but it could and should be extended to CBMs and integrated into our framework.

Conclusion

We proposed a unified framework for debugging concept-based CBMs that fixes bugs by acquiring corrective feedback from a human supervisor. Our key insight is that bugs can affect both how the concepts are defined and how they are aggregated, and that both elements have to be high-quality for a CBM to be effective. We proposed a three-step procedure to achieve this, shown how existing attribution

losses are unsuitable for CBMs and proposed a new loss that is robust to changes in the learned concepts, and illustrated how the same schema can be used to correct the concepts themselves by interacting with their labels and explanations. A thorough empirical validation of these ideas is currently underway.

Acknowledgments

We are grateful to the anonymous reviewers for helping us to improve the manuscript. This research has received funding from the European Union’s Horizon 2020 FET Proactive project “WeNet - The Internet of us”, grant agreement No. 823783, and from the “DELPhi - Discovering Life Patterns” project funded by the MIUR Progetti di Ricerca di Rilevante Interesse Nazionale (PRIN) 2017 – DD n. 1062 del 31.05.2019. The research of ST and AP was partially supported by TAILOR, a project funded by EU Horizon 2020 research and innovation programme under GA No 952215.

References

- Al-Shedivat, M.; Dubey, A.; and Xing, E. P. 2020. Contextual Explanation Networks. *J. Mach. Learn. Res.*, 21: 194–1.
- Alvarez-Melis, D.; and Jaakkola, T. S. 2018. Towards robust interpretability with self-explaining neural networks. In *Proceedings of the 32nd International Conference on Neural Information Processing Systems*, 7786–7795.
- Baehrens, D.; Schroeter, T.; Harmeling, S.; Kawanabe, M.; Hansen, K.; and Müller, K.-R. 2010. How to explain individual classification decisions. *The Journal of Machine Learning Research*, 11: 1803–1831.
- Bahadori, M. T.; and Heckerman, D. 2021. Debiasing Concept-based Explanations with Causal Analysis. In *International Conference on Learning Representations*.
- Barnett, A. J.; Schwartz, F. R.; Tao, C.; Chen, C.; Ren, Y.; Lo, J. Y.; and Rudin, C. 2021. IAIA-BL: A Case-based Interpretable Deep Learning Model for Classification of Mass Lesions in Digital Mammography. *arXiv preprint arXiv:2103.12308*.
- Chen, C.; Li, O.; Tao, D.; Barnett, A.; Rudin, C.; and Su, J. K. 2019. This Looks Like That: Deep Learning for Interpretable Image Recognition. *Advances in Neural Information Processing Systems*, 32: 8930–8941.
- Chen, Z.; Bei, Y.; and Rudin, C. 2020. Concept whitening for interpretable image recognition. *Nature Machine Intelligence*, 2(12): 772–782.
- Dombrowski, A.-K.; Alber, M.; Anders, C.; Ackermann, M.; Müller, K.-R.; and Kessel, P. 2019. Explanations can be manipulated and geometry is to blame. *Advances in Neural Information Processing Systems*, 32: 13589–13600.
- Fidler, S.; and Leonardis, A. 2007. Towards scalable representations of object categories: Learning a hierarchy of parts. In *2007 IEEE Conference on Computer Vision and Pattern Recognition*, 1–8. IEEE.
- Gou, J.; Yu, B.; Maybank, S. J.; and Tao, D. 2021. Knowledge distillation: A survey. *International Journal of Computer Vision*, 129(6): 1789–1819.
- Hase, P.; Chen, C.; Li, O.; and Rudin, C. 2019. Interpretable image recognition with hierarchical prototypes. In *Proceedings of the AAAI Conference on Human Computation and Crowdsourcing*, volume 7, 32–40.
- Hoffmann, A.; Fanconi, C.; Rade, R.; and Kohler, J. 2021. This Looks Like That... Does it? Shortcomings of Latent Space Prototype Interpretability in Deep Networks. *arXiv preprint arXiv:2105.02968*.
- Jebara, T.; Kondor, R.; and Howard, A. 2004. Probability product kernels. *The Journal of Machine Learning Research*, 5: 819–844.
- Koh, P. W.; Nguyen, T.; Tang, Y. S.; Musmann, S.; Pierson, E.; Kim, B.; and Liang, P. 2020. Concept bottleneck models. In *International Conference on Machine Learning*, 5338–5348. PMLR.
- Lage, I.; and Doshi-Velez, F. 2020. Learning Interpretable Concept-Based Models with Human Feedback. *arXiv preprint arXiv:2012.02898*.
- Lakkaraju, H.; and Bastani, O. 2020. “How do I fool you?” Manipulating User Trust via Misleading Black Box Explanations. In *Proceedings of the AAAI/ACM Conference on AI, Ethics, and Society*, 79–85.
- Lapuschkin, S.; Wäldchen, S.; Binder, A.; Montavon, G.; Samek, W.; and Müller, K.-R. 2019. Unmasking Clever Hans predictors and assessing what machines really learn. *Nature communications*, 10(1): 1–8.
- Lertvittayakumjorn, P.; Specia, L.; and Toni, F. 2020. FIND: human-in-the-loop debugging deep text classifiers. In *Conference on Empirical Methods in Natural Language Processing*, 332–348.
- Li, O.; Liu, H.; Chen, C.; and Rudin, C. 2018. Deep learning for case-based reasoning through prototypes: A neural network that explains its predictions. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 32.
- Lipton, Z. C. 2018. The Mythos of Model Interpretability. *Queue*, 16(3): 31–57.
- Losch, M.; Fritz, M.; and Schiele, B. 2019. Interpretability beyond classification output: Semantic Bottleneck Networks. *arXiv preprint arXiv:1907.10882*.
- Nauta, M.; Jutte, A.; Provoost, J.; and Seifert, C. 2020. This Looks Like That, Because... Explaining Prototypes for Interpretable Image Recognition. *arXiv preprint arXiv:2011.02863*.
- Nauta, M.; van Bree, R.; and Seifert, C. 2021. Neural Prototype Trees for Interpretable Fine-grained Image Recognition. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 14933–14943.
- Popordanoska, T.; Kumar, M.; and Teso, S. 2020. Machine guides, human supervises: Interactive learning with global explanations. *arXiv preprint arXiv:2009.09723*.
- Ribeiro, M. T.; Singh, S.; and Guestrin, C. 2016. “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*, 1135–1144.

- Ross, A. S.; Hughes, M. C.; and Doshi-Velez, F. 2017. Right for the right reasons: training differentiable models by constraining their explanations. In *Proceedings of the 26th International Joint Conference on Artificial Intelligence*, 2662–2670.
- Rudin, C. 2019. Stop explaining black box machine learning models for high stakes decisions and use interpretable models instead. *Nature Machine Intelligence*, 1(5): 206–215.
- Rymarczyk, D.; Struski, Ł.; Tabor, J.; and Zieliński, B. 2020. ProtoPSHare: Prototype Sharing for Interpretable Image Classification and Similarity Discovery. *arXiv preprint arXiv:2011.14340*.
- Schramowski, P.; Stammer, W.; Teso, S.; Brugger, A.; Herbert, F.; Shao, X.; Luigs, H.-G.; Mahlein, A.-K.; and Kersting, K. 2020. Making deep neural networks right for the right scientific reasons by interacting with their explanations. *Nature Machine Intelligence*, 2(8): 476–486.
- Selvaraju, R. R.; Lee, S.; Shen, Y.; Jin, H.; Ghosh, S.; Heck, L.; Batra, D.; and Parikh, D. 2019. Taking a hint: Leveraging explanations to make vision and language models more grounded. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, 2591–2600.
- Shao, X.; Skryagin, A.; Schramowski, P.; Stammer, W.; and Kersting, K. 2021. Right for Better Reasons: Training Differentiable Models by Constraining their Influence Function. In *Proceedings of Thirty-Fifth AAAI Conference on Artificial Intelligence (AAAI)*.
- Sixt, L.; Granz, M.; and Landgraf, T. 2020. When explanations lie: Why many modified bp attributions fail. In *International Conference on Machine Learning*, 9046–9057. PMLR.
- Stammer, W.; Schramowski, P.; and Kersting, K. 2021. Right for the Right Concept: Revising Neuro-Symbolic Concepts by Interacting with their Explanations. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 3619–3629.
- Sundararajan, M.; Taly, A.; and Yan, Q. 2017. Axiomatic attribution for deep networks. In *International Conference on Machine Learning*, 3319–3328. PMLR.
- Szegedy, C.; Zaremba, W.; Sutskever, I.; Bruna, J.; Erhan, D.; Goodfellow, I.; and Fergus, R. 2013. Intriguing properties of neural networks. *arXiv preprint arXiv:1312.6199*.
- Teso, S. 2019. Toward Faithful Explanatory Active Learning with Self-explainable Neural Nets. In *Proceedings of the Workshop on Interactive Adaptive Learning (IAL 2019)*, 4–16.
- Teso, S.; Bontempelli, A.; Giunchiglia, F.; and Passerini, A. 2021. Interactive Label Cleaning with Example-based Explanations. *arXiv preprint arXiv:2106.03922*.
- Teso, S.; and Kersting, K. 2019. Explanatory interactive machine learning. In *Proceedings of the 2019 AAAI/ACM Conference on AI, Ethics, and Society*, 239–245.
- Van den Broeck, G.; Lykov, A.; Schleich, M.; and Suciu, D. 2021. On the tractability of SHAP explanations. In *Proceedings of AAAI*.
- Yeh, C.-K.; Kim, B.; Arik, S.; Li, C.-L.; Pfister, T.; and Ravikumar, P. 2020. On Completeness-aware Concept-Based Explanations in Deep Neural Networks. *Advances in Neural Information Processing Systems*, 33.
- Zhang, Z.; Singh, J.; Gadiraju, U.; and Anand, A. 2019. Dissonance between human and machine understanding. *Proceedings of the ACM on Human-Computer Interaction*, 3(CSCW): 1–23.
- Zylberajch, H.; Lertvittayakumjorn, P.; and Toni, F. 2021. HILDIF: Interactive Debugging of NLI Models Using Influence Functions. *Workshop on Interactive Learning for Natural Language Processing*, 1.