

Document Version

Final published version

Citation (APA)

Lemeire, J., & Buijsman, S. (2025). Defining and Evaluating the Degrees of Abstraction in Explanations with Kolmogorov Complexity. In F. A. Oliehoek, M. Kok, & S. Verwer (Eds.), *Artificial Intelligence and Machine Learning: 35th Benelux Conference, BNAIC/Benelearn 2023, Delft, The Netherlands, November 8–10, 2023, Revised Selected Papers* (pp. 40-53). (Communications in Computer and Information Science; Vol. 2187 CCIS). Springer.
https://doi.org/10.1007/978-3-031-74650-5_3

Important note

To cite this publication, please use the final published version (if applicable).
Please check the document version above.

Copyright

In case the licence states "Dutch Copyright Act (Article 25fa)", this publication was made available Green Open Access via the TU Delft Institutional Repository pursuant to Dutch Copyright Act (Article 25fa, the Taverne amendment). This provision does not affect copyright ownership.
Unless copyright is transferred by contract or statute, it remains with the copyright holder.

Sharing and reuse

Other than for strictly personal use, it is not permitted to download, forward or distribute the text or part of it, without the consent of the author(s) and/or copyright holder(s), unless the work is under an open content license such as Creative Commons.

Takedown policy

Please contact us and provide details if you believe this document breaches copyrights.
We will remove access to the work immediately and investigate your claim.

Green Open Access added to TU Delft Institutional Repository

'You share, we take care!' - Taverne project

<https://www.openaccess.nl/en/you-share-we-take-care>

Otherwise as indicated in the copyright section: the publisher is the copyright holder of this work and the author uses the Dutch legislation to make this work public.



Defining and Evaluating the Degrees of Abstraction in Explanations with Kolmogorov Complexity

Jan Lemeire^{1,2}  and Stefan Buijsman³  

¹ Department of Industrial Sciences (INDI), Vrije Universiteit Brussel (VUB),
Pleinlaan 2, 1050 Brussels, Belgium

jan.lemeire@vub.be

² Department of Electronics and Informatics (ETRO), Vrije Universiteit Brussel
(VUB), Pleinlaan 2, 1050 Brussels, Belgium

³ Section of Ethics and Philosophy of Technology, TU Delft, Jaffalaan 5, 2628 BX
Delft, The Netherlands

s.n.r.buijsman@tudelft.nl

Abstract. What variables should be used to get explanations (of AI systems) that are easily interpretable? The challenge to find the right degree of abstraction in explanations, also called the ‘variables problem’, has been actively discussed in the philosophy of science. The challenge is striking the right balance between specificity and generality. Concepts such as proportionality and exhaustivity are investigated and discussed. We propose a new and formal definition based on Kolmogorov complexity and argue that this corresponds to our intuitions about the right level of abstraction. First, we require that variables are uniform, so that they cannot be decomposed into less abstract variables without increasing the Kolmogorov complexity. Next, uniform variables are optimal for an explanation if they can compose the explanation without increasing its Kolmogorov complexity. For this, the concepts K-decomposability and K-composability of sets are defined. Explanations of a certain instance should encompass a maximal set of instances without being K-decomposable. Although Kolmogorov complexity is uncomputable and depends on the choice of programming language, we show that it can be used effectively to evaluate and reason about explanations, such as in the evaluation of XAI methods.

Keywords: Explainability · Explainable AI · Kolmogorov complexity · Abstraction

1 Introduction

How do we best explain a particular outcome of a binary function in terms of the properties of the input? One of the challenges in answering this question is finding the right variables to use in these explanations. Intuitively, we prefer an



Fig. 1. Images of coffee that are correctly identified by a trained Convolutional Neural Network.

explanation that is *not too specific, nor too abstract*. Consider a Convolutional Neural Network (CNN) that is trained to recognize coffee in images and assume the network successfully recognizes coffee in all the images shown in Fig. 1. An explanation for the identification of coffee in image (a) is the dark brown color and the foam. But if image (b) also leads to a positive identification, the property ‘dark brown’ is too specific; ‘brown’ is sufficient and seems to better capture the behaviour of the network. Similarly, image (c) is linked to coffee by the shape of the cup and the steam, but a very specific description of the cup may lead us to believe that the system generalizes less than it does if the more abstract cup in image (d) is also classified as a coffee cup. This can be observed again in image (e) where the ‘flower pattern’ in the foam may not be necessary for identification if image (f) is also classified correctly. Vice versa, ‘rounded shapes’ may be too general a variable for the explanation, even if it matches these examples. The final image (g) illustrates this again, where one may use ‘food’ or ‘medium-sized objects’ to describe what is on the image, but these are likely too abstract compared to options such as ‘croissant’ and ‘coffee cup’.

However, it is a challenge to specify formally what an optimal degree of abstraction is and on which concepts, called *variables* in the philosophical literature, an explanation should be based. This is a general problem for theories of explanation [5, 8], but one that reasserts itself in the field of XAI [6]. Especially in the philosophy of science there has been earlier work on precisely this question, which we survey in Sect. 2. In Sect. 3 we show that current definitions for the upper bound of the level of abstraction are falling short. Our aim in this paper is to build upon this work by giving a formal definition of this degree of abstraction using descriptive complexity. We therefore first introduce Kolmogorov complexity theory in Sect. 4. In Sect. 5 we provide an alternative formal specification for optimal degrees of abstraction of variables in an explanation. We then show in Sect. 6 how this definition applies and resolves the current problems. In the final

section we illustrate how this theoretical discussion applies to XAI methods for explanations.

2 Related Work

There is a wide-ranging literature on explanations of AI systems [1]. Methods showing how important different features were for the output [17, 19] are one option, as are methods extracting rules (e.g. decision trees) to describe the (local) behaviour of the AI system and counterfactuals showing what should be changed to the input to achieve the desired output [14]. And while there is still disagreement about how we should define explanations, both in the computer science literature [11] and in the philosophical literature on explanation [3], the overall goal on all of these definitions is to let recipients of explanations better understand the AI system.

In all of these cases, too, explanations require the use of variables: either the input variables of the system, or a set of (often more abstract) variables that are closer to the variables humans are used to working with. Examples of XAI methods doing the latter are so-called concept-based explainability methods, such as Concept Activation Vectors [10, 24] which attempt to extract (some of) the patterns that a convolutional neural network uses to arrive at the output classification. Alternative methods use crowd workers to attribute concepts to highlighted regions in images [2, 4], thus abstracting from highlighted individual pixels to more abstract concepts. This makes explanations not only more interpretable, as humans are more used to reasoning with concepts such as chairs and tables than we are with sets of pixel values. It also makes explanations more general, as more abstract variables typically cover a wider set of cases.

This is important, as a common standard for the quality of an explanation is how general the explanation is [6]. In other words, “powerful explanations should, just like any predictor, generalize as much as possible” [14, p.36]. However, finding the point where an explanation has generalized *as much as possible* is difficult. The ‘variables problem’ [8, 21] in the philosophy of explanation shows the challenge of identifying optimal degrees of abstraction for explanations. To illustrate with an example commonly used in philosophical literature, there is an intuitive sense that of the following three factually correct explanations the second is the best, being neither too specific nor too general:

- (1) The pigeon pecked because it was presented with a scarlet stimulus
- (2) The pigeon pecked because it was presented with a red stimulus
- (3) The pigeon pecked because it was presented with something stimulating

Specifying why this is so is non-trivial, but by now two approaches can be found. Blanchard [5] suggests that we opt for the most abstract variables that are still specific when compared to even less abstract variables, where abstraction and specificity are defined as follows:

An explanation with explanans variable(s) e_1 is more abstract than an explanation with explanans variable(s) e_2 when the actual value of e_1 is implied by the actual value of e_2 , but not vice versa

An explanation with explanans variable(s) e_1 is more specific than an explanation with explanans variable(s) e_2 when e_2 is a function f of e_1 and other variables e_3, \dots, e_n such that ...neither e_1 nor $G(e_2) = G(f(e_1, e_3 \dots e_n))$ change value if the variables e_3, \dots, e_n are varied. G refers to the explanation provided by the variables.

The proposed definitions apply as follows: using *red* leads to a better explanation than using *scarlet* because it is more abstract (if *scarlet* = 1 then *red* = 1, but not vice versa) without being more specific. On the other hand, using *red* leads to a better explanation than using *something stimulating* because it is more specific (*something stimulating* can be seen as a function $f(\textit{red}, \textit{food}, \textit{tickle}) = \textit{red} \vee \textit{food} \vee \textit{tickle}$ where the value of f remains the same as long as *red* = 1).

Woodward [23] takes a slightly different approach to the same problem, stating that a requirement of proportionality instead motivates the choice of variable. This principle of *proportionality* states that “other things being equal, we should prefer those causal claims/explanations that more fully represent or exhibit those patterns of dependence that hold” [23, p. 247]. It then functions as follows: using *scarlet* suggests the following relation: if *scarlet* is set to 1 then *peck* = 1, if *scarlet* is set to 0 then *peck* = 0 (all other things being equal, so without the presentation of other kinds of stimuli that will make the bird peck) The latter part is false, because we can change the colour of the stimulus to other shades of red (meeting the requirements of setting *scarlet* to 0) while the bird will still peck. Hence, the explanation that if *red* is set to 1 then *peck* = 1, if *red* is set to 0 then *peck* = 0 (all other things being equal) is better because it better represents how the pecking depends on the colour; changing the colour to anything other than red will, all else being equal, entail that the bird stops pecking at it.

Our approach instead links the choice of degree of abstraction to descriptive complexity, based on the intuitive idea that we should choose variables in our explanations that minimize the complexity of these explanations (given equal accuracy). This may be linked to philosophical discussions on effective conversational communication in common social situations. Grice’s four maxims of conversation [12]¹ describe the rationality behind what people expect from effective communication, and stress conveying as much as possible while as relevant and brief as possible. One way to comply with these norms of communication is to minimize descriptive complexity through variable choice.

To make this idea of using the descriptive complexity of an explanation more precise we need an objective complexity measure. Although it has its limitations (discussed in Sect. 4.2) we will use Kolmogorov complexity in this paper. Alternatives are of course available, such as Minimum Description Length (MDL) [13], but Kolmogorov complexity is well-suited due to its generality as MDL is mainly used for finding the optimal model within an a priori chosen model class. That being said, our approach will also work with other approaches to

¹ <https://www.sas.upenn.edu/~haroldfs/drawing/grice.html> for a summary of the 4 principles.

descriptive complexity, as one can simply replace the Kolmogorov measure in our definitions with another way to formalize descriptive complexity. In order to utilize the formal precision that Kolmogorov complexity brings we now need to further formalize the problem of variable choice.

3 Formalization of the Challenge

With the context of XAI as part of the motivation for formalizing the discussion on abstraction, we will consider a binary classifier b which outputs 0 or 1 for each input $x \in X$. x is a multi-dimensional feature vector. The set of all inputs for which b outputs 1 is called the *positive subset*, which we denote with S_b . b is also called the *indicator function* for set S_b . An explanation of an instance $b(x)$ for a particular x is then based on the properties of the input that ensured the outcome. For this paper, we assume that an explanation defines a *sufficient condition* for b , in line with definitions of explanations from [22] and [6] that see them as generalizations describing a set of outputs. Zooming in on a particular instance, this points us to the idea that an explanation G for $x \in S_b$ corresponds to a description of a subset S_G of the positive subset S_b ; it defines a set with elements that are all identified by the classifier. The question we address in this paper is then (1) the optimal choice of subset S_G , how large the subset S_G should be, and (2) which features to use to define the subset.

To illustrate these two points, consider a scenario where a pigeon pecks at both red and yellow stimuli inspired by the work in the philosophical literature. Here, the choice will be formalized based on what we consider to be the relevant S_G to describe this behaviour. We can opt for one explanation per colour, separating the behaviour into two possible scenario's and arriving at the explanation:

- (4) The pigeon pecked because it was presented with a red stimulus OR The pigeon pecked because it was presented with a yellow stimulus

Or we could use a broader subset S_G along with a single composed variable $redow = red \vee yellow$ to generate the following explanation:

- (5) The pigeon pecked because it was presented with a redow stimulus

If we compare these two explanations there seems to be a clear preference for (4), where no new abstract variable is introduced to cover the two cases in one go. However, judging by the criteria of [5] we should in fact prefer (5). The variable is more abstract ($yellow = 1$ implies $redow = 1$ and $red = 1$ implies $redow = 1$) but not more specific (red changes values if the value of $yellow$ is changed and vice versa). We will propose a definition that makes a clear distinction between (4) and (5), and gives a preference for (4).

This definition is based on the idea that the crucial difference between (4) and (5) is the uniformity of the explanatory domain S_G , as red is a uniform set (informally, one can test membership directly, for the formal definition see Definition 5) whereas $redow$ is not (one needs to look whether the stimulus is

either red, or yellow to determine its value). To further illustrate this idea consider a second variation on the same example. Here, there is a range of different colours that a pigeon responds to. What they all have in common is that they are bright colours, such as red, orange and yellow. So, we again have two options for an explanation (assuming here that in this case we decide to opt for a domain S_G that covers all these colour shades in a single explanation). Either we use a disjunction of less abstract variables or we use the more abstract variable *bright colour*:

- (6) The pigeon pecked because it was presented with a bright red or a bright orange or a bright yellow stimulus
- (7) The pigeon pecked because it was presented with a bright coloured stimulus

Here, it seems that the abstract variable in (7) is preferable whereas the abstract variable *redow* in (5) is not. As both are specifiable as disjunctions of incompatible colour concepts this raises the question: what is the difference? In our view, it is that there is a separate, undecomposable way to define *bright colour*. Specifically, colours can be defined using the HSV colour space (https://en.wikipedia.org/wiki/HSL_and_HSV), where the V-component defines the brightness using a single numerical value. There is no such unifying measure for *redow*, which could be why we find this a less illuminating concept to use. This brings us back to the descriptive complexity of the explanation.

The underlying idea behind the formal definitions that we introduce below is that variables should be both general and, at the same time, undecomposable. In other words, the variable should track a single property that can be tested for in a *canonical* manner. Explanation should moreover only contain the relevant information that makes the instance positive, prompting our choice for S_G . An explanation that is too specific contains information about the properties that are not essential for the classifier. On the other hand, an explanation that is too general contains information that is not essential for the classification of the given instance. Descriptive complexity, formalized here using Kolmogorov complexity, can help to make these ideas more precise.

4 The Kolmogorov Complexity of Functions and Sets

To introduce Kolmogorov complexity, we consider first how we can describe an indicator function of a set. Most arbitrary sets (e.g. by randomly picking elements) can only be described by an enumeration of its elements because the elements have no properties in common. In most cases, however, we are interested in sets that mean something, of which the elements have properties in common on which the indicator function can be built. Then, the implementation of the indicator function will become shorter than a literal enumeration. This can be formalized by algorithmic information or ‘Kolmogorov complexity’, a concept put forward as an objective measure of complexity.

4.1 Definition of Kolmogorov Complexity

First we define the Kolmogorov complexity (KC) of a single object:

Definition 1. For a binary sequence $x \in \{0, 1\}^*$, the algorithmic information $K(x)$ (or ‘Kolmogorov complexity’) is defined as the length of the shortest program on a universal Turing machine that generates x and then stops:

$$K(x) = \min_{p: \mathcal{U}(p)=x} l(p) \quad (1)$$

with \mathcal{U} a universal computer, and $l(\cdot)$ the length in bits of a binary sequence.

The shortest program is denoted with p_x^* .

To illustrate this definition, consider the following two sequences of 1000 bits:

- 01111000011001100111 ... 00001111100100011101
- 00010001000100010001 ... 00010001000100010001

The first string is arbitrary without any patterns, while the second repeats “0001”. $K(x)$ is maximal for the random string, namely around 1000 bits. The shortest program literally encodes the string. The second string can be described by program REPEAT 250 TIMES"0001"+ and needs far fewer bits. The program exploits the ‘regularities’ (patterns) of the string to *compress* its description. It is these regularities that make up the meaningful information we are interested in. This same idea can then be applied to indicator functions:

Definition 2. The *Kolmogorov complexity of a binary function* b that takes as argument $x \in X$ and returns 0 or 1, is defined as the length of the shortest program p^* that when executed by a universal Turing machine together with any argument $x \in X$ returns the same output as $b(x)$: $\mathcal{U}(p^*, x) = b(x)$. The shortest program is denoted as p_b^* ,

Definition 3. The *Kolmogorov complexity of a set* $S \subseteq X$ is defined as the Kolmogorov complexity of the indicator function of S .

4.2 Limitations and Practical Use of Kolmogorov Complexity

There are two problems to apply Kolmogorov complexity to practical problems [13]. First, Kolmogorov complexity is not computable. It is proven that there is no algorithm that given a bitstring will output the length of the shortest program and halts [7]. For a lot of cases, however, the shortest program is indisputable, as will be shown in the discussed examples. Still, for more intricate programs it is not trivial, as for example in the case of neural networks trained to detect objects – which quickly use millions of parameters. Instead of trying to identify the absolute shortest implementation (and with that the absolute best concept to use in the explanation), we will therefore use the definitions to *compare and validate implementations* in the same way as explanations are compared in philosophical literature.

Second, Kolmogorov complexity depends on the choice of programming language (or Turing machine) up to a constant. Since one programming language can be translated into another one with a program of length C , the difference of describing x in both languages, can be maximally be C , a constant which does not depend on x . Therefore, theorems often have to incorporate this constant [15]. This can be seen in the additivity rule (which we need later) for the joint Kolmogorov complexity, which has the following formulation:

$$K(x, y) \stackrel{\pm}{=} K(x) + K(y|p_x^*), \quad (2)$$

where $K(y|p_x^*)$ denotes the conditional Kolmogorov complexity of y , given the shortest program p_x^* of x . As usual in algorithmic information theory, $\stackrel{\pm}{=}$ denotes equality up to a constant that is independent of the string x , but does depend on the chosen Turing machine. Since information is symmetric [9], see also [16, Eq. 2.1]:

$$K(x) + K(y|p_x^*) \stackrel{\pm}{=} K(y) + K(x|p_y^*), \quad (3)$$

we can write that:

$$K(x, y) \stackrel{\pm}{=} K(y, x) \stackrel{\pm}{=} K(y) + K(x|p_y^*). \quad (4)$$

5 Abstraction and Undecomposable Concepts

We will now formalize our proposed definitions.

5.1 Formal Definition of Uniformity

Formalizing this idea of having a unifying measure and undecomposable definition available for a concept we can appeal to Kolmogorov complexity to define when a concept meets this requirement. To make this translation we have to interpret concepts as sets, where the set has as members every element to which the concept applies. The question of whether the corresponding concept is appropriately uniform can then be approached in terms of the Kolmogorov complexity of the description of the set:

Definition 4. *A set S is K -decomposable if there exist different and non-empty subsets S_1 and S_2 such that:*

- $S = S_1 \cup S_2$, and
- $K(S) \stackrel{\pm}{=} K(S_1) + K(S_2|p_{S_1}^*)$.

The conditional in the second term of the last equation indicates that the identification of the S_1 by $p_{S_1}^*$ can be reused for describing S_2 . If the description of both sets is nevertheless of equal complexity as S , it signifies that S_1 and S_2 contain no additional information that is not already required to describe S . K -decomposability thus refers to the possibility of decomposing a set into multiple

sets keeping the descriptive complexity invariant. The description of the total set can be decomposed into a separate description of subsets. Note that this definition has the same form as the additivity rule (Eq. (2) in Sect. 4.2), which always holds for Kolmogorov complexities. In the case of K -decomposability, however, we only get additivity if the set decomposition does not bring in new information on the right side of the equation that is not present on the left side: information that is required to identify S_1 and S_2 but which is not required for identifying S . Moreover, by Eq. (3), the definition is symmetric, the roles of S_1 and S_2 can be swapped.

Applied to the example of *something stimulating* we can see that it is in fact K -decomposable. The three variables *red*, *food*, and *tickle* are identified with three separate functions and so $K(S) \stackrel{\pm}{=} K(\text{red}) + K(\text{food}|p_{\text{red}}^*) + K(\text{tickle}|p_{\text{red}}^*, p_{\text{food}}^*)$. The Kolmogorov complexity of the set S is the sum of the KCs of the three subsets. On the other hand, a concept that is appropriately uniform cannot be decomposed. Consider the set of all squares. To decompose this set, we would have to segregate the squares according to a certain criterion. For example, we could apply a threshold on their size to distinguish small from large squares. But then we have to include this criterion in the indicator functions of both subsets, which makes the total description larger than the original one and invalidates the decomposition. Thus, we can plausibly say that ‘square’ is not K -decomposable.

Using the notion of K -decomposable sets we can then define when a variable V is uniform, in accordance with the informal characterization above:

Definition 5. V is a **uniform** variable if the set S_V that corresponds to V is not K -decomposable

In other words, the variables that we are looking for are those that attach to a unified characteristic, such as *red* or *brightness*, guaranteed by the fact that the concepts used are not K -decomposable. Importantly, for any explanation and element x there is a wide range of uniform variables that can be used. A specific x could be both red (one uniform variable), and a rectangle (a second uniform variable) and a large object (a third uniform variable) at the same time. Furthermore, variables at various levels of abstraction can be uniform: *scarlet* is uniform, as is *red* and *colour*. Which uniform variables one chooses then depends on the specific explanation (and to be more precise the domain of that explanation), to which we turn next.

Here, we will say that an explanation $G(x)$ of $b(x)$ regarding element x aims to cover as large a set as possible while still using only patterns (of dependence) relevant to x .

Definition 6. An explanation $G(x)$ is called *optimal* if it has domain $S_G \subseteq S_b$, where S_G is a maximal non- K -decomposable set which contains x

For optimal explanations of x there is then a guarantee that S_G contains as many inputs as possible, while it does not include irrelevant information for x (as in this case it would be possible to K -decompose S_G). Note that there is

always a non-K-decomposable set that contains x , namely the trivial set, which only contains x .

5.2 Formal Definition of Optimal Variable in an Explanation

To arrive at our final definition of optimal variables for an explanation we then also need the notion of *K-composability*. The idea here is that in explanations we often use more than one variable to capture the set of inputs X for whose outputs Y the explanation is supposed to provide additional insight. The interaction of these different variables needs to be accounted for, as they should be complementary. To capture this aspect we therefore define *K-composability* as follows:

Definition 7. *A set S is K-composable if there exist different and non-empty subsets S_1 and S_2 such that:*

- $S_1 \setminus S_2 \neq \emptyset$,
- $S_2 \setminus S_1 \neq \emptyset$,
- $S = S_1 \cap S_2$, and
- $K(S) \stackrel{\pm}{=} K(S_1) + K(S_2|p_{S_1}^*)$.

Using both K-decomposability (based on the union of sets) and K-composability (based on the intersection of sets) we can then define when variables are of an optimal degree of abstraction for an explanation of an input x . Our definition states that these optimal variables together identify this subset S_G of inputs in a complementary fashion, building on the definition of a uniform variable.

Definition 8. *Uniform variables V_1, \dots, V_n , associated with set S_{V_1}, \dots, S_{V_n} , are **optimal variables in explanation** $G(x)$ of input x for $b(x)$ if the set S_G associated with the explanation is such that the sets S_{V_1}, \dots, S_{V_n} corresponding to uniform V_1, \dots, V_n K-compose S_G .*

Figure 2 helps to visualize what this definition states. Our choice of S_G as the maximal subset of S_b that is non-K-decomposable and contains x is illustrated on the left. On the right we see how an explanation of x is then built up using optimal variables. The composability requirement states that if we use different variables in an explanation then they have to overlap, to together characterize S_G . However, they have to do so in complementary fashion (and with a minimal number of variables). So, these sets will be similar to those seen on the right-hand side in the image. For example, if S_G is the set of all big, red rectangles then it is K-composed of the sets S_{V_1} : ‘rectangles’, S_{V_2} : ‘big objects’ and S_{V_3} : ‘red objects’. The requirement that the difference sets are non-empty helps to exclude the possibility to further compose the set of rectangles based on the set of ‘quadrilaterals’ or ‘polygons’. This is not a valid composition of ‘rectangles’ by our definition since ‘rectangles’ minus ‘quadrilaterals’ is the empty set. This, together with the requirement that the Kolmogorov complexity does not increase through K-composition, helps prevent the move to more abstract concepts.

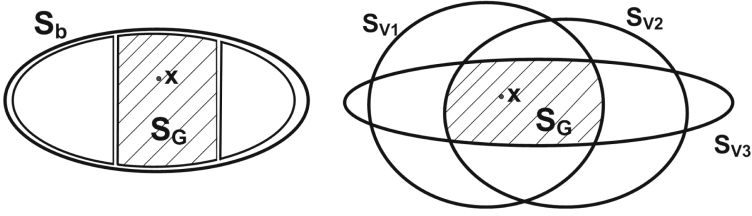


Fig. 2. An explanation G for X explains inputs in subset S_G of the positive subset S_b (left), and does so using the intersection of global variables S_{V1} , S_{V2} and S_{V3} (right).

6 Application of the Definition to the Pigeon Case

If we apply this proposed definition to the examples depicted in Sect. 3, we see that it tracks exactly the judgements we are inclined to make. According to our definition, we should prefer *red* over *scarlet* because the resulting explanation covers a larger subset (namely all red stimuli rather than only the scarlet stimuli), while *red* can be defined in simpler (i.e. shorter) terms than as a disjunction of the different shades of red. The subset S_G is in this case the set of all red stimuli, assuming that x is a specific red stimulus. While we could describe this with variables of different shades of red (which are uniform variables), this increases the Kolmogorov complexity of the set as red is not K -decomposable. Hence, we should prefer *red* over a disjunction of shades of red. We should also prefer it over *scarlet*, as *scarlet* does not K -compose S_G on its own. Finally, more abstract variables such as *colour* are ruled out as composition of more abstract variables is more complex than simply using *red* (violating the last condition of K -composition). However, had the bird pecked only at scarlet stimuli, our definition would state that *scarlet* is an optimal variable to use. In that case, S_b would not have contained other shades of red and so likewise $S_G \subseteq S_b$ would have been restricted to the specific set of scarlet stimuli, which are then captured by the uniform variable *scarlet*.

Furthermore, we should prefer *red* over *something stimulating* in the situation where the bird pecks at a wider range of stimuli. Despite the broader reach of *something stimulating* it is K -decomposable in terms of *red*, *food* and *tickle*. As a result, we first fix S_G as one of the maximal non- K -decomposable subsets, in the case of a red stimulus this will be *red*. This means that *something stimulating*, the variable that covers all S_b , is ruled as being too abstract. Instead, we should look at the minimal number of non- K -decomposable sets that together K -compose the smaller set S_G , which is simply *red* again. Should we want a more general explanation of the behaviour of the pigeon then we can simply go for the disjunction of the explanations corresponding to our K -decomposed subsets: $red \vee food \vee tickle$. For the follow-up examples we get again the desired results: *red* is preferable to *redow* because *redow* is K -decomposable in terms of *red* and *yellow* (which affects the choice of S_G), whereas *bright coloured* is preferable because it is not K -decomposable, again assuming that the set of positive instances that we

aim to explain is in the first case that of red and yellow stimuli and in the second case that of bright coloured stimuli.

To consider this in the pigeon example just discussed, we can imagine a setting in which the brain of the pigeon is studied and neural signals are measured. Based on these measurements it is observed that when a particular part of the brain gets stimulated it cause the pecking. In such a context, with the knowledge of what’s going on in the brain, ‘something stimulating’ can provide a good (uniform) explanation for the pecking. For each of the different stimuli (a bright color, food or tickle), a similar process in the brain can be observed. This shows that an explanation might depend on the features available to the binary classifier b .

7 Explanations in AI

A popular option in XAI for explaining black-box algorithms is to fit decision trees to them, which aim to approximate the input-output relation of b [14, 20]. They offer a human-understandable description, thanks to the explicit variables and clear decision paths (for trees that are not too large). Can they represent/describe the level of abstraction we defined in this paper? Decision trees are based on clauses that form conditions on the input variables by conjunctions, negations, and disjunctions. Also rule-based systems are based on such clauses.

Figure 3 shows the decision tree for the Pigeon example (4) in which the color might be red or yellow. The decision tree provides an explanation for all positive instances, where each node is an optimal variable for an explanation and each leaf represents the right level of abstraction for an explanation of an individual outcome (an S_G -set).

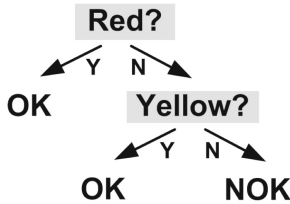


Fig. 3. Decision tree for the Pigeon2 example (Red or Yellow). (Color figure online)

In this example, the decision tree is the shortest description of the partitioning. But this is not always true: for more regular structures, a shorter description is possible. Consider the partitioning of a chessboard into 64 squares. Describing all black squares separately results in a large tree. An algorithm can do this in a more succinct way by exploiting the regularities of a chessboard. The KC of the black squares is smaller than a literal enumeration of all squares.

Likewise, the set of images containing a pattern, such as a rectangle which is recognized by a NN, cannot be explained succinctly by a decision tree. The NN employs multiple layers of various operations applied to the image pixels to achieve the recognition. This cannot be described by simple clauses. Patterns are the regularities that reduce the KC, while constraints on parameters do not reduce the KC. Consider the set S_b of the rectangles of a certain size and color. An explanation for S_b is K-composed of 3 optimal variables: the rectangular shape, the size and the color. The variables can be extracted without increasing the KC. The first one describes the pattern. The second and third are constraints. Such constraints can be formed by conditions on the input variables, but also conditions on the parameters of the patterns. Conditions can be described succinctly by a decision tree or rules. Patterns, however, cannot. To overcome this challenge, [18] propose to use decision trees containing *prototypes* that are representative for a set of similar instances. By checking against the prototype in the node it is possible to classify a case using more abstract concepts.

8 Conclusion

How abstract should variables in explanations be? We have proposed an account based on Kolmogorov complexity which, although not computable, gives us a formal definition of optimal degrees of abstraction. As shown in Sect. 5, our formal definition handles the examples in the philosophical literature well. We have done so by first defining the notion of a uniform, i.e. undecomposable, variable. Which of these uniform variables is optimal for a given explanation is then based on what variables can be combined to characterize the patterns of dependence captured by the explanation with minimal Kolmogorov complexity. Ultimately, therefore, we approach the problem of abstraction by arguing that optimal degrees of abstraction are those which lead to the least complex description of the patterns and constraints in the explanation. As abstraction is precisely meant to simplify description, we consider it a natural link to say that optimal degrees of abstraction are those which optimally reduce the complexity of descriptions.

References

1. Arrieta, A.B., et al.: Explainable artificial intelligence (XAI): concepts, taxonomies, opportunities and challenges toward responsible AI. *Inf. Fusion* **58**, 82–115 (2020)
2. Balayn, A., Soilis, P., Lofi, C., Yang, J., Bozzon, A.: What do you mean? Interpreting image classification with crowdsourced concept extraction and analysis. In: *Proceedings of the Web Conference 2021*, pp. 1937–1948 (2021)
3. Beisbart, C., R az, T.: Philosophy of science at sea: clarifying the interpretability of machine learning. *Philos. Compass* **17**(6) (2022)
4. Biswas, S., Corti, L., Buijsman, S., Yang, J.: CHIME: causal human-in-the-loop model explanations. In: *Proceedings of the AAAI Conference on Human Computation and Crowdsourcing*, vol. 10, pp. 27–39 (2022)

5. Blanchard, T.: Explanatory abstraction and the goldilocks problem: interventionism gets things just right. *Br. J. Philos. Sci.* (2020)
6. Buijsman, S.: Defining explanation and explanatory depth in XAI. *Mind. Mach.* **32**(3), 563–584 (2022)
7. Chaitin, G.J., Arslanov2, A., Calude3, C.: Program-size complexity computes the halting problem. *Bull. EATCS* **57** (1995)
8. Franklin-Hall, L.R.: High-level explanation and the interventionist’s ‘variables problem’. *Br. J. Philos. Sci.* (2016)
9. Gács, P.: On the symmetry of algorithmic information. In: *Soviet Mathematics - Doklady*, vol. 15, pp. 1477–1480 (1974)
10. Ghorbani, A., Wexler, J., Zou, J.Y., Kim, B.: Towards automatic concept-based explanations. *Adv. Neural Inf. Process. Syst.* **32** (2019)
11. Gilpin, L.H., Paley, A.R., Alam, M.A., Spurlock, S., Hammond, K.J.: “Explanation” is not a technical term: the problem of ambiguity in XAI. *arXiv preprint arXiv:2207.00007* (2022)
12. Grice, H.P.: Logic and conversation. In: Cole, P., Morgan, J.L. (eds.) *Syntax and Semantics: Vol. 3: Speech Acts*, pp. 41–58. Academic Press (1975)
13. Grünwald, P.: *The Minimum Description Length Principle*. MIT Press, Cambridge (2007)
14. Guidotti, R., Monreale, A., Ruggieri, S., Turini, F., Giannotti, F., Pedreschi, D.: A survey of methods for explaining black box models. *ACM Comput. Surv.* **51**(5) (2018)
15. Li, M., Vitányi, P.: *An Introduction to Kolmogorov Complexity and its Applications*. Springer, Heidelberg (1997)
16. Li, M., Chen, X., Li, X., Ma, B., Vitányi, P.M.: The similarity metric. *IEEE Trans. Inf. Theory* **50**(12), 3250–3264 (2004)
17. Lundberg, S.M., Lee, S.I.: A unified approach to interpreting model predictions. *Advances Neural Inf. Process. Syst.* **30** (2017)
18. Nauta, M., Van Bree, R., Seifert, C.: Neural prototype trees for interpretable fine-grained image recognition. In: *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pp. 14933–14943 (2021)
19. Ribeiro, M.T., Singh, S., Guestrin, C.: “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*, pp. 1135–1144 (2016)
20. Sagi, O., Rokach, L.: Explainable decision forest: transforming a decision forest into an interpretable tree. *Inf. Fusion* **61**, 124–138 (2020)
21. Weatherston, B.: Explanation, idealisation and the goldilocks problem. *Res.* **84**(2), 461–473 (2012)
22. Woodward, J.: *Making Things Happen: A Theory of Causal Explanation*. Oxford University Press, Oxford (2005)
23. Woodward, J.: Explanatory autonomy: the role of proportionality, stability, and conditional irrelevance. *Synthese* **198**, 237–265 (2021)
24. Yeh, C.K., Kim, B., Arik, S., Li, C.L., Pfister, T., Ravikumar, P.: On completeness-aware concept-based explanations in deep neural networks. *Adv. Neural. Inf. Process. Syst.* **33**, 20554–20565 (2020)