A review of possible effects of cognitive biases on interpretation of rule-based machine learning models

Tomáš Kliegr^{a,*}, Štěpán Bahník^b, Johannes Fürnkranz^c

Abstract

While the interpretability of machine learning models is often equated with their mere syntactic comprehensibility, we think that interpretability goes beyond that, and that human interpretability should also be investigated from the point of view of cognitive science. In particular, the goal of this paper is to discuss to what extent cognitive biases may affect human understanding of interpretable machine learning models, in particular of logical rules discovered from data. Twenty cognitive biases are covered, as are possible debiasing techniques that can be adopted by designers of machine learning algorithms and software. Our review transfers results obtained in cognitive psychology to the domain of machine learning, aiming to bridge the current gap between these two areas. It needs to be followed by empirical studies specifically focused on the machine learning domain.

Keywords: cognitive bias, cognitive illusion, machine learning, interpretability, rule induction

1. Introduction

This paper aims to investigate the possible effects of cognitive biases on human understanding of machine learning models, in particular inductively learned rules. We use the term "cognitive bias" as a representative for various cognitive phenomena that materialize themselves in the form of occasionally irrational reasoning patterns, which are thought to allow humans to make fast judgments and decisions. Their cumulative effect on human reasoning should not be underestimated as "cognitive biases seem reliable, systematic, and difficult to eliminate" [79]. The effect of some cognitive biases is more pronounced when people do

^{*}Corresponding author

not have well-articulated preferences [154], which is often the case in explorative machine learning.

Previous works have analysed the impact of cognitive biases on multiple types of human behavior and decision making. A specific example is the seminal book "Social cognition" by Kunda [85], which is concerned with the impact of cognitive biases on social interaction. Another, more recent work by Serfas [133] focused on the context of capital investment. Closer to the domain of machine learning, in their article entitled "Psychology of Prediction", Kahneman and Tversky [80] warned that cognitive biases can lead to violations of the Bayes theorem when people make fact-based predictions under uncertainty. These results directly relate to inductively learned rules, since these are associated with measures such as confidence and support expressing the (un)certainty of the prediction they make. Despite some early works [99, 100] showing the importance of study of cognitive phenomena for rule induction and machine learning in general, there has been a paucity of follow-up research. In previous work [51], we have evaluated a selection of cognitive biases in the very specific context of whether minimizing the complexity or length of a rule will also lead to increased interpretability, which is often taken for granted in machine learning research.

In this paper, we attempt to systematically relate cognitive biases to the interpretation of machine learning results. To that end, we review twenty cognitive biases that can distort interpretation of inductively learned rules. The review is intended to help to answer questions such as: How do cognitive biases affect human understanding of symbolic machine learning models? What could help as a "debiasing antidote"?

This paper is organized as follows. Section 2 provides a brief review of related work published at the intersection of rule learning and psychology. Section 3 motivates our study on the example of the insensitivity to sample size effect. Section 4 describes the criteria that we applied to select a subset of cognitive biases into our review, which eventually resulted in twenty biases. These biases and their disparate effects and causes are covered in detail in Section 5. Section 6 provides a concise set of recommendations aimed at developers of rule learning algorithms and user interfaces. In Section 7 we state the limitations of our review and outline directions for future work. The conclusions summarize the contributions of the paper.

2. Background and Related Work

We selected individual rules as learnt by many machine learning algorithms as the object of our study. Focusing on simple artefacts—individual rules—as opposed to entire models such as rule sets or rule lists allows a deeper, more focused analysis since a rule is a small self-contained item of knowledge. Making a small change in one rule, such as adding a new condition, allows to test the effect of an individual factor. In this section, we first motivate our work by putting it into the context of prior research on related topics. Then, we proceed by a brief introduction to inductive rule learning (Section 2.2) and a brief recapitulation of previous work in cognitive science on the subject of

decision rules (Section 2.3). Finally, we introduce cognitive biases (Section 2.4) and and rule plausibility (Section 2.5, which is a measure of rule comprehension.

2.1. Motivation

In the following three paragraphs, we discuss our motivation for this review, and summarize why we think this work is relevant to the larger artificial intelligence community.

Explaining "black box" models with rules. While neural networks and ensembles of decision trees are increasingly becoming the prevalent type of representation used in machine learning, it might be at first surprising that our review focuses almost exclusively on decision rules. The reason is that rules are widely used as a means for communicating explanations of a variety of machine learning approaches, since many types of models can be converted to rules, or rules can be extracted from them [5]. Decision trees can be represented as a rule model, when one rule is generated for each path from the root of the tree to a leaf. Guidotti et al. [66] provides a review of research on using rules for explaining some "black-box models", including neural networks, support vector machines and tree ensembles.

Embedding cognitive biases to learning algorithms. The applications of cognitive biases go beyond explaining existing machine learning models. For example, Taniguchi et al. [145] demonstrate how a cognitive bias can be embedded into a machine learning algorithm, achieving superior performance on small datasets compared to commonly used machine learning algorithms with "generic" inductive bias.

Paucity of research on cognitive biases in artificial intelligence. Several recent position and review papers on explainability in Artificial Intelligence (xAI) recognize that cognitive biases play an important role in explainability research [101, 117]. To our knowledge, the only systematic treatment of psychological phenomena applicable to machine learning is provided by the review of Miller [101], which focuses on reasons and thought processes that people apply during explanation selection, such as causality, abnormality and the use of counterfactuals. This authoritative review observes that there are currently no studies that look at cognitive biases in the context of selecting explanations. Because of the paucity of applicable research focusing on machine learning, the review of Miller [101] — same as the present paper — takes the first step of applying influential psychological studies to explanation in the xAI context without accompanying experimental validation specific to machine learning. While Miller [101] summarizes main reasoning processes that drive generation and understanding of explanations, our review focuses specifically on cognitive biases as psychological phenomena that can distort interpretation of machine learning models, if not properly accounted for.

```
IF A AND B THEN C
    confidence=c and support=s

IF veil is white AND odor is foul THEN mushroom is poisonous
    confidence = 90%, support = 5%
```

Figure 1: Inductively learned rule

2.2. Decision Rules in Machine Learning

While neural networks and ensembles of decision trees are increasingly becoming the prevalent type of representation used in machine learning, our review focuses almost solely on decision rules because many types of models can be converted to rules, or rules can be extracted from them, and rules are therefore widely used as a means for communicating explanations of a variety of machine learning approaches.

An example of an inductively learned decision rule, which is a subject of the presented review, is shown in Figure 1. Following the terminology of Fürnkranz et al. [50], A, B, C represent literals, i.e., Boolean expressions which are composed of attribute name (e.g., veil) and its value (e.g., white). The conjunction of literals on the left side of the rule is called antecedent or $rule\ body$, the single literal predicted by the rule is called consequent or $rule\ head$. Literals in the body are sometimes referred to as conditions throughout the text, and the consequent as the target. While this rule definition is restricted to conjunctive rules, other definitions, e.g., the formal definition given by Slowinski et al. [138], also allow for negation and disjunction as connectives.

Rules on the output of rule learning algorithms are most commonly characterized by two parameters, confidence and support. The confidence of a rule—sometimes also referred to as precision—is defined as a/(a+b), where a is the number objects that match both the conditions of the rule as well as the consequent, and b is the number of objects that match the antecedent but not the consequent. The support of a rule is either defined as a/N, where N is the number of all objects (relative support), or simply as a (absolute support). A related measure is coverage, which is the total number of objects that satisfy the body of the rule (a+b). In the special case of learning rules for the purpose of building a classifier, the consequent of a rule consists only of a single literal, the so-called class. In this case, a is also known as the number of true positives, and b as the false positives.

Some rule learning frameworks, in particular association rule learning [1, 171], require the user to set thresholds for minimum confidence and support. Only rules with confidence and support values meeting or exceeding these thresholds are included on the output of rule learning and presented to the user.

2.3. Decision Rules in Cognitive Science

Rules are used in commonly embraced models of human reasoning in cognitive science [139, 112, 120]. They also closely relate to Bayesian inference, which also frequently occurs in models of human reasoning. Consider the first rule of Figure 1. This rule can be interpreted as a hypothesis corresponding to the logical implication $A \wedge B \Rightarrow C$. We can express the plausibility of such a hypothesis in terms of Bayesian inference as the conditional probability $\Pr(C|A,B)$. This corresponds to the confidence of the rule, as used in machine learning and as defined above, and to *strength of evidence*, a term used by cognitive scientists [151].

Given that $\Pr(C|A,B)$ is a probability estimate computed on a sample, another relevant piece of information for determining the plausibility of the hypothesis is the robustness of this estimate. This corresponds to the number of instances for which the rule has been observed to be true. The size of the sample (typically expressed as ratio) is known as rule support in machine learning and as weight of the evidence in cognitive science [151].

Psychological research on hypothesis testing in rule discovery tasks has been performed in cognitive science at least since the 1960s. The seminal article by Wason [161] introduced what is widely referred to as Wason's 2-4-6 task. Participants are given the sequence of numbers 2, 4 and 6 and asked to find out the rule that generated this sequence. In search for the hypothesized rule they provide the experimenter other sequences of numbers, such as 3-5-7, and the experimenter answers whether the provided sequence conforms to the rule, or not. While the target rule is simple "ascending sequence", people find it difficult to discover this specific rule, presumably because they use the positive test strategy, a strategy of testing a hypothesis by examining evidence confirming the hypothesis at hand rather then searching for disconfirming evidence [82].

2.4. Cognitive Bias

According to the Encyclopedia of Human Behavior [164], the term cognitive bias was introduced in the 1970s by Amos Tversky and Daniel Kahneman [151], and is defined as a

"systematic error in judgment and decision-making common to all human beings which can be due to cognitive limitations, motivational factors, and/or adaptations to natural environments."

The narrow initial definition of cognitive bias as a shortcoming of human judgment was criticized by German psychologist Gerd Gigerenzer, who started in the late 1990s the "Fast and frugal heuristic" program to emphasize ecological rationality (validity) of judgmental heuristics [61]. According to this research

¹ Interestingly, balancing the likelihood of the judgment and the weight of the evidence in the assessed likelihood was already studied by Keynes [81] (according to Camerer and Weber [22]).

program, cognitive biases often result from an application of a heuristic in an environment for which it is not suited rather than from problems with heuristics themselves, which work well in usual contexts.

In the present view, we define cognitive biases and associated phenomena broadly. We include cognitive biases related to thinking, judgment, and memory. We also include descriptions of thinking strategies and judgmental heuristics that may result in cognitive biases, even if they are not necessarily biases themselves.

Debiasing. An important aspect related to the study of cognitive biases is the validation of strategies for mitigating their effects in cases when they lead to incorrect judgment. A number of such debiasing techniques has been developed, with researchers focusing intensely on the clinical and judicial domains (cf. e.g. [89, 27, 95]), apparently due to costs associated with erroneous judgment in these domains. Nevertheless, general debiasing techniques can often be derived from such studies.

The choice of an appropriate debiasing technique typically depends on the type of error induced by the bias, since this implies an appropriate debiasing strategy [6]. Larrick [88] recognizes the following three categories: psychophysically-based error, association-based error, and strategy-based error. The first two are attributable to unconscious, automatic processes, sometimes referred to as "System 1". The last one is attributed to reasoning processes (System 2) [36]. For biases attributable to System 1, the most generic debiasing strategy is to shift processing to the conscious System 2 [92], [134, p. 491].

Another perspective on debiasing is provided by Croskerry et al. [27], who organize debiasing techniques by their way of functioning, rather than the bias they address, into the following three categories: educational strategies, work-place strategies and forcing functions. While Croskerry et al. [27] focused on clinicians, our review of debiasing aims to be used as a starting point for analogous guidelines for an audience of machine learning practitioners. For example, the general workplace strategies applicable in the machine learning context include group decision making, personal accountability, and planning time-out sessions to help slowing down.

Function and validity of cognitive biases. In the introduction, we briefly characterized cognitive biases as seemingly irrational reasoning patterns that are thought to allow humans to make fast and risk-averse decisions. In fact, the function of cognitive biases is subject of scientific debate. According to the review of functional views by Pohl [122], there are three fundamental positions among researchers. The first group considers them as dysfunctional errors of the system, the second group as faulty by-products of otherwise functional processes, and the third group as adaptive and thus functional responses. According to Pohl [122], most researchers are in the second group, where cognitive biases are considered to be "built-in errors of the human information-processing systems".

In this work, we consider cognitive biases as strategies that evolved to improve the fitness and chances of survival of the individual in particular situations

or are consequences of such strategies. This defense of biases is succinctly expressed by Haselton and Nettle [68]: "Both the content and direction of biases can be predicted theoretically and explained by optimality when viewed through the long lens of evolutionary theory. Thus, the human mind shows good design, although it is design for fitness maximization, not truth preservation."

According to the same paper, empirical evidence shows that cognitive biases are triggered or strengthened by environmental cues and context [68]. Given that the interpretation of machine learning results is a task unlike the simple automatic cognitive processes to which a human mind is adapted, cognitive biases are likely to have an influence upon it.

2.5. Measures of Interpretability, Perceived and Objective Plausibility

We claim that cognitive biases can affect the interpretation of rule-based models. However, how does one measure interpretability? According to our literature review, there is no generally accepted measure of interpretability of machine learning models. Model size, which was used in several studies, has recently been criticized [46, 142, 51] primarily on the grounds that the model's syntactic size does not capture any aspect of the model's semantics. A particular problem related to semantics is the compliance to pre-existing expert knowledge, such as domain-specific monotonicity constraints.

In our work, we embrace the concept of plausibility to measure interpretability [52]. The word 'plausible' is defined according to the Oxford Dictionary of US English as "seeming reasonable or probable" and according to the Cambridge dictionary of UK English as "seeming likely to be true, or able to be believed". We can link the inductively learned rule to the concept of "hypothesis" used in cognitive science. There is a body of work in cognitive science on analyzing the perceived plausibility of hypotheses [57, 58, 4].

In a recent review of interpretability definitions by Bibal and Frénay [17], the term plausibility is not explicitly covered, but a closely related concept of justifiability is stated to depend on interpretability. Martens et al. [94] define justifiability as "intuitively correct and in accordance with domain knowledge". By adopting plausibility, we address the concern expressed in Freitas [46] regarding the need to reflect domain semantics when interpretability is measured.

We are aware of the fact that if a decision maker finds a rule plausible, it does not necessarily mean that the rule is correctly understood, it can be quite the contrary in many cases. Nevertheless, we believe that the *alignment of the perceived plausibility with objective, data-driven, plausibility of a hypothesis* should be at the heart of an effort that strives for interpretable machine learning.

3. Motivational Example

It is well known in machine learning that chance rules with a deceptively high confidence can appear in the output of rule learning algorithms [8]. For this reason, the rule learning process typically outputs both confidence and support for the analyst to make an informed choice about merits of each rule.

Example.

- IF a film is released in 2006 AND the language of the film is English THEN Rating is good, confidence = 80%, support = 10%.
- IF a film is released in 2006 AND the director was John Smith THEN Rating is good, confidence = 90%, support = 1%.

In the example above, both rules are associated with values of confidence and support to inform about the strength and weight of evidence for both rules. While the first rule is less strong (80% vs 90% correct), its weight of the evidence is ten times higher than of the second rule.

According to the insensitivity to sample size effect [151] there is a systematic bias in human thinking that makes humans overweigh the strength of evidence (confidence) and underweigh the weight of evidence (support). The bias has been also shown in psychologists knowledgable in statistics [149] and thus is likely to be applicable to the widening number of professions that use rule learning to obtain insights from data.

The analysis of relevant literature from cognitive science not only reveals applicable biases, but also sometimes provides methods for limiting their effect (debiasing). The standard way used in rule learning software for displaying rule confidence and support metrics is to use percentages, as in our example. Extensive research in psychology has shown that if frequencies are used instead, then the number of errors in judgment drops [60, 62]. Reflecting these suggestions, the first rule in our example could be presented as follows:

Example.

• IF a film is released in 2006 AND the language of the film is English THEN Rating is good.

In our data, there are 100 movies which match the conditions of this rule. Out of these, 80 are predicted correctly as having good rating.

Rules can be presented in different ways (as shown), and depending on the way the information is presented, humans may perceive their plausibility differently. In this particular example, confidence is no longer conveyed as a percentage "80%" but using expression "80 out of 100". Support is presented as an absolute number (100) rather than as a percentage (10%).

A correct understanding of machine learning models can be difficult even for experts. In this section, we tried to motivate why addressing cognitive biases can play an important role in making the results of inductive rule learning more understandable. In the remainder of this paper, the bias applied to our example will be revisited in greater depth, along with 19 other biases.

4. Selection Criteria

A number of cognitive biases have been discovered, experimentally studied, and extensively described in the literature. As Pohl [122] states in a recent authoritative book on cognitive illusions: "There is a plethora of phenomena showing that we deviate in our thinking, judgment and memory from some objective and arguably correct standard."

We first selected a subset of biases which would be reviewed. To select applicable biases, we considered those that have some relation to the following properties of inductively learned rules: 1. rule length (the number of literals in an antecedent), 2. rule interest measures (especially support and confidence), 3. position (ordering) of conditions in a rule and ordering of rules in the rule list, 4. specificity and predictive power of conditions (correlation with a target variable), 5. use of additional logical connectives (conjunction, disjunction, negation), 6. treatment of missing information (inclusion of conditions referring to missing values), and 7. conflict between rules in the rule list.

Through selection of appropriate learning heuristics, the rule learning algorithm can influence these properties. For example, most heuristics implement some form of a trade-off between the coverage or support of a rule, and its implication strength or confidence [49, 50].

5. Review of Cognitive Biases

In this section, we cover a selection of twenty cognitive biases. For all of them, we include a short description including an example of a study demonstrating the bias and its proposed explanation. We pay particular attention to their potential effect on the interpretability of rule learning results, which has not been covered in previous works. For all cognitive biases we also suggest a debiasing technique that could be effective in aligning the perceived plausibility of the rule with its objective plausibility.

In a recent scientometric survey of research on cognitive biases in information systems [44], no papers are mentioned that aim at machine learning. For general information systems research, the authors claim that "most articles' research goal [is] to provide an explanation of the cognitive bias phenomenon rather than to develop ways and strategies for its avoidance or targeted use". In contrast, our review aims at advancement of the field beyond explanation of applicable phenomena, by discussing specific debiasing techniques.

An overview of the main features of the reviewed cognitive biases is presented in Table 1. Note that the debiasing techniques that we describe have only limited grounding in applied psychological research and require further validation, since as Lilienfeld et al. [92] observe, there is a general paucity of research on debiasing in psychological literature, and the existing techniques suffer from a lack of theoretical coherence and a mixed research evidence concerning their efficacy.

| phenomenon | implications for rule-learning | debiasing technique |
|------------------------------------|--|--|
| Representativeness Heuristic | Overestimate the probability of condition representative of consequent | Use natural frequencies instead of ratios or probabilities |
| Averaging Heuristic | Probability of antecedent as the average of probabilities of conditions | Reminder of probability theory |
| Disjunction Fallacy | Prefer more specific conditions over less specific | Inform on taxonomical relation between conditions; explain benefits of higher support |
| Base-rate Neglect | Emphasis on confidence, neglect for support | Express confidence and support in natural frequencies |
| Insensitivity to Sample Size | Analyst does not realize the increased reliability of confidence | Present support as absolute number rather than percentage; |
| | estimate with increasing value of support | use support to compute confidence (reliability) intervals for the value of confidence |
| Availability Heuristic | Ease of recollection of instances matching the rule | Explain to analyst why instances matching the particular rule are (not) easily recalled |
| Reiteration Effect | Presentation of redundant rules or conditions increases plausibility | rule pruning; clustering; explaining overlap |
| Confirmation Bias | Rules confirming analyst's prior hypothesis are "cherry picked" | Explicit guidance to consider evidence for and against hypothesis; education about the bias; interfaces making users slow down |
| Mere Exposure Effect | Repeated exposure (even subconscious) results in increased preference | Changes to user interfaces that limit subliminal presentation of rules |
| Overconfidence and underconfidence | Rules with small support and high confidence are "overrated" | Present less information when not relevant via pruning, feature selection, limiting rule length; actively present conflicting rules/knowledge. |
| Recognition Heuristic | Recognition of attribute or its value increases preference | More time; knowledge of attribute/value |
| Information Bias | belief that more information (rules, conditions) will improve decision making even if it is irrelevant | Communicate attribute importance |
| Ambiguity Aversion | Prefer rules without unknown conditions | Increase user motivation; instruct users to provide textual justifications |
| Confusion of the Inverse | Confusing the difference between the confidence of the rule Pr(consequent antecedent) with Pr(antecedent consequent) | Training in probability theory; unambiguous wording |
| Misunderstanding of "and" | "and" is understood as disjunction | Unambiguous wording; visual representation |
| Context and Tradeoff Contrast | Preference for a rule is influenced by other rules | Removal of rules, especially of those that are strong, yet irrelevant |
| Negativity Bias | Words with negative valence in the rule make it appear more important | Review words with negative valence in data, and possibly replace with neutral alternatives |
| Primacy Effect | Information presented first has the highest impact | Education on the bias; resorting; rule annotation |
| Weak Evidence Effect | Condition only weakly perceived as predictive of target decreases plausibility | Numerical expression of strength of evidence; omission of weak predictors (conditions) |
| Unit Bias | Conditions are perceived to have same importance | Inform on discriminatory power of conditions |

Table 1: Summary of analysis of cognitive biases.

5.1. Conjunction Fallacy and Representativeness Heuristic

The conjunction fallacy refers to a judgment that is inconsistent with the conjunction rule – the probability of a conjunction, Pr(A, B), cannot exceed the probability of its constituents, Pr(A) and Pr(B). It is often illustrated with the "Linda" problem in the literature [153]. In the Linda problem, depicted in Figure 2, subjects are asked to compare conditional probabilities Pr(F, B|L) and Pr(B|L), where B refers to "bank teller", F to "active in feminist movement" and L to the description of Linda [10].

Linda is 31 years old, single, outspoken, and very bright. She majored in philosophy. As a student, she was deeply concerned with issues of discrimination and social justice, and also participated in anti-nuclear demonstrations.

Which is more probable?

- (a) Linda is a bank teller.
- (b) Linda is a bank teller and is active in the feminist movement.

Figure 2: Linda problem

Multiple studies have shown that people tend to consistently select the second hypothesis as more probable, which is in conflict with the conjunction rule. In other words, it always holds for the Linda problem that

$$\Pr(F, B|L) \le \Pr(B|L).$$

Preference for the alternative $F \wedge B$ (option (b) in Figure 2) is thus always a logical fallacy. For example, Tversky and Kahneman [153] report that 85% of their subjects indicated (b) as the more probable option for the Linda problem. The conjunction fallacy has been shown across multiple settings (hypothetical scenarios, real-life domains), as well as for various kinds of subjects (university students, children, experts, as well as statistically sophisticated individuals) [146].

The conjunction fallacy is often explained by use of the representativeness heuristic [79]. The representativeness heuristic refers to the tendency to make judgments based on similarity, based on the rule "like goes with like", which is typically used to determine whether an object belongs to a specific category. When people use the representativeness heuristic, "probabilities are evaluated by the degree to which A is representative of B, that is by the degree to which A resembles B" [151]. This heuristic provides people with means for assessing a probability of an uncertain event. It is used to answer questions such as "What is the probability that object A belongs to class B? What is the probability that event A originates from process B?" [151].

The representativeness heuristic is not the only explanation for the results of the conjunction fallacy experiments. Hertwig et al. [70] hypothesized that the fallacy is caused by "a misunderstanding about conjunction", in other words by a different interpretation of "probability" and "and" by the subjects than assumed by the experimenters. The validity of this alternative hypothesis has been subject to criticism [146], nevertheless some empirical evidence suggests that the problem of correct understanding of "and" is of particular importance to rule learning [52].

Recent research has provided several explanations for conjunctive and disjunctive (cf. Section 5.4) fallacies, such as configural weighting and adding theory [111], applying principles of quantum cognition [21] and inductive confirmation theory [147]. In the following, we will focus on the CWA theory. CWA essentially assumes that the causes of conjuctive and disjunctive fallacies relate to the fact that decision makers perform weighted average instead of multiplication of the component probabilities. For conjunctions, weights are set so that more weight is assigned to the lower component probability. For disjunctive probabilities, more weight is assigned to the likely component. This assumption was verified in at least one study [41]. For more discussion of the related averaging heuristic, cf. Section 5.3.

Implications for rule learning. Rules are not composed only of conditions, but also of an outcome (the value of a target variable in the consequent). A higher number of conditions generally allows the rule to filter a purer set of objects with respect to the value of the target variable than a smaller number of conditions. Application of representativeness heuristic can affect the human perception of rule plausibility in that rules that are more "representative" of the user's mental image of the concept may be preferred even in cases when their objective discriminatory power may be lower.

Debiasing techniques. A number of factors that decrease the proportion of subjects exhibiting the conjunction fallacy have been identified: Charness et al. [24] found that the number of participants committing the fallacy is reduced under a monetary incentive. Such an addition was reported to drop the fallacy rate in their study from 58% to 33%. The observed rate under a monetary incentive suggests smaller importance of this problem for important real-life decisions. Zizzo et al. [172] found that unless the decision problem is simplified, neither monetary incentives nor feedback can ameliorate the fallacy rate. A reduced task complexity is a precondition for monetary incentives and feedback to be effective.

Stolarz-Fantino et al. [143] observed that the rate of fallacies is reduced but still strongly present when the subjects receive training in logic. Gigerenzer and Goldstein [60] as well as Gigerenzer and Hoffrage [62] showed that the rate of fallacies can be reduced or even eliminated by presenting the problems in terms of frequencies rather than probabilities.

Nilsson et al. [111] present a computer simulation showing that when the component probabilities are not precisely known, averaging often provides equally

good alternative to the normative computation of probabilities (cf. also Juslin et al. [76]). This computational model could be possibly adopted to detect high risk of fallacy, corresponding to the case when the deviation between the perceived probability and the normative probability is high.

5.2. Misunderstanding of "and"

The misunderstanding of "and" refers to a phenomenon affecting the syntactic comprehensibility of the logical connective "and". As discussed by Hertwig et al. [70], "and" in natural language can express several relationships, including temporal order, causal relationship, and most importantly, can also indicate a collection of sets² as well as their intersection. People can therefore interpret "and" in a different meaning than intended.

For example, according to the two experiments reported by Hertwig et al. [70], the conjunction "bank teller and active in the feminist movement" used in the Linda problem (cf. Section 5.1) was found by about half of subjects as ambiguous—they explicitly asked the experimenter how "and" was to be understood. Furthermore, when participants indicated how they understood "and" by shading Venn diagrams, it turned out that about a quarter of them interpreted "and" as union rather than intersection, which is usually assumed by experimenters using the Linda problem.

Implications for rule learning. The formation of conjunctions via "and" is a basic building block of rules. Its correct understanding is thus important for effective communication of results of rule learning. Existing studies suggest that the most common type of error is understanding "and" as a union rather than intersection. In such a case, a rule containing multiple "ands" will be perceived as having a higher support than it actually has. Each additional condition will be incorrectly perceived as increasing the coverage of the rule. This implies higher perceived plausibility of the rule. Misunderstanding of "and" will thus generally increase the preference of rules with more conditions.

Debiasing techniques. According to Sides et al. [135] "and" ceases to be ambiguous when it is used to connect propositions rather than categories. The authors give the following example of a sentence which is not prone to misunderstanding: "IBM stock will rise tomorrow and Disney stock will fall tomorrow." A similar wording of rule learning results may be, despite its verbosity, preferred.

Mellers et al. [97] showed that using "bank tellers who are feminists" or "feminist bank tellers" rather than "bank tellers and feminists" as a category in the Linda problem (Figure 2) might reduce the likelihood of committing the conjunction fallacy. It follows that using different wording such as "and also" might also help reduce the danger of a misunderstanding of "and".

Representations that visually express the semantics of "and" such as decision trees may be preferred over rules, which do not provide such visual guidance.³

²As in "He invited friends and colleagues to the party"

 $^{^3}$ We find limited grounding for this proposition in the following: Conditions connected with

5.3. Averaging Heuristic

While the conjunction fallacy is most commonly explained by operation of the representativeness heuristic, the averaging heuristic provides an alternative explanation: it suggests that people evaluate the probability of a conjuncted event as the average of probabilities of the component events [38]. As reported by Fantino et al. [38], in their experiment "approximately 49% of variance in subjects' conjunctions could be accounted for by a model that simply averaged the separate component likelihoods that constituted a particular conjunction."

Implications for rule learning. When applying the averaging heuristic, an analyst may not fully realize the consequences of the presence of a low-probability condition for the overall likelihood of the set of conditions in the antecedent of the rule.

Consider the following example: Let us assume that the learning algorithm only adds independent conditions that have a probability of 0.8, and we compare a 3-condition rule to a 2-condition rule. Averaging would evaluate both rules equally, because both have an average probability of 0.8. A correct computation of the joint probability, however, shows that the longer rule is considerably less likely $(0.8^3 \text{ vs. } 0.8^2 \text{ because all conditions are assumed to be independent)}$.

Averaging can also affect same-length rules. Fantino et al. [38] derive from their experiments on the averaging heuristic that humans tend to judge "unlikely information [to be] relatively more important than likely information." Continuing our example, if we compare the above 2-condition rule with another rule with two features with more diverse probability values, e.g., one condition has 1.0 and the other has 0.6, then averaging would again evaluate both rules the same, but in fact the correct interpretation would be that the rule with equal probabilities is more likely than the other $(0.8^2 > 1.0 \times 0.6)$. In this case, the low 0.6 probability in the new rule would "knock down" the normative conjoint probability below the one of the rule with two 0.8 conditions.

Debiasing techniques. Experiments conducted by Zizzo et al. [172] showed that prior knowledge of probability theory, and a direct reminder of how probabilities are combined, are effective tools for decreasing the incidence of the conjunction fallacy, which is the hypothesized consequence of the averaging heuristic. A specific countermeasure for the biases caused by linear additive integration (weighted averaging) is the use of logarithm formats. Experiments conducted by Juslin et al. [77] show that recasting probability computation in terms of logarithm formats, thus requiring additive rather than multiplicative integration, improves probabilistic reasoning.

an arch in a tree are to be interpreted as simultaneously valid (i.e., arch means conjunction). A recent empirical study on comprehensibility of decision trees [119] does not consider ambiguity of this notation to be a systematic problem among the surveyed users.

5.4. Disjunction Fallacy

The disjunction fallacy refers to a judgment that is inconsistent with the disjunction rule, which states that the probability Pr(X) cannot be higher than the probability Pr(Z), where $Z = X \cup Y$ is a union of event X with another event Y.

In experiments reported by Bar-Hillel and Neter [11], X and Z were nested pairs of categories, such as Switzerland and Europe. Subjects read descriptions of people such as "Writes letter home describing a country with snowy wild mountains, clean streets, and flower decked porches. Where was the letter written?" It follows that since Europe contains Switzerland, Europe must be more likely than Switzerland. However, Switzerland was chosen as the more likely place by about 75% of the participants [11].

The disjunction fallacy is considered as another consequence of the representativeness heuristic [11]: "Which of two events—even nested events—will seem more probable is better predicted by their representativeness than by their scope, or by the level in the category hierarchy in which they are located." The description in the example is more representative of Switzerland than of Europe, so when people use representativeness as the basis for their judgment, they judge Switzerland to be a more likely answer than Europe, even though this judgment breaks the disjunction rule.

Implications for rule learning. In the context of data mining, it can be the case that the feature space is hierarchically ordered. The analyst can thus be confronted with rules containing attributes (literals) on multiple levels of granularity. Following the disjunction fallacy, the analyst will generally prefer rules containing more specific attributes, which can result in preference for rules with fewer backing instances and thus in weaker statistical validity.

Debiasing techniques. When asked to assign categories to concepts (such as land of origin of a letter) under conditions of certainty, people are known to prefer a specific category to a more general category that subsumes it, but only if the specific category is considered representative [11]: "whenever an ordering of events by representativeness differs from their ordering by set inclusion, there is a potential for an extension fallacy to occur." From this observation a possible debiasing strategy emerges: making the analysts aware of the taxonomical relation of the individual attributes and their values. For example, the user interface can work with the information that Europe contains Switzerland, possibly actively notifying the analyst on the risk of falling for the disjunctive fallacy. This intervention can be complemented by "training in rules" [88]. In this case, the analysts should be explained the benefits of larger supporting sample associated with more general attributes.

5.5. Base-rate Neglect

People tend to underweigh the evidence provided by base rates, which results in the so-called *base-rate neglect*. For example, Kahneman and Tversky [80] gave participants a description of a person who was selected randomly from a group

and asked them whether the person is an engineer or a lawyer. Participants based their judgment mostly on the description of the person and paid little consideration to the occupational composition of the group, even though the composition as provided as part of the task and should play a significant role in the judgment.

Kahneman and Tversky [80] view the base-rate neglect as a possible consequence of the representativeness heuristic [79]. When people base their judgment of an occupation of a person mostly on similarity of the person to a prototypical member of the occupation, they ignore other relevant information such as base rates, which results in the base-rate neglect.

Implications for rule learning. The application of the base rate neglect suggests that when facing two otherwise identical rules with different values of confidence and support metrics, an analyst's preferences will be primarily shaped by the confidence of the rule. Support corresponds to "base rate", which is sometimes almost completely ignored [80].

It follows that by increasing preference for higher confidence, the base-rate neglect will generally contribute to a positive correlation between rule length and plausibility, since longer rules can better adapt to a particular group in data and thus have a higher confidence than a more general, shorter rules. This is in contrast to the general bias for simple rules that are implemented by state-of-the-art rule learning algorithms, because simple rules tend to be more general, have a higher support, and are thus statistically more reliable.

Debiasing techniques. Gigerenzer and Hoffrage [62] show that representations in terms of natural frequencies, rather than conditional probabilities, facilitate the computation of cause's probability. Confidence is typically presented as percentage in current software systems. The support rule quality metric is sometimes presented as a percentage and sometimes as a natural number. It would foster correct understanding if analysts are consistently presented natural frequencies in addition to percentages.

5.6. Insensitivity to Sample Size

People tend to underestimate the increased benefit of higher robustness of estimates that are made on a larger sample, which is called insensitivity to sample size. The insensitivity to sample size effect can be illustrated by the so-called hospital problem. In this problem, subjects are asked which hospital is more likely to record more days in which more than 60 percent of the newborns are boys. The options are a larger hospital, a smaller hospital, or both hospitals with about a similar probability. The correct expected answer—the smaller hospital—was chosen only by 22% of participants in an experiment reported by Tversky and Kahneman [151]. Insensitivity to sample size may be another bias resulting from use of the representativeness heuristic [79]. When people use the representativeness heuristic, they compare the proportion of newborns who are boys to the proportion expected in the population, ignoring other relevant

information. Since the proportion is similarly representative of the whole population for both hospitals, most of the participants believed that both hospitals are equally likely to record days in which more than 60 percents of the newborns are boys [151].

Implications for rule learning. This effect implies that analysts may be unable to appreciate the increased reliability of the confidence estimate with increasing value of support, i.e., they may fail to appreciate that the strength of the connection between antecedent and consequent of a rule rises with an increasing number of observations. If confronted with two rules, where one of them has a slightly higher confidence and the second rule a higher support, this cognitive bias suggests that the analyst will prefer the rule with higher confidence (all other factors equal).

In the context of this bias, it is important to realize that population size is statistically irrelevant for determination of sample size for large populations [26]. However, previous research [9] has shown that the perceived sample accuracy can incorrectly depend on the sample-to-population ratio rather than on the absolute sample size. For a small population, a 10% sample can be considered as more reliable than 1% sample drawn from much larger population.

This observation has substantial consequences for the presentation of rule learning results. The support of a rule is typically presented as a percentage of the dataset size. Assuming that support relates to sample size and number of instances in the dataset to population size, it follows that the presentation of support as a percentage (relative support) induces the insensitivity to sample size effect. The recommended alternative is to present support as an absolute number (absolute support).

Debiasing techniques. There have been successful experiments with providing decision aids to overcome the insensitivity to sample size bias. In particular, Kachelmeier and Messier Jr [78] experimented with providing auditors a formula for computing appropriate sample size for substantive tests of details based on the description of a case and tolerable error. Provision of the aid resulted in larger sample sizes being selected by the auditors in comparison to intuitive judgment without the aid. Similarly, as the auditor can choose the sample size, a user of an association rule learning algorithm can specify the minimum support threshold. To leverage the debiasing strategy validated by Kachelmeier and Messier Jr [78], the rule learning interface should also inform the user of the effects of chosen support threshold on the accuracy of the confidence estimate of the resulting rules. For algorithms and workflows where the user cannot influence the support of a discovered rule, relevant information should be available as a part of rule learning results. In particular, the value of rule support can be used to compute a confidence interval for the value of confidence. Such supplementary information is already provided by Bayesian decision lists [90], a recently proposed algorithmic framework positively evaluated with respect to interpretability (cf., e.g., [28]).

5.7. Confirmation Bias and Positive Test Strategy

Confirmation bias refers to the notion that people tend to look for evidence supporting the current hypothesis, disregarding conflicting evidence. According to Evans [35, p. 552] confirmation bias is "the best known and most widely accepted notion of inferential error of human reasoning."

Research suggests that even neutral or unfavorable evidence can be interpreted to support existing beliefs, or, as Trope et al. [148, p. 115–116] put it, "the same evidence can be constructed and reconstructed in different and even opposite ways, depending on the perceiver's hypothesis."

A closely related phenomenon is the *positive test strategy* (PTS) described by Klayman and Ha [82]. This reasoning strategy suggests that when trying to test a specific hypothesis, people examine cases which they expect to confirm the hypothesis rather than the cases which have the best chance of falsifying it. The difference between PTS and confirmation bias is that PTS is applied to test a candidate hypothesis while confirmation bias is concerned with hypotheses that are already established [121, p. 93]. The experimental results of Klayman and Ha [82] show that under realistic conditions, PTS can be a very good heuristic for determining whether a hypothesis is true or false, but it can also lead to systematic errors if applied to an inappropriate task.

Implications for rule learning. This bias can have a significant impact depending on the purpose for which the rule learning results are used. If the analyst has some prior hypothesis before she obtains the rule learning results, according to the confirmation bias she will tend to "cherry pick" rules confirming this prior hypothesis and disregard rules that contradict it. Given that some rule learners may output contradicting rules, the analyst may tend to select only the rules conforming to the hypothesis, disregarding applicable rules with the opposite conclusion, which could otherwise turn out to be more relevant.

Debiasing techniques. Delaying final judgment and slowing down work has been found to decrease confirmation bias in several studies [140, 118]. User interfaces for rule learning should thus give the user not only the opportunity to save or mark interesting rules, but also allow the user to review and edit the model at a later point in time. An example rule learning system with this specific functionality is EasyMiner [158].

Wolfe and Britt [169] successfully experimented with providing subjects with explicit guidelines for considering evidence both for and against a hypothesis. Provision of "balanced" instructions to search evidence for and against a given hypothesis reduced the incidence of myside bias, an effect closely related to confirmation bias, from 50% exhibited by the control group to a significantly lower 27.5%.

Similarly, providing explicit guidance combined with modifications of the user interface of the system presenting the rule learning results could also be

⁴Cited according to Nickerson [110].

considered. The assumption that educating users about cognitive illusions can be an effective debiasing technique for positive test strategy has been empirically validated on a cohort of adolescents by Barberia et al. [12].

5.8. Availability Heuristic

The availability heuristic is a judgmental heuristic in which a person evaluates the frequency of classes or the probability of events by the ease with which relevant instances come to mind. This heuristic is explained by its discoverers, Tversky and Kahneman [150], as follows: "That associative bonds are strengthened by repetition is perhaps the oldest law of memory known to man. The availability heuristic exploits the inverse form of this law, that is, it uses the strength of the association as a basis for the judgment of frequency." The availability heuristic is not itself a bias, but it may lead to biased judgments when availability is not a valid cue.

In one of the original experiments, participants were asked whether the letter "R" appears more frequently on the first or third position in English texts [150]. About 70% of participants answered incorrectly that it appears more frequently on the first position, presumably because they estimated the frequency by recalling words containing "R" and it is easier to recall words starting with R than words with R on the third position.

While original research did not distinguish between the number of recollected instances and ease of the recollection, later studies showed that to determine availability, it is sufficient to assess the ease with which instances or associations could be brought to mind; it is not necessary to count all the instances one is able to come up with [131].

Implications for rule learning. An application of the availability heuristic in rule learning would be based on the ease of recollection of instances (examples) matching the *complete* rule (all conditions and consequent) by the analyst. Rules containing conditions for which instances can be easily recalled would be found more plausible compared to rules not containing such conditions. As an example, consider the rule pair

```
R_1: IF latitude \leq 44.189 AND longitude \leq 6.3333 AND longitude > 1.8397 THEN Unemployment is high R_2: IF population \leq 5 million THEN Unemployment is high.
```

It is arguably easier to recall specific countries matching the second rule, than countries matching the conditions of the first rule.

It is conceivable that the availability heuristic could also be applied in case when the easily recalled instances match *only some of the conditions* in the antecedent of the rule, such as only latitude in the example above. The remaining conditions would be ignored.

On the other hand, such a bias can also be implemented as a bias into rule learning algorithms. Often, in particular in cases where many candidate conditions are available, such as datasets with features derived from the semantic web [127], the same information can be encoded in rules that use different sets

of conditions. For example, Gabriel et al. [53] proposed an algorithm that gives preference to selecting conditions that are semantically coherent. A similar technique could be used for realizing a preference for attributes that are easier to recall for human analysts.

Debiasing techniques. Several studies have found that people use ease of recollection in judgment only when they cannot attribute it to a source that should not influence their judgment [130]. Alerting an analyst to the reason why instances matching the conditions in the rule under consideration are easily recalled should therefore reduce the impact of the availability heuristic as long as the reason is deemed irrelevant to the task at hand.

5.9. Reiteration Effect, Effects of Validity and Illusiory Truth

The reiteration effect describes the phenomenon that repeated statements tend to become more believable [71, 116]. For example, in one experiment, Hasher et al. [69] presented subjects with general statements and asked them to asses their validity. Part of the statements were false and part were true. The experiment was conducted in several sessions, where some of the statements were repeated in subsequent sessions. The average perceived validity of both true and false repeated statements rose between the sessions, while for non-repeated statements it dropped slightly.

The effect is usually explained by use of processing fluency in judgment. Statements that are processed fluently (easily) tend to be judged as true and repetition makes processing easier. A recent alternative account argues that repetition makes the referents of statements more coherent and people judge truth based on coherency [155].

The reiteration effect is also known under different labels, such as "frequency-validity" or "illusory truth" [71, 195]. However, some research suggests that these are not identical phenomena. For example, the *truth effect* "disappears when the actual truth status is known" [122, p. 253], which does not hold for validity effect in general. There is also a clear distinction between the effects covered here, and the mere exposure effect covered in Section 5.10: the truth effect has been found largely independent of duration of stimulus exposure [29, p. 245].

Implications for rule learning. In the rule learning context, a repeating statement which becomes more believable corresponds to the entire rule or possibly a "subrule" consisting of the consequent of the rule and a subset of conditions in its antecedent. A typical rule learning result contains multiple rules that are substantially overlapping. If the analyst is exposed to multiple similar statements, the reiteration effect will increase the analyst's belief in the *repeating* subrule. Especially in the area of association rule learning, a very large set of redundant rules—covering the same, or nearly same set of examples—is routinely included in the output.

Schwarz et al. [132] suggest that mere 30 minutes of delay can be enough for information originally seen as negative to have positive influence. Applying this in a data exploration task, consider an analyst who is presented a large number of "weak" rules corresponding to highly speculative patterns of data. Even if the analyst rejects the rule—for example based on the presented metrics, pre-existing domain knowledge or common sense—the validity and truthfulness effects will make the analyst more prone to accept a similar rule later.

Debiasing techniques. The reiteration effect can be suppressed already on the algorithmic level by ensuring that rule learning output does not contain redundant rules. This can be achieved by pruning algorithms [47]. Another possible technique is presenting the result of rule learning in several layers, where only clusters of rules ("rule covers") summarizing multiple sub rules are presented at first [114]. The user can expand the cluster to obtain more similar rules. A more recent algorithm that can be used for summarizing multiple rules is the meta-learning method proposed by [16].

Several lessons can be learnt from Hess and Hagen [73], who studied the role of the reiteration effect for spreading of gossip. Interestingly, already simple reiteration was found to increase gossip veracity, but only for those who found the gossip relatively uninteresting. Multiple sources of gossip were found to increase its veracity, especially when these sources were independent. Information that explained the gossip by providing benign interpretation decreased the veracity of gossip. These findings suggest that it is important to explain to the analyst which rules share the same source, i.e. what is the overlap in their coverage in terms of specific instances. Second, explanations can be improved by utilisation of recently proposed techniques that use domain knowledge to filter or explain rules, such as expert deduction rules proposed by Rauch [126].

The research related to debiasing validity and truth effects has been largely centered around the problem of debunking various forms of misinformation (cf., e.g., [132, 91, 32]). The current largely accepted recommendation is that to correct a misinformation, it is best to address it directly – repeat the misinformation along with arguments against it [91, 32]. This can be applied, for example, in incremental machine learning settings, when the results of learning are revised when new data arrive, or when mining with formalized domain knowledge. Generally, when the system has knowledge of the analyst being previously presented a rule (a hypothesis), which is falsified following the current state of knowledge, the system can explicitly notify the analyst, listing the rule in question and explaining why it does not hold.

5.10. Mere Exposure Effect

According to the mere exposure effect, repeated exposure to an object results in an increased preference (liking, affect) for that object. When a concrete stimulus is repeatedly exposed, the preference for that stimulus increases logarithmically as a function of the number of exposures [20]. The size of the mere exposure effect also depends on whether the stimulus the subject is exposed to is exactly the same as in prior exposure or only similar to it [103]—the same stimuli are associated with larger mere exposure effect. The mere exposure effect is another consequence of increased fluency of processing associated with

repeated exposure (cf. Section 5.9) [167]. While the reitaration effect referred to the use of processing fluency in judgment of truth, the mere exposure effect relates to the positive feeling that is associated with fluent processing.

Duration of the exposure below 1 second produces the strongest effects, with increasing time of exposure the effect drops and repeating exposures decrease the mere exposure effect. The liking induced by the effect drops more quickly with increasing exposures when the presented stimuli is simple (e.g., an ideogram) as opposed to complex (e.g., a photograph) [20]. A recent meta analysis suggests that there is an inverted-U shaped relation between exposure and affect [104].

Implications for rule learning. The extent to which the mere exposure effect can affect the interpretation of rule leaning results is limited by the fact that that its magnitude decreases with extended exposure to the stimuli. It can be expected that the analysts inspect the rule learning results for a much longer period of time than the 1 second below which exposure results in the strongest effects [20]. However, it is not unusual for rule-based models to be composed of several thousand rules [3]. When the user scrolls through a list of rules, each rule can be shown only for a fraction of a second. The analyst is not aware of having seen the rule, yet the rule can influence the analyst's judgment through the mere exposure effect.

The mere exposure effect can also play a role when rules from the text mining or sentiment analysis domains are interpreted. The initial research of the mere exposure effect by Zajonc [170] included experimental evidence on the positive correlation between word frequency and affective connotation of the word. From this it follows that a rule containing frequently occurring words can induce the mere exposure effect.

Debiasing techniques. While there is a considerable body of research focusing on the mere exposure effect, our literature survey did not result in any directly applicable debiasing techniques. Only recently, Becker and Rinck [15] reported the first reversal of the mere exposure effect. This was achieved by presenting threatening materials (spider pictures) to people fearful of spiders in an unpleasant detection situation. This result, although interesting, is difficult to transpose to the domain of rules.

Nevertheless, there are some conditions known to decrease the mere exposure effect that can be utilized in machine learning interfaces. The effect is strongest for repeated, "flash-like" presentation of information. A possible workaround is to avoid subliminal exposure completely, by changing the mode of operation of the corresponding user interfaces. One attempt at a user interface to rule learning respecting these principles is the EasyMiner system [159]. In EasyMiner, the user precisely formulates the mining task as a query against data. This restricts the number of rules that are discovered and the user is consequently exposed to.

5.11. Overconfidence and underconfidence

A decision maker's judgment is normally associated with belief that the judgment is true, i.e., with confidence in the judgment. Griffin and Tversky [65]

argue that confidence in judgment is based on a combination of the strength of evidence and its weight (credibility). According to their studies, people tend to combine strength with weight in suboptimal ways, resulting in the decision maker being too much or too little confident about the hypothesis at hand than would be normatively appropriate given the available information. This discrepancy between the normative confidence and the decision maker's confidence is called *overconfidence* or *underconfidence*.

People use the provided data to assess a hypothesis, but they insufficiently regard the quality of the data. Griffin and Tversky [65] describe this manifestation of bounded rationality as follows: "If people focus primarily on the warmth of the recommendation with insufficient regard for the credibility of the writer, or the correlation between the predictor and the criterion, they will be overconfident when they encounter a glowing letter based on casual contact, and they will be underconfident when they encounter a moderately positive letter from a highly knowledgeable source."

Implications for rule learning. Research has revealed systematic patterns of overconfidence and underconfidence [65, p. 426]: If the estimated difference between two hypotheses is large, it is easy to say which one is better and there is a pattern of underconfidence. As the degree of difficulty rises (the difference between the normative confidence of two competing hypotheses is decreasing), there is an increasing pattern of overconfidence.

The strongest overconfidence was recorded for problems where the weight of evidence is low and the strength of evidence is high. This directly applies to rules with high value of confidence and low value of support. The empirical results related to the effect of difficulty therefore suggest that the predictive ability of such rules will be substantially overrated by analysts. This is particularly interesting because rule learning algorithms often suffer from a tendency to unduly prefer overly specific rules that have a high confidence on small parts of the data to more general rules that have a somewhat lower confidence, a phenomenon also known as overfitting. The above-mentioned results seem to indicate that humans suffer from a similar problem (albeit presumably for different reasons), which, e.g., implies that a human-in-the-loop solution may not alleviate this problem.

Debiasing techniques. Research applicable to debiasing of overconfidence originated in 1950', but most initial efforts to reduce overconfidence have failed [40, 7]. Some recent research focuses on the hypothesis that the feeling of confidence reflects factors indirectly related to choice processes [45, 67]. For example, in a sport betting experiment performed by Hall et al. [67], participants underweighted statistical cues while betting, when they knew the names of players. This research leads to the conclusion that "more knowledge can decrease accuracy and simultaneously increase prediction confidence" [67]. Applying this to debiasing in the rule learning context, presenting less information can be achieved by reducing the number of rules and removing some conditions in the remaining rules. This can be achieved by a number of methods, such as feature

selection to external setting of maximum antecedent length, which is permitted by some algorithms. Also, rules and conditions that do not pass a statistical significance test can be removed from the output.

As with other biases, research on debiasing overconfidence points at the importance of educating the experts on principles of subjective probability judgment and the associated biases [25]. Shafir [134, p. 487] recommends to debias overconfidence (in policy making) by making the subject hear both sides of an argument. In the rule learning context, this would correspond to the user interface making rules and knowledge easily accessible, which is in "unexpectedness" or "exception" relation with the rule in question, as, e.g., experimented with in frameworks postprocessing association rule learning results [83].

5.12. Recognition Heuristic

Pachur et al. [116] define the recognition heuristic as follows: "For twoalternative choice tasks, where one has to decide which of two objects scores higher on a criterion, the heuristic can be stated as follows: If one object is recognized, but not the other, then infer that the recognized object has a higher value on the criterion." In contrast with the availability heuristic, which is based on ease of recall, the recognition heuristic is based only on the fact that a given object is recognized. The two heuristics could be combined. When only one object in a pair is recognized, then the recognition heuristic would be used for judgment. If both objects are recognized, then the speed of the recognition could influence the choice [72].

The use of this heuristic could be seen from an experiment performed by Goldstein and Gigerenzer [64], which focused on estimating which of two cities in a presented pair is more populated. People using the recognition heuristic would say that the city they recognize has a higher population. The median proportion of judgments complying to the recognition heuristic was 93%. It should be noted that the application of this heuristic is in this case ecologically justified since recognition will be related to how many times the city appeared in a newspaper report, which in turn is related to the city size [14].

Implications for rule learning. The recognition heuristic can manifest itself by preference for rules containing a recognized attribute name or value in the antecedent of the rule. Analysts processing rule learning results are typically shown many rules, contributing to time pressure. This can further increase the impact of the recognition heuristic.

Empirical results reported by Michalkiewicz et al. [98] indicate that people with higher cognitive ability use the recognition heuristic more when it is successful and less when it is not. The work of Pohl et al. [123] shows that people adapt their decision strategy with respect to the more general environment rather than the specific items they are faced with. Considering that the application of the recognition heuristic can in some situations lead to better results than the use of available knowledge, the recognition heuristic may not necessarily have overly negative impacts on the intepretation of rule learning results.

Debiasing techniques. Under time pressure people assign a higher value to recognized objects than to unrecognized objects. This happens also in situations when recognition is a poor cue [115]. Changes to user interfaces that induce "slowing down" could thus help to address this bias. As to the alleviation of effects of recognition heuristic in situations where it is ecologically unsuitable, Pachur and Hertwig [115] note that suspension of the heuristic requires additional time or direct knowledge of the "criterion variable". In typical real-world machine learning tasks, the data can include a high number of attributes that even experts are not acquainted with in detail. When these are recognized (but not understood), even experts may be liable to the recognition heuristic. When information on the meaning of individual attributes and literals is made easily accessible, we conjecture that the application of the recognition heuristic can be suppressed.

5.13. Information Bias

Information bias refers to the tendency to seek more information to improve the perceived validity of a statement even if the additional information is not relevant or helpful. The typical manifestation of the information bias is evaluating questions as worth asking even when the answer cannot affect the hypothesis that will be accepted [13].

For example, Baron et al. [13] asked subjects to assess to what degree a medical test is suitable for deciding which of three diseases to treat. The test detected a chemical, which was with a certain probability associated with each of the three diseases. These probabilities varied across the cases. Even though in some of the cases an outcome of the test would not change the most likely disease and thus the treatment, people tended to judge the test as worth doing. While information bias is primarily researched in the context of information acquisition [109, 107], some scientists interpret this more generally as judging features with zero probability gain as useful, having potential to change one's belief [108, p. 158].

Implications for rule learning. Many rule learning algorithms allow the user to select the size of the generated model – in terms of the number of rules that will be presented, as well as by setting the maximum length of conditions of the generated rules. Either as part of the feature selection, or when defining constraints for the learning, the users decide which attributes are relevant. These can then appear among conditions of the discovered rules.

According to the information bias, people will be prone to setup the task so that they receive more information – resulting in larger rule list with longer rules containing attributes with little information value.

It is unclear if the information effect applies also to the case when the user is readily presented with more information, rather then given the possibility to request more information. Given the proximity of these two scenarios, we conjecture that information bias (or some related bias) will make people prefer more information to less, even if it is obviously not relevant.

According to the information bias, a rule containing additional (redundant) condition may be preferred to a rule not containing this condition.

Debiasing techniques. While informing people about the diagnosticity of considered questions does not completely remove the information bias, it reduces it [13]. To this end, communicating attribute importance can help guide the analyst in the task definition phase.

Although existing algorithms and systems already provide ways for determining the importance of individual rules, for example via values of confidence, support, and lift, the cues on the importance of individual conditions in rule antecedent are typically not provided. While feature importance is computed within many learning algorithms, it is often used only internally. Exposing this information to the user can help counter the information bias.

5.14. Ambiguity Aversion

Ambiguity aversion refers to the tendency to prefer known risks over unknown risks. This is often illustrated by the Ellsberg paradox [34], which shows that humans tend to systematically prefer a bet with known probability of winning over a bet with not precisely known probability of winning, even if it means that their choice is systematically influenced by irrelevant factors.

As argued by Camerer and Weber [22], ambiguity aversion is related to the information bias: the demand for information in cases when it has no effect on decision can be explained by the aversion to ambiguity — people dislike having missing information.

Implications for rule learning. The ambiguity aversion may have profound implications for rule learning. The typical data mining task will contain a number of attributes the analyst has no or very limited knowledge of. The ambiguity aversion will manifest itself in a preference for rules that do not contain ambiguous conditions.

Debiasing techniques. An empirically proven way to reduce ambiguity aversion is accountability – "the expectation on the side of the decision maker of having to justify her decisions to somebody else" [156]. This debiasing technique is hypothesized to work through higher cognitive effort that is induced by accountability.

This can be applied in the rule learning context by requiring the analysts to provide justifications for why they evaluated a specific discovered rule as interesting. Such explanation can be textual, but also can have a structured form. To decrease demands on the analyst, the explanation may only be required only if a conflict with existing knowledge has been automatically detected, for example, using approach proposed by Rauch [126].

Since the application of the ambiguity aversion can partly stem from the lack of knowledge of the conditions included in the rule, it is conceivable this bias would be alleviated if description of the meaning of the conditions is made easily accessible to the analyst, as demonstrated in e.g. [83].

5.15. Confusion of the Inverse

This effect corresponds to confusing the probability of cause and effect, or, formally, confidence of an implication $A \to B$ with its inverse $B \to A$, i.e., $\Pr(B|A)$ is confused with the inverse probability $\Pr(A \mid B)$. For example, Villejoubert and Mandel [157] showed in an experiment that about half of the participants estimating the probability of membership in a class gave most of their estimates that corresponded to the inverse probability.

Implications for rule learning. The confusion of the direction of an implication sign has significant consequences on the interpretation of a rule. Already Michalski [100] noted that there are two different kinds of rules, discriminative and characteristic. *Discriminative rules* can quickly discriminate an object of one category from objects of other categories. A simple example is the rule

IF trunk THEN elephant

which states that an animal with a trunk is an elephant. This implication provides a simple but effective rule for recognizing elephants among all animals.

Characteristic rules, on the other hand, try to capture all properties that are common to the objects of the target class. A rule for characterizing elephants could be

IF elephant THEN heavy, large, grey, bigEars, tusks, trunk.

Note that here the implication sign is reversed: we list all properties that are implied by the target class, i.e., by an animal being an elephant. From the point of understandability, characteristic rules are often preferable to discriminative rules. For example, in a customer profiling application, we might prefer to not only list a few characteristics that discriminate one customer group from the other, but are interested in all characteristics of each customer group.

Characteristic rules are very much related to formal concept analysis [165, 55]. Informally, a concept is defined by its intent (the description of the concept, i.e., the conditions of its defining rule) and its extent (the instances that are covered by these conditions). A formal concept is then a concept where the extension and the intension are Pareto-maximal, i.e., a concept where no conditions can be added without reducing the number of covered examples. In Michalski's terminology, a formal concept is both discriminative and characteristic, i.e., a rule where the head is equivalent to the body.

The confusion of the inverse thus seems to imply that humans will not clearly distinguish between these types of rules, and, in particular, tend to interpret an implication as an equivalence. From this, we can infer that characteristic rules, which add all possible conditions even if they do not have additional discriminative power, may be preferable to short discriminative rules.

This confusion may manifest itself strongest in the area of association rule learning, where an attribute can be of interest to the analyst both in the antecedent and consequent of a rule.

Debiasing techniques. Edgell et al. [33] studied the influence of the effect of training of analysts in probabilistic theory with the conclusion that it is not effective in addressing the confusion of the inverse fallacy.

Werner et al. [163, p. 195] point at a concern regarding use of language liable to misinterpretation in statistical textbooks teaching fundamental concepts such as independence. The authors illustrate the misinterpretation on the statement whenever Y has no effect on X as "This statement is used to explain that two variables, X and Y, are independent and their joint distribution is simply the product of their margins. However, for many experts, the term 'effect' might imply a causal relationship." From this it follows that representations of rules should strive for unambiguous meaning of the wording of the implication construct. The specific recommendations provided by Díaz et al. [31] for teaching probability can also be considered in the next generation of textbooks aimed at the data science audience.

5.16. Context and Tradeoff Contrast Effects

People evaluate objects in relation to other available objects, which may lead to various effects of context of presentation of a choice. For example, in one of the experiments described by Tversky and Simonson [154], subjects were asked to choose between two microwave ovens (Panasonic priced 180 USD and Emerson priced 110 USD), both a third off the regular price. The number of subjects who chose Emerson was 57% and 43% chose Panasonic. Another group of subjects was presented the same problem with the following manipulation: A more expensive Panasonic valued at 200 USD (10% off the regular price) was added to the list of possible options. The newly added device was described to look as inferior to the other Panasonic, but not to the Emerson device. After this manipulation, only 13% chose the more expensive Panasonic, but the number of subjects choosing the less expensive Panasonic rose from 43% to 60%. That is, even though the additional option was dominated by the cheaper Panasonic device and it should have been therefore irrelevant to the relative preference of the other ovens, its addition changed the preference in favor of the better Panasonic device. The experiment thus shows that selection of one of the available alternatives, such as products or job candidates, can be manipulated by addition or deletion of alternatives that are otherwise irrelevant. Tversky and Simonson [154] attribute the tradeoff effect to the fact that "people often do not have a global preference order and, as a result, they use the context to identify the most 'attractive' option."

It should be noted that according to Tversky and Simonson [154] if people have well-articulated preferences, the background context has no effect on the decision.

Implications for rule learning. The effect could be illustrated on the inter-rule comparison level. In the base scenario, a constrained rule learning yields only a rule R_1 with a confidence value of 0.7. Due to the relatively low value of confidence, the user does not find the rule very plausible. By lowering the minimum confidence threshold, multiple other rules predicting the same target

class are discovered and shown to the user. These other rules, inferior to R_1 , would increase the plausibility of R_1 by the tradeoff contrast effect.

Debiasing techniques. Marketing professionals sometimes introduce more expensive versions of the main product, which induces the tradeoff contrast. The presence of a more expensive alternative with little added value increases sales of the main product [136]. Somewhat similarly, a rule learning algorithm can have on its output rules with very high confidence, sometimes even 1.0, but very low values of support. Removal of such rules can help to debias the analysts.

The influence of context can in some cases improve communication [136, p. 293]. An attempt at making contextual attributes explicit in the rule learning context was made by Gamberger and Lavrač [54], who introduced *supporting factors* as a means for complementing the explanation delivered by conventional learned rules. Essentially, supporting factors are additional attributes that are not part of the learned rule, but nevertheless have very different distributions with respect to the classes of the application domain. In line with the results of Kononenko [84], medical experts found that these supporting factors increase the plausibility of the found rules.

5.17. Negativity Bias

According to the negativity bias, negative evidence tends to have a greater effect than neutral or positive evidence of equal intensity [129].

For example, the experiments by Pratto and John [125] investigated whether the valence of a word (desirable or undesirable trait) has effect on the time required to identify the color in which the word appears on the screen. The results showed that the subjects took longer to name the color of an undesirable word than for a desirable word. The authors argued that the response time was higher for undesirable words because undesirable traits get more attention. Information with negative valence is given more attention partly because people seek diagnostic information, and negative information is more diagnostic [137]. Some research suggests that negative information is better memorized and subsequently recognized [128, 113].

Implications for rule learning. An interesting applicable discovery shows that negativity is an "attention magnet" [42, 113]. This implies that a rule predicting a class phrased with negative valence will get more attention than a rule predicting a class phrased with words with positive valence.

Debiasing techniques. Putting a higher weight to negative information may in some situations be a valid heuristic. What needs to be addressed are cases, when the relevant piece of information is positive and a less relevant piece of information is negative [74, 152]. It is therefore advisable that any such suspected cases are detected in the data preprocessing phase, and the corresponding attributes or values are replaced with more neutral sounding alternatives.

5.18. Primacy Effect

Once people form an initial assessment of plausibility (favorability) of an option, its subsequent evaluations will reflect this initial disposition.

Bond et al. [19] investigated to what extent changing the order of information which is presented to a potential buyer affects the propensity to buy. For example, in one of the experiments, if the positive information (product description) was presented as first, the number of participants indicating they would buy the product was 48%. When the negative information (price) was presented first, this number decreased to 22%. Bond et al. [19] argue that the effect is caused by distortion of interpretation of new information in the direction of the already held opinion. The information presented first not only influences disproportionately the final opinion, but it also influences interpretation of novel information.

Implications for rule learning. Following the primacy effect, the analyst will favor rules that are presented as first in the rule model. Largest negative effects of this bias are likely to occur, when such ordering is not observed, for example, when rules are presented in the order in which they were discovered by a breadth-first algorithm. In this case, mental contamination is another applicable bias related to the primacy effect (or in general order effects). This refers to the case when a presented hypothesis can influence subsequent decision making by its content, even if the subject is fully aware of the fact that the presented information is purely speculative [43]. Note that our application scenario differs from [43] and some other related research, in that cognitive psychology mostly investigated the effect of asking a hypothetical question, while we are concerned with considering the plausibility of a presented hypothesis (inductively learnt rule). Fitzsimons and Shiv [43] found that respondents are not able to prevent the contamination effects of the hypothetical questions and that the bias increases primarily when the hypothetical question is relevant. This bias is partly attributed to the application of expectations related to conversational maxims [63].

Debiasing techniques. Three types of debiasing techniques were examined by Mumma and Wilson [105] in the context of clinical-like judgments. The bias inoculation intervention involves direct training on the applicable bias or biases, consisting of information on the bias, strategies for adjustment, as well as completing several practical assignments. The second technique was considerthe-opposite debiasing strategy, which sorts the information according to diagnosticity before it is reviewed. The third strategy evaluated was simply taking notes when reviewing each cue before the final judgment was made. Interestingly, bias inoculation, a representative of direct debiasing techniques, was found to be the least effective. Consider-the-opposite and taking notes were found to work equally well.

To this end, a possible debiasing strategy can be founded in presentation of the most relevant rules first. Similarly, the conditions within the rules can be ordered by predictive power. Some rule learning algorithms, such as CBA [93], readily take advantage of the primacy effect, since they naturally create rule models that contain rules sorted by their strength. Other algorithms order rules so that more general rules (i.e., rules that cover more examples) are presented first. This typically also corresponds to the order in which rules are learned with the commonly used separate-and-conquer or covering strategies [48]. Simply reordering the rules output by these algorithms may not work in situations, when rules compose a rule list that is automatically processed for prediction purposes.⁵ In order to take advantage of the note taking debiasing strategy, the user interface can support the analyst in annotating the individual rules.

Lau and Coiera [89] provide a reason for optimism concerning the debiasing effect stemming from the proposed changes to user interface of machine learning tools. Their paper showed debiasing effect of similar changes implemented in a user interface to an information retrieval system used by consumers to find health information. Three versions of the system were compared: a baseline "standard" search interface, anchor debiasing interface, which asked the users to annotate the read documents as providing evidence for/against/neutral the proposition in question. Finally, the order debiasing interface reordered the documents to neutralize the primacy bias by creating a "counteracting order bias". This was done by randomly reshuffling a part of the documents. When participants used the baseline and anchor debiasing interface, the order effect was present. On the other hand, the use of the order debiasing interface eliminated the order effect [89].

5.19. Weak Evidence Effect

According to the weak evidence effect, presenting weak evidence in favor of an outcome can actually decrease the probability that a person assigns to the outcome. For example, in an experiment in the area of forensic science reported by Martire et al. [96], it was shown that participants presented with evidence weakly supporting guilt tended to "invert" the evidence, thereby counterintuitively reducing their belief in the guilt of the accused. Fernbach et al. [39] argue that the effect occurs because people give undue weight to the weak evidence and fail to take into account alternative evidence that more strongly favors the hypothesis at hand.

Implications for rule learning. The weak evidence effect can be directly applied to rules: the evidence is represented by the rule antecedent; the consequent corresponds to the outcome. The analyst can intuitively interpret each of the conditions in the antecedent as a piece of evidence in favor of the outcome.

⁵One technique that can positively influence comprehensibility of the rule list is prepending (adding to the beginning) a new rule to the previously learned rules [162]. The intuition behind this argument is that there are often simple rules that would cover many of the positive examples, but also cover a few negative examples that have to be excluded as exceptions to the rule. Placing the simple general rule near the end of the rule list allows us to handle exceptions with rules that are placed before the general rule and keep the general rule simple.

Typical of many machine learning problems is the uneven contribution of individual attributes to the prediction. Let us assume that the analyst is aware of the prediction strength of the individual attributes. If the analyst is to choose from a rule containing only one strong condition (predictor) and another rule containing a strong predictor and a weak (weak enough to trigger this effect) predictor, according to the weak evidence effect the analyst should choose the shorter rule with one predictor.

Debiasing techniques. Martire et al. [95] performed an empirical study aimed at evaluating what mode of communication of the strength of evidence is most resilient to the weak evidence effect. The surveyed modes of expression were numerical, verbal, a table, and a visual scale. It should be noted that the study was performed in the specific field of assessing evidence by a juror in a trial and the verbal expressions were following standards proposed by the Association of Forensic Science Providers [166].⁶ The results clearly suggested that numerical expressions of evidence are most suitable for expressing uncertainty.

Likelihood ratios studied by Martire et al. [95] are conceptually close to the lift metric, used to characterize association rules. While lift is still typically presented as a number in machine learning user interfaces, there has been research towards communicating rule learning results in natural language since at least 2005 [144]. With recent resurgence of interest in interpretable models, the use of natural language has been taken up by commercial machine learning services, such as BigML, which allow to generate predictions via spoken questions and answers using Amazon Alexa voice service. Similarly, machine learning interfaces increasingly rely on visualizations. The research on debiasing of the weak evidence effect suggests that when conveying machine learning results using modern means, such as transformation to natural language or through visualizations, care must be taken when numerical information is communicated.

Martire et al. [95] also observe high level of miscommunication associated with low-strength verbal expressions. In these instances, it is "appropriate to question whether expert opinions in the form of verbal likelihood ratios should be offered at all" [95]. Transposing this result to the machine learning context, we suggest to consider intentional omission of weak predictors from rules either directly by the rule learner or as part of feature selection.

5.20. Unit Bias

The unit bias refers to the tendency to give each unit similar weight while ignoring or underweighing the size of the unit [56].

Geier et al. [56] offered people various food items in two different sizes on different days and observed how this would affect consumption of the food. They found that people are larger amount of food when the size of a single unit of

 $^{^6}$ These provide guidelines on translation of numerical likelihood ratios into verbal formats. For example, likelihood "> 1-10" is translated as "weak or limited", and likelihood of "1000-10,000" as "strong".

⁷https://bigml.com/tools/alexa-voice

the food item was big than when it was small. A possible explanation is that people ate one unit of food at a time without taking into account how big it was. Because the food was not consumed in larger amounts at any single occasion, but was rather eaten intermittently, the behavior led to higher consumption when a unit of food was larger.

Implications for rule learning. Unit bias was so far primarily studied for quite different purposes than is the domain of machine learning. Nevertheless, as we will argue in the following, it can be very relevant for the domain of rule learning.

From a technical perspective, the number of conditions in rules is not important. What matters is the actual discriminatory power of the individual conditions, which can vary substantially. However, following the application of unit bias, people can view conditions as units of similar importance, disregarding their sometimes vastly different discriminatory and predictive power.

Debiasing techniques. One of the common ways how regulators address unhealthy food consumption patterns related to varying sizes of packaging is introduction of mandatory labelling of the size and calorie contents. Following an analogy to clearly communicating the size of food item, informing analysts about the discriminatory power of the individual conditions may alleviate unit bias. Such indicator can be generated automatically, for example, by listing the number of instances in the entire dataset that meet the condition.

6. Recommendations for Rule Learning Algorithms and Software

This section provides a concise list of considerations that is aimed to raise awareness among machine learning practitioners regarding the availability of measures that could potentially suppress effect of cognitive biases on comprehension of rule-based models. We expect part of the list to be useful also for other symbolic machine learning models, such as decision trees. In our recommendations, we focus on systems that present the rule model to a human user, which we refer to as the analyst. We consider two basic roles the analyst can have in the process: approval of the complete classification model ("interpretable classifiation task"), and selection of interesting rules ("nugget discovery").

6.1. Representation of a rule

The interpretation of natural language expressions used to describe a rule can lead to systematic distortions. Our review revealed the following recommendations applicable to individual rules:

1. Syntactic elements. There are several cognitive studies indicating that AND is often misunderstood [70], [59, p. 95-96]. The results of our experiments [52] support the conclusion that AND needs to be presented unambiguously in the rule learning context. Research has shown that "and" ceases to be ambiguous when it is used to connect propositions

rather than categories. Similarly, the communication of the implication construct IF THEN connecting antecedent and consequent should be made unambiguous.

Another important syntactic construct is negation (NOT). While processing of negation has not been included among the surveyed biases, our review of literature (cf. Section 7.5) suggests that its use should be discouraged on the grounds that its processing requires more cognitive effort, and because the fact that a specific information was negated may not be remembered in the long term.

2. Conditions. Attribute-value pairs comprising conditions are typically either formed of words with semantics meaningful to the user, or of codes that are not directly meaningful. When conditions contain words with negative valence, these need to be reviewed carefully, since negative information is known to receive more attention and is associated with higher weight than positive information. A number of biases can be triggered or strengthened by the lack of understanding of attributes and their values appearing in rules. Providing easily accessible information on conditions in the rules, including their predictive power, can thus prove as an effective debiasing technique.

People have the tendency to put higher emphasis on information they are exposed to first. By ordering the conditions by strength, machine learning software can conform to human conversational maxims. The output could also visually delimit conditions in the rules based on their significance or predictive stength.

3. **Interestingness measures.** The values of interestingness measures should be communicated using numerical expressions. Alternate verbal expressions, with wordings such as "strong relationship" replacing specific numerical values, are discouraged because there is some evidence that such verbal expressions are prone to miscommunication.

Currently, rule interest measures are typically represented as probabilities (confidence) or ratios (lift), whereas results in cognitive science indicate that natural frequencies are better understood.

The tendency of humans to ignore base rates and sample sizes (which closely relate to rule support) is a well established fact in cognitive science. Results of our experiments on inductively learned rules also provide evidence for this conclusion [52]. Our proposition is that this effect can be addressed by presenting confidence (reliability) intervals for the values of measures of interest, where applicable.

6.2. Rule models

In many cases, rules are not presented in isolation to the analyst, but instead within a collection of rules comprising a rule model. Here, we relate the results of our review to the following aspects of rule models:

4. Model size. An experiment by Poursabzi-Sangdeh et al. [124] found that people are better able to simulate results of a larger regression model composed of eight coefficients than of a smaller model composed of two coefficients. The results indicate that removal of any unnecessary variables could improve model interpretability even though the experiment did not find a difference in the trust in the model based on the number of coefficients it consisted of. Similarly to regression models, rule models often incorporate output that is considered as marginally relevant. This can take a form of (nearly) redundant rules or (nearly) redundant conditions in the rule. Our analysis shows that such redundancies can induce a number of biases, which may be accountable for misinteretation of the model. Size of a rule model can be reduced by utilizing various pruning techniques, or by using learning algorithms that allow the user to set or influence size of the resulting model. Examples of such approaches include those proposed by Letham et al. [90], Lakkaraju et al. [87], Wang et al. [160]. The Interpretable Decision Sets algorithm [87] can additionally optimize for diversity and non-overlap of discovered rules, directly countering the reiteration effect.

Another potentially effective approach to discarding some rules can be using domain knowledge or constraints set by the user to remove the strong (e.g., highly confident), yet "obvious" rules confirming common knowledge.⁸ Removal of weak rules could help to address the tradeoff contrast as well as the weak evidence effect.

5. Rule grouping. The rule learning literature has seen multiple attempts to develop methods for grouping similar rules, often by clustering. Our review suggests that presenting clusters of similar rules can help to reduce cognitive biases caused by reiteration.

Algorithms that learn rule lists provide mandatory ordering of rules, while the rule order in rule-set learning algorithms is not important. In either case, the rule order as presented to the user will affect perception of the model due to conversational maxims and the primacy effect. It is recommended to sort the presented rules by strength. However, due to paucity of applicable research, it is unclear which particular definition of rule strength would lead to the best results in terms of bias mitigation.

6.3. User Engagement

Some results of our review suggest that increasing user interaction can help counter some biases. Some specific suggestions for machine learning user interfaces (UIs) follow:

7. **Domain knowledge**. Selectively presenting domain knowledge "conflicting" with the considered rule can help to invoke the 'consider-the-opposite'

⁸For example, it is well-known that diastolic blood pressure rises with body mass index (DBP↑↑BMI). Rules confirming this relationship might be removed [83].

debiasing strategy. Other research has shown that the plausibility of a model depends on compliance to monotonicity constraints [46]. We thus suggest that UIs make background information on discovered rules easily accessible.

- 8. Eliciting rule annotation. Activating the deliberate "System 2" is one of the most widely applicable debiasing strategies. One way to achieve this is to require accountability, e.g., through visual interfaces motivating users to annotate selected rules, which would induce the 'note taking' debiasing strategy. Giving people additional time to consider the problem has been in some cases shown as an effective debiasing strategy. This can be achieved by making the selection process (at least) two stage, allowing the user to revise the selected rules.
- 9. User search for rules rather than scroll. Repeating rules can affect users via the mere exposure effect even if they are exposed to them even for a short moment, e.g., when scrolling a rule list. The user interfaces should thus deploy alternatives to scrolling in discovered rules, such as search facilities.

6.4. Bias inoculation

In some studies, basic education about specific biases, such as brief tutorials, decreased the fallacy rate. This debiasing strategy has been called *bias inoculation* in the literature.

10. Education on specific biases. Several studies have shown that providing explicit guidance and education on formal logic, hypothesis testing, and critical assessment of information can reduce fallacy rates in some tasks. However, the effect of psychoeducational methods is still a subject of dispute [92], and cannot be thus recommended as a sole or sufficient measure.

7. Limitations and Future Work

Our goal was to examine whether cognitive biases can affect the interpretation of machine learning models and to propose possible remedies if they do. Since this field is untapped from the machine learning perspective, we tried to approach the problem holistically. Our work yielded a number of partial contributions, rather than a single profound result. We mapped applicable cognitive biases, identified prior works on their suppression, and proposed how these could be transferred to machine learning. In the following, we outline some promising direction of future work.

7.1. Validation through human-subject experiments

All the identified shortcomings of human judgment pertaining to the interpretation of inductively learned rules are based on empirical cognitive science research. For each cognitive bias, we provided a justification how it would relate to machine learning. Due to the absence of applicable prior research in the intersection between cognitive science and machine learning, this justification is mostly based on authors' experience in machine learning.

A critical next step is empirical validation of the selected cognitive biases. We have already described several user experiments aimed at validating selected cognitive biases in Fürnkranz et al. [52]. Some other machine learning researchers have reported human-subject experiments that do not explicitly refer to cognitive biases, yet the cognitive phenomena they investigate may correspond to a known cognitive bias. One example is a study by Lage et al. [86] (cf. also extended version in [106]), which investigated the effect of the number of cognitive chunks (conditions) in a rule on response time. While the main outcome confirms the intuition that higher complexity results in higher response times, this study has also revealed several unexpected patterns, such as that defining a new concept and reusing it leads to a higher response time than repeating the description whenever that concept implicitly appears, even though this repetition means that subjects have to read more lines. The findings could possibly be attributed to fluency in judgement, a cognitive phenomenon assumed to underlie multiple cognitive biases.

Despite the existence of several early studies, much more concentrated and systematic effort is needed to yield insights on the size of effect individual biases can have on understanding of machine learning models.

7.2. Role of Domain Knowledge

It has been long recognized that external knowledge plays an important rule in the rule learning process. Already Mitchell [102] recognized at least two distinct roles external knowledge can play in machine learning: it can constrain the search for appropriate generalizations, and guide learning based on the intended use of the learned generalizations. Interaction with domain knowledge has played an important role in multiple stages of the machine learning process. For example, it can improve semi-supervised learning [23], and in some applications it is vital to convert discovered rules back into domain knowledge [50, p. 288]. Some results also confirm the common intuition that compliance to constraints valid in the given domain increases the plausibility of the learned models [46].

Our review shows that domain knowledge can be one of the important instruments in the toolbox aimed at debiasing interpretation of discovered rules. To give a specific example, the presence or strength of the validity effect depends on the familiarity of the subject with the topic area from which the information originates [18]. Future work should focus on a systematic review of the role of domain knowledge on activation or inhibition of cognitive phenomena applicable to interpretability of rule learning results.

7.3. Individual Differences

The presence of multiple cognitive biases and their strengths have been linked to specific personality traits. For example, overconfidence and the rate of conjunctive fallacy have been shown to be inversely related to numeracy [168]. According to Juslin et al. [77], the application of the averaging heuristic rather than the normative multiplication of probabilities seems to depend on the working memory capacity and/or high motivation.

Some research can even be interpreted as indicating that data analysts can be more susceptible to the myside bias than the general population. An experiment reported by Wolfe and Britt [169] shows that subjects who defined good arguments as those that can be proved by facts (this stance, we assume, would also apply to many data analysts) were more prone to exhibiting the myside bias. Stanovich et al. [141] show that the incidence of myside bias is surprisingly not related to general intelligence. This suggests that even highly intelligent analysts can be affected. Albarracín and Mitchell [2] propose that the susceptibility to the confirmation bias can depend on one's personality traits. They also present a diagnostic tool called "defense confidence scale" that can identify individuals who are prone to confirmational strategies. Further research into personality traits of users of machine learning outputs, as well as into development of appropriate personality tests, would help to better target education focused on debiasing.

7.4. Incorporating Additional Biases

There are about 24 cognitive biases covered in *Cognitive Illusions*, the authoritative overview of cognitive biases by Pohl [122], and even 51 different biases are covered by Evans et al. [37]. While doing the initial selection of cognitive biases to study, we tried to identify those most relevant for machine learning research matching our criteria. In the end, our review focused on a selection of 20 cognitive biases (effects, illusions). Future work might focus on expanding the review with additional relevant biases, such as labelling and overshadowing effects [122].

7.5. Extending Scope Beyond Biases

There is a number of cognitive phenomena affecting the interpretability of rules, which are not classified as cognitive biases. Remarkably, since 1960 there is a consistent line of work by psychologists studying cognitive processes related to rule induction, which is centred around the so-called *Wason's 2-4-6 problem* [161]. Cognitive science research on rule induction in humans has so far not been noticed in the rule learning subfield of machine learning.¹⁰ It was out of the scope of the objectives of this review to conduct an analysis of the significance

⁹This tendency is explained by Wolfe and Britt [169] as follows: "For people with this belief, facts and support are treated uncritically. . . . More importantly, arguments and information that may support another side are not part of the schema and are also ignored."

¹⁰Based on our analysis of cited reference search in Google Scholar for [161].

of these results for rule learning, nevertheless we believe that such investigation could bring interesting insights for cognitively-inspired design of rule learning algorithms.

Another promising direction for further work is research focused on the interpretation of negations ("not"). Experiments conducted by Jiang et al. [75] show that the mental processes involved in processing negations slow down reasoning. Negation can be also sometimes ignored or forgotten [30], as it decreases veracity of long-term correct remembrance of information.

Most rule learning algorithms are capable of generating rules containing negated literals. For example, a healthy company can be represented as status = not(bankrupt).

Our precautionary suggestion based on interpretation of results obtained in general studies performed in experimental psychology [30] and neurolinguistics [75] is that artificial learning systems should refrain, wherever feasible, from the use of negation in the discovered rules that are to be presented to the user. Due the adverse implications of the use of negation on cognitive load and remembrance, empirical research focused interpretability of negation in machine learning is urgently needed.

8. Conclusion

To our knowledge, cognitive biases have not yet been discussed in relation to the interpretability of machine learning results. We thus initiated this review of research published in cognitive science with the intent of providing a psychological basis to further research in inductive rule learning algorithms, and to the way their results are communicated. Our review covered twenty cognitive biases, heuristics, and effects that can give rise to systematic errors when inductively learned rules are interpreted.

For most biases and heuristics included in our review, psychologists have proposed "debiasing" measures. Application of prior empirical results obtained in cognitive science allowed us to propose several methods that could be effective in suppressing these cognitive phenomena when machine learning models are interpreted.

Overall, in our review, we processed only a fraction of potentially relevant psychological studies of cognitive biases, but we were unable to locate a single study focused on machine learning. Future research should thus focus on empirical evaluation of effects of cognitive biases in the machine learning domain.

Acknowledgments

TK was supported by long term institutional support of research activities. ŠB and TK were supported by grant IGA 33/2018 by Faculty of Informatics and Statistics, University of Economics, Prague. An initial version of this review was published as a part of TK's PhD thesis at Queen Mary University of London.

References

References

- Agrawal, R., Mannila, H., Srikant, R., Toivonen, H., Verkamo, A. I., 1995.
 Fast discovery of association rules. In: Fayyad, U. M., Piatetsky-Shapiro,
 G., Smyth, P., Uthurusamy, R. (Eds.), Advances in Knowledge Discovery
 and Data Mining. AAAI Press, pp. 307–328.
- [2] Albarracín, D., Mitchell, A. L., 2004. The role of defensive confidence in preference for proattitudinal information: How believing that one is strong can sometimes be a defensive weakness. Personality and Social Psychology Bulletin 30 (12), 1565–1584.
- [3] Alcala-Fdez, J., Alcala, R., Herrera, F., 2011. A fuzzy association rule-based classification model for high-dimensional problems with genetic rule selection and lateral tuning. IEEE Transactions on Fuzzy Systems 19 (5), 857–872.
- [4] Anderson, J., Fleming, D., 2016. Analytical procedures decision aids for generating explanations: Current state of theoretical development and implications of their use. Journal of Accounting and Taxation 8 (5), 51.
- [5] Andrews, R., Diederich, J., Tickle, A. B., 1995. Survey and critique of techniques for extracting rules from trained artificial neural networks. Knowledge-based systems 8 (6), 373–389.
- [6] Arkes, H. R., 1991. Costs and benefits of judgment errors: Implications for debiasing. Psychological Bulletin 110 (3), 486.
- [7] Arkes, H. R., Christensen, C., Lai, C., Blumer, C., 1987. Two methods of reducing overconfidence. Organizational Behavior and Human Decision Processes 39 (1), 133–144.
- [8] Azevedo, P. J., Jorge, A. M., 2007. Comparing rule measures for predictive association rules. In: Proceedings of the 18th European Conference on Machine Learning (ECML-07). Springer, Warsawa, Poland, pp. 510–517.
- [9] Bar-Hillel, M., 1979. The role of sample size in sample evaluation. Organizational Behavior and Human Performance 24 (2), 245–257.
- [10] Bar-Hillel, M., 1991. Commentary on Wolford, Taylor, and Beck: The conjunction fallacy? Memory & Cognition 19 (4), 412–414.
- [11] Bar-Hillel, M., Neter, E., 1993. How alike is it versus how likely is it: A disjunction fallacy in probability judgments. Journal of Personality and Social Psychology 65 (6), 1119.
- [12] Barberia, I., Blanco, F., Cubillas, C. P., Matute, H., 2013. Implementation and assessment of an intervention to debias adolescents against causal illusions. PLoS One 8 (8), e71303.

- [13] Baron, J., Beattie, J., Hershey, J. C., 1988. Heuristics and biases in diagnostic reasoning: II congruence, information, and certainty. Organizational Behavior and Human Decision Processes 42 (1), 88–110.
- [14] Beaman, C. P., McCloy, R., Smith, P. T., 2006. When does ignorance make us smart? Additional factors guiding heuristic inference. In: Proceedings of the Cognitive Science Society. Vol. 28.
- [15] Becker, E. S., Rinck, M., 2016. Reversing the mere exposure effect in spider fearfuls: Preliminary evidence of sensitization. Biological Psychology 121, 153–159.
- [16] Berka, P., 2018. Comprehensive concept description based on association rules: A meta-learning approach. Intelligent Data Analysis 22 (2), 325– 344.
- [17] Bibal, A., Frénay, B., 2016. Interpretability of machine learning models and representations: an introduction. In: Proceedings of the 24th European Symposium on Artificial Neural Networks (ESANN). pp. 77–82.
- [18] Boehm, L. E., 1994. The validity effect: A search for mediating variables. Personality and Social Psychology Bulletin 20 (3), 285–293.
- [19] Bond, S. D., Carlson, K. A., Meloy, M. G., Russo, J. E., Tanner, R. J., 2007. Information distortion in the evaluation of a single option. Organizational Behavior and Human Decision Processes 102 (2), 240–254.
- [20] Bornstein, R. F., 1989. Exposure and affect: overview and meta-analysis of research, 1968–1987. Psychological Bulletin 106 (2), 265.
- [21] Bruza, P. D., Wang, Z., Busemeyer, J. R., 2015. Quantum cognition: a new theoretical approach to psychology. Trends in Cognitive Sciences 19 (7), 383–393.
- [22] Camerer, C., Weber, M., 1992. Recent developments in modeling preferences: Uncertainty and ambiguity. Journal of Risk and Uncertainty 5 (4), 325–370.
- [23] Carlson, A., Betteridge, J., Kisiel, B., Settles, B., Hruschka Jr, E. R., Mitchell, T. M., 2010. Toward an architecture for never-ending language learning. In: Proceedings of the 24th AAAI Conference on Artificial Intelligence. Vol. 5. Atlanta, Atlanta, Georgia, p. 3.
- [24] Charness, G., Karni, E., Levin, D., 2010. On the conjunction fallacy in probability judgment: New experimental evidence regarding Linda. Games and Economic Behavior 68 (2), 551 556.
- [25] Clemen, R. T., Lichtendahl, K. C., 2002. Debiasing expert overconfidence: A Bayesian calibration model. In: Sixth International Conference on Probablistic Safety Assessment and Management (PSAM6).

- [26] Cochran, W. G., 2007. Sampling techniques. John Wiley & Sons.
- [27] Croskerry, P., Singhal, G., Mamede, S., 2013. Cognitive debiasing 2: impediments to and strategies for change. BMJ Quality & Safety, bmjqs—2012.
- [28] De Laat, P. B., 2017. Algorithmic decision-making based on machine learning from big data: Can transparency restore accountability? Philosophy & Technology, 1–17.
- [29] Dechêne, A., Stahl, C., Hansen, J., Wänke, M., 2010. The truth about the truth: A meta-analytic review of the truth effect. Personality and Social Psychology Review 14 (2), 238–257.
- [30] Deutsch, R., Kordts-Freudinger, R., Gawronski, B., Strack, F., 2009. Fast and fragile: A new look at the automaticity of negation processing. Experimental Psychology 56 (6), 434.
- [31] Díaz, C., Batanero, C., Contreras, J. M., 2010. Teaching independence and conditional probability. Boletín de Estadística e Investigación Operativa 26 (2), 149–162.
- [32] Ecker, U. K., Hogan, J. L., Lewandowsky, S., 2017. Reminders and repetition of misinformation: Helping or hindering its retraction? Journal of Applied Research in Memory and Cognition 6 (2), 185–192.
- [33] Edgell, S. E., Harbison, J., Neace, W. P., Nahinsky, I. D., Lajoie, A. S., 2004. What is learned from experience in a probabilistic environment? Journal of Behavioral Decision Making 17 (3), 213–229.
- [34] Ellsberg, D., 1961. Risk, ambiguity, and the Savage axioms. The Quarterly Journal of Economics 75 (4), 643–669.
- [35] Evans, J. S. B., 1989. Bias in Human Reasoning: Causes and Consequences. Lawrence Erlbaum Associates, Inc.
- [36] Evans, J. S. B., Stanovich, K. E., 2013. Dual-process theories of higher cognition: Advancing the debate. Perspectives on psychological science 8 (3), 223–241.
- [37] Evans, J. S. B., et al., 2007. Hypothetical thinking: Dual processes in reasoning and judgement. Vol. 3. Psychology Press.
- [38] Fantino, E., Kulik, J., Stolarz-Fantino, S., Wright, W., 1997. The conjunction fallacy: A test of averaging hypotheses. Psychonomic Bulletin & Review 4 (1), 96–101.
- [39] Fernbach, P. M., Darlow, A., Sloman, S. A., 2011. When good evidence goes bad: The weak evidence effect in judgment and decision-making. Cognition 119 (3), 459–467.

- [40] Fischoff, B., 1981. Debiasing. Tech. rep., Decision Research, Eugene, OR.
- [41] Fisk, J. E., 2002. Judgments under uncertainty: Representativeness or potential surprise? British Journal of Psychology 93 (4), 431–449.
- [42] Fiske, S. T., 1980. Attention and weight in person perception: The impact of negative and extreme behavior. Journal of Personality and Social Psychology 38 (6), 889.
- [43] Fitzsimons, G. J., Shiv, B., 2001. Nonconscious and contaminative effects of hypothetical questions on subsequent decision making. Journal of Consumer Research 28 (2), 224–238.
- [44] Fleischmann, M., Amirpur, M., Benlian, A., Hess, T., 2014. Cognitive biases in information systems research: A scientometric analysis. In: Proceedings of the 22st European Conference on Information Systems (ECIS 2014). Tel Aviv, Israel.
- [45] Fleisig, D., 2011. Adding information may increase overconfidence in accuracy of knowledge retrieval. Psychological Reports 108 (2), 379–392.
- [46] Freitas, A. A., 2014. Comprehensible classification models: a position paper. ACM SIGKDD Explorations 15 (1), 1–10.
- [47] Fürnkranz, J., 1997. Pruning algorithms for rule learning. Machine Learning 27 (2), 139–172.
- [48] Fürnkranz, J., 1999. Separate-and-conquer rule learning. Artificial Intelligence Review 13 (1), 3–54.
- [49] Fürnkranz, J., Flach, P. A., 2005. Roc nrule learningtowards a better understanding of covering algorithms. Machine Learning 58 (1), 39–77.
- [50] Fürnkranz, J., Gamberger, D., Lavrač, N., 2012. Foundations of Rule Learning. Springer-Verlag.
- [51] Fürnkranz, J., Kliegr, T., Paulheim, H., 2018. On cognitive preferences and the interpretability of rule-based models. CoRR abs/1803.01316. URL http://arxiv.org/abs/1803.01316
- [52] Fürnkranz, J., Kliegr, T., Paulheim, H., 2018. On cognitive preferences and the interpretability of rule-based models. arXiv preprint arXiv:1803.01316.
- [53] Gabriel, A., Paulheim, H., Janssen, F., 2014. Learning semantically coherent rules. In: Proceedings of the 1st International Workshop on Interactions between Data Mining and Natural Language Processing colocated with The European Conference on Machine Learning and Principles and Practice of Knowledge Discovery in Databases (DMNLP@PKDD/ECML). CEUR Workshop Proceedings, Nancy, France, pp. 49–63.

- [54] Gamberger, D., Lavrač, N., 2003. Active subgroup mining: A case study in coronary heart disease risk group detection. Artificial Intelligence in Medicine 28 (1), 27–57. URL http://dx.doi.org/10.1016/S0933-3657(03)00034-4
- [55] Ganter, B., Wille, R., 1999. Formal Concept Analysis Mathematical Foundations. Springer.
- [56] Geier, A. B., Rozin, P., Doros, G., 2006. Unit bias a new heuristic that helps explain the effect of portion size on food intake. Psychological Science 17 (6), 521–525.
- [57] Gettys, C. F., Fisher, S. D., Mehle, T., 1978. Hypothesis generation and plausibility assessment. Tech. rep., Decision Processes Laboratory, University of Oklahoma, Norman, annual report TR 15-10-78 (AD A060786.
- [58] Gettys, C. F., Mehle, T., Fisher, S., 1986. Plausibility assessments in hypothesis generation. Organizational Behavior and Human Decision Processes 37 (1), 14–33.
- [59] Gigerenzer, G., 2001. Content-blind norms, no norms, or good norms? A reply to Vranas. Cognition 81 (1), 93–103.
- [60] Gigerenzer, G., Goldstein, D. G., 1996. Reasoning the fast and frugal way: models of bounded rationality. Psychological Review 103 (4), 650.
- [61] Gigerenzer, G., Goldstein, D. G., 1999. Fast and frugal heuristics. In: Simple heuristics that make us smart. Oxford University Press, pp. 75–95.
- [62] Gigerenzer, G., Hoffrage, U., 1995. How to improve Bayesian reasoning without instruction: frequency formats. Psychological Review 102 (4), 684.
- [63] Gigerenzer, G., Hoffrage, U., 1999. Overcoming difficulties in Bayesian reasoning: A reply to Lewis and Keren (1999) and Mellers and McGraw (1999). Psychological Review (106), 425–430.
- [64] Goldstein, D. G., Gigerenzer, G., 1999. The recognition heuristic: How ignorance makes us smart. In: Simple heuristics that make us smart. Oxford University Press, pp. 37–58.
- [65] Griffin, D., Tversky, A., 1992. The weighing of evidence and the determinants of confidence. Cognitive Psychology 24 (3), 411–435.
- [66] Guidotti, R., Monreale, A., Ruggieri, S., Turini, F., Giannotti, F., Pedreschi, D., 2018. A survey of methods for explaining black box models. ACM computing surveys (CSUR) 51 (5), 93.

- [67] Hall, C. C., Ariss, L., Todorov, A., 2007. The illusion of knowledge: When more information reduces accuracy and increases confidence. Organizational Behavior and Human Decision Processes 103 (2), 277–290.
- [68] Haselton, M. G., Nettle, D., 2006. The paranoid optimist: An integrative evolutionary model of cognitive biases. Personality and Social Psychology Review 10 (1), 47–66.
- [69] Hasher, L., Goldstein, D., Toppino, T., 1977. Frequency and the conference of referential validity. Journal of Verbal Learning and Verbal Behavior 16 (1), 107–112.
- [70] Hertwig, R., Benz, B., Krauss, S., 2008. The conjunction fallacy and the many meanings of and. Cognition 108 (3), 740–753.
- [71] Hertwig, R., Gigerenzer, G., Hoffrage, U., 1997. The reiteration effect in hindsight bias. Psychological Review 104 (1), 194.
- [72] Hertwig, R., Herzog, S. M., Schooler, L. J., Reimer, T., 2008. Fluency heuristic: A model of how the mind exploits a by-product of information retrieval. Journal of Experimental Psychology: Learning, Memory, and Cognition 34 (5), 1191.
- [73] Hess, N. H., Hagen, E. H., 2006. Psychological adaptations for assessing gossip veracity. Human Nature 17 (3), 337–354.
- [74] Huber, M., 2010. From mindless to mindful decision making: Reflecting on prescriptive processes. Ph.D. thesis, University of Colorado at Boulder.
- [75] Jiang, Z.-q., Li, W.-h., Liu, Y., Luo, Y.-j., Luu, P., Tucker, D. M., 2014. When affective word valence meets linguistic polarity: Behavioral and erp evidence. Journal of Neurolinguistics 28, 19–30.
- [76] Juslin, P., Nilsson, H., Winman, A., 2009. Probability theory, not the very guide of life. Psychological Review 116 (4), 856.
- [77] Juslin, P., Nilsson, H., Winman, A., Lindskog, M., 2011. Reducing cognitive biases in probabilistic reasoning by the use of logarithm formats. Cognition 120 (2), 248–267.
- [78] Kachelmeier, S. J., Messier Jr, W. F., 1990. An investigation of the influence of a nonstatistical decision aid on auditor sample size decisions. Accounting Review, 209–226.
- [79] Kahneman, D., Tversky, A., 1972. Subjective probability: A judgment of representativeness. Cognitive Psychology 3 (3), 430–454.
- [80] Kahneman, D., Tversky, A., 1973. On the psychology of prediction. Psychological Review 80 (4), 237.
- [81] Keynes, J. M., 1922. A Treatise on Probability. Macmillan & Co.

- [82] Klayman, J., Ha, Y.-W., 1987. Confirmation, disconfirmation, and information in hypothesis testing. Psychological Review 94 (2), 211.
- [83] Kliegr, T., Svátek, V., Ralbovský, M., Šimůnek, M., Dec. 2011. SEWEBAR-CMS: Semantic analytical report authoring for data mining results. Journal of Intelligent Information Systems 37 (3), 371–395. URL http://dx.doi.org/10.1007/s10844-010-0137-0
- [84] Kononenko, I., 1993. Inductive and Bayesian learning in medical diagnosis. Applied Artificial Intelligence 7, 317–337.
- [85] Kunda, Z., 1999. Social Cognition: Making Sense of People. MIT press.
- [86] Lage, I., Chen, E., He, J., Narayanan, M., Kim, B., Gershman, S., Doshi-Velez, F., 2018. An evaluation of the human-interpretability of explanation. In: 32nd Conference on Neural Information Processing Systems (NIPS 2018). Montral, Canada.
- [87] Lakkaraju, H., Bach, S. H., Leskovec, J., 2016. Interpretable decision sets: A joint framework for description and prediction. In: Proceedings of the 22Nd ACM SIGKDD International Conference on Knowledge Discovery and Data Mining. KDD '16. ACM, New York, NY, USA, pp. 1675–1684.
- [88] Larrick, R. P., 2004. Debiasing. Blackwell Handbook of Judgment and Decision Making, 316–338.
- [89] Lau, A. Y., Coiera, E. W., 2009. Can cognitive biases during consumer health information searches be reduced to improve decision making? Journal of the American Medical Informatics Association 16 (1), 54–65.
- [90] Letham, B., Rudin, C., McCormick, T. H., Madigan, D., et al., 2015. Interpretable classifiers using rules and Bayesian analysis: Building a better stroke prediction model. The Annals of Applied Statistics 9 (3), 1350–1371.
- [91] Lewandowsky, S., Ecker, U. K., Seifert, C. M., Schwarz, N., Cook, J., 2012. Misinformation and its correction: Continued influence and successful debiasing. Psychological Science in the Public Interest 13 (3), 106–131.
- [92] Lilienfeld, S. O., Ammirati, R., Landfield, K., 2009. Giving debiasing away: Can psychological research on correcting cognitive errors promote human welfare? Perspectives on Psychological Science 4 (4), 390–398.
- [93] Liu, B., Hsu, W., Ma, Y., 1998. Integrating classification and association rule mining. In: Proceedings of the Fourth International Conference on Knowledge Discovery and Data Mining. KDD'98. AAAI Press, pp. 80–86.
- [94] Martens, D., Vanthienen, J., Verbeke, W., Baesens, B., 2011. Performance of classification models from a user perspective. Decision Support Systems 51 (4), 782–793.

- [95] Martire, K., Kemp, R., Sayle, M., Newell, B., 2014. On the interpretation of likelihood ratios in forensic science evidence: Presentation formats and the weak evidence effect. Forensic Science International 240, 61–68.
- [96] Martire, K. A., Kemp, R. I., Watkins, I., Sayle, M. A., Newell, B. R., 2013. The expression and interpretation of uncertain forensic science evidence: verbal equivalence, evidence strength, and the weak evidence effect. Law and Human Behavior 37 (3), 197.
- [97] Mellers, B., Hertwig, R., Kahneman, D., 2001. Do frequency representations eliminate conjunction effects? an exercise in adversarial collaboration. Psychological Science 12 (4), 269–275.
- [98] Michalkiewicz, M., Arden, K., Erdfelder, E., 2018. Do smarter people employ better decision strategies? the influence of intelligence on adaptive use of the recognition heuristic. Journal of Behavioral Decision Making 31 (1), 3–11.
- [99] Michalski, R. S., 1969. On the quasi-minimal solution of the general covering problem. In: Proceedings of the V International Symposium on Information Processing (FCIP 69)(Switching Circuits). Yugoslavia, Bled, pp. 125–128.
- [100] Michalski, R. S., 1983. A theory and methodology of inductive learning. In: Machine Learning. Springer, pp. 83–134.
- [101] Miller, T., 2019. Explanation in artificial intelligence: Insights from the social sciences. Artificial Intelligence 267, 1-38.
- [102] Mitchell, T. M., 1980. The Need for Biases in Learning Generalizations. Department of Computer Science, Laboratory for Computer Science Research, Rutgers Univ. New Jersey.
- [103] Monahan, J. L., Murphy, S. T., Zajonc, R. B., 2000. Subliminal mere exposure: Specific, general, and diffuse effects. Psychological Science 11 (6), 462–466.
- [104] Montoya, R. M., Horton, R. S., Vevea, J. L., Citkowicz, M., Lauber, E. A., 2017. A re-examination of the mere exposure effect: The influence of repeated exposure on recognition, familiarity, and liking. Psychological Bulletin 143 (5), 459.
- [105] Mumma, G. H., Wilson, S. B., 1995. Procedural debiasing of primacy/anchoring effects in clinical-like judgments. Journal of Clinical Psychology 51 (6), 841–853.
- [106] Narayanan, M., Chen, E., He, J., Kim, B., Gershman, S., Doshi-Velez, F., 2018. How do humans understand explanations from machine learning systems? an evaluation of the human-interpretability of explanation. arXiv preprint arXiv:1802.00682.

- [107] Nelson, J. D., 2005. Finding useful questions: On Bayesian diagnosticity, probability, impact, and information gain. Psychological Review 112 (4), 979.
- [108] Nelson, J. D., 2008. Towards a rational theory of human information acquisition. Oxford University Press Oxford, UK, pp. 143–163.
- [109] Nelson, J. D., McKenzie, C. R., Cottrell, G. W., Sejnowski, T. J., 2010. Experience matters: Information acquisition optimizes probability gain. Psychological Science 21 (7), 960–969.
- [110] Nickerson, R. S., 1998. Confirmation bias: A ubiquitous phenomenon in many guises. Review of General Psychology 2 (2), 175.
- [111] Nilsson, H., Winman, A., Juslin, P., Hansson, G., 2009. Linda is not a bearded lady: Configural weighting and adding as the cause of extension errors. Journal of Experimental Psychology: General 138 (4), 517.
- [112] Nisbett, R. E., 1993. Rules for Reasoning. Psychology Press.
- [113] Ohira, H., Winton, W. M., Oyama, M., 1998. Effects of stimulus valence on recognition memory and endogenous eyeblinks: Further evidence for positive-negative asymmetry. Personality and Social Psychology Bulletin 24 (9), 986–993.
- [114] Ordonez, C., Ezquerra, N., Santana, C. A., 2006. Constraining and summarizing association rules in medical data. Knowledge and Information Systems 9 (3), 1–2.
- [115] Pachur, T., Hertwig, R., 2006. On the psychology of the recognition heuristic: Retrieval primacy as a key determinant of its use. Journal of Experimental Psychology: Learning, Memory, and Cognition 32 (5), 983.
- [116] Pachur, T., Todd, P. M., Gigerenzer, G., Schooler, L., Goldstein, D. G., 2011. The recognition heuristic: A review of theory and tests. Frontiers in Psychology 2, 147.
- [117] Páez, A., 2019. The pragmatic turn in explainable artificial intelligence (xai). Minds and Machines, 1–19.
- [118] Parmley, M. C., 2006. The effects of the confirmation bias on diagnostic decision making. Ph.D. thesis, Drexel University.
- [119] Piltaver, R., Lustrek, M., Gams, M., Martincic-Ipsic, S., 2016. What makes classification trees comprehensible? Expert Systems with Applications 62, 333 346.
- [120] Pinker, S., 2015. Words and Rules: The Ingredients of Language. Basic Books.

- [121] Pohl, R., 2004. Cognitive Illusions: A Handbook on Fallacies and Biases in Thinking, Judgement and Memory. Psychology Press.
- [122] Pohl, R., 2017. Cognitive Illusions: A Handbook on Fallacies and Biases in Thinking, Judgement and Memory. Psychology Press, 2nd ed.
- [123] Pohl, R. F., Michalkiewicz, M., Erdfelder, E., Hilbig, B. E., 2017. Use of the recognition heuristic depends on the domains recognition validity, not on the recognition validity of selected sets of objects. Memory & Cognition 45 (5), 776–791.
- [124] Poursabzi-Sangdeh, F., Goldstein, D. G., Hofman, J. M., Vaughan, J. W., Wallach, H., 2018. Manipulating and measuring model interpretability. arXiv preprint arXiv:1802.07810.
- [125] Pratto, F., John, O. P., 2005. Automatic vigilance: The attention-grabbing power of negative social information. Social Cognition: Key Readings 250.
- [126] Rauch, J., May 2018. Expert deduction rules in data mining with association rules: a case study. Knowledge and Information Systems.
- [127] Ristoski, P., de Vries, G. K. D., Paulheim, H., 2016. A Collection of Benchmark Datasets for Systematic Evaluations of Machine Learning on the Semantic Web. Springer International Publishing, Cham, pp. 186–194. URL https://doi.org/10.1007/978-3-319-46547-0_20
- [128] Robinson-Riegler, G. L., Winton, W. M., 1996. The role of conscious recollection in recognition of affective material: Evidence for positive-negative asymmetry. The Journal of General Psychology 123 (2), 93–104.
- [129] Rozin, P., Royzman, E. B., 2001. Negativity bias, negativity dominance, and contagion. Personality and Social Psychology Review 5 (4), 296–320.
- [130] Schwarz, N., 2004. Metacognitive experiences in consumer judgment and decision making. Journal of Consumer Psychology 14 (4), 332–348.
- [131] Schwarz, N., Bless, H., Strack, F., Klumpp, G., Rittenauer-Schatka, H., Simons, A., 1991. Ease of retrieval as information: Another look at the availability heuristic. Journal of Personality and Social psychology 61 (2), 195.
- [132] Schwarz, N., Sanna, L. J., Skurnik, I., Yoon, C., 2007. Metacognitive experiences and the intricacies of setting people straight: Implications for debiasing and public information campaigns. Advances in Experimental Social Psychology 39, 127–161.
- [133] Serfas, S., 2011. Cognitive biases in the capital investment context. In: Cognitive Biases in the Capital Investment Context. Springer, pp. 95–189.

- [134] Shafir, E., 2013. The behavioral foundations of public policy. Princeton University Press.
- [135] Sides, A., Osherson, D., Bonini, N., Viale, R., 2002. On the reality of the conjunction fallacy. Memory & Cognition 30 (2), 191–198.
- [136] Simonson, I., Tversky, A., 1992. Choice in context: Tradeoff contrast and extremeness aversion. Journal of Marketing Research 29 (3), 281.
- [137] Skowronski, J. J., Carlston, D. E., 1989. Negativity and extremity biases in impression formation: A review of explanations. Psychological Bulletin 105 (1), 131.
- [138] Slowinski, R., Brzezinska, I., Greco, S., 2006. Application of Bayesian confirmation measures for mining rules from support-confidence Paretooptimal set. Proceedings of the 7th International Conference on Artificial Intelligence and Soft Computing (ICAISC 2006), 1018–1026.
- [139] Smith, E. E., Langston, C., Nisbett, R. E., 1992. The case for rules in reasoning. Cognitive Science 16 (1), 1–40.
- [140] Spengler, P. M., Strohmer, D. C., Dixon, D. N., Shivy, V. A., 1995. A scientist-practitioner model of psychological assessment: Implications for training, practice and research. The Counseling Psychologist 23 (3), 506– 534.
- [141] Stanovich, K. E., West, R. F., Toplak, M. E., 2013. Myside bias, rational thinking, and intelligence. Current Directions in Psychological Science 22 (4), 259–264.
- [142] Stecher, J., Janssen, F., Fürnkranz, J., 2016. Shorter rules are better, aren't they? In: Proceedings of the 19th International Conference on Discovery Science (DS-16). Bari, Italy, pp. 279–294.
- [143] Stolarz-Fantino, S., Fantino, E., Kulik, J., 1996. The conjunction fallacy: Differential incidence as a function of descriptive frames and educational context. Contemporary Educational Psychology 21 (2), 208–218.
- [144] Strossa, P., Černý, Z., Rauch, J., 2005. Reporting data mining results in a natural language. In: Foundations of Data Mining and Knowledge Discovery. Springer, pp. 347–361.
- [145] Taniguchi, H., Sato, H., Shirakawa, T., 2018. A machine learning model with human cognitive biases capable of learning from small and biased datasets. Scientific reports 8 (1), 7397.
- [146] Tentori, K., Crupi, V., 2012. On the conjunction fallacy and the meaning of and, yet again: A reply to Hertwig, Benz, and Krauss (2008). Cognition 122 (2), 123–134.

- [147] Tentori, K., Crupi, V., Russo, S., 2013. On the determinants of the conjunction fallacy: Probability versus inductive confirmation. Journal of Experimental Psychology: General 142 (1), 235.
- [148] Trope, Y., Gervey, B., Liberman, N., 1997. Wishful thinking from a pragmatic hypothesis-testing perspective. Lawrence Erlbaum Mahway, NJ, pp. 105–31.
- [149] Tversky, A., Kahneman, D., 1971. Belief in the law of small numbers. Psychological Bulletin 76 (2), 105.
- [150] Tversky, A., Kahneman, D., 1973. Availability: A heuristic for judging frequency and probability. Cognitive Psychology 5 (2), 207–232.
- [151] Tversky, A., Kahneman, D., 1974. Judgment under uncertainty: Heuristics and biases. Science 185 (4157), 1124–1131.
- [152] Tversky, A., Kahneman, D., 1981. The framing of decisions and the psychology of choice. Science 211 (4481), 453–458.
- [153] Tversky, A., Kahneman, D., 1983. Extensional versus intuitive reasoning: the conjunction fallacy in probability judgment. Psychological Review 90 (4), 293.
- [154] Tversky, A., Simonson, I., 1993. Context-dependent preference. Management Science 39 (10), 1179–1189.
- [155] Unkelbach, C., Rom, S. C., 2017. A referential theory of the repetitioninduced truth effect. Cognition 160, 110–126.
- [156] Vieider, F. M., 2009. The effect of accountability on loss aversion. Acta Psychologica 132 (1), 96–101.
- [157] Villejoubert, G., Mandel, D. R., 2002. The inverse fallacy: An account of deviations from Bayess theorem and the additivity principle. Memory & Cognition 30 (2), 171–178.
- [158] Vojíř, S., Zeman, V., Kuchař, J., Kliegr, T., 2018. Easyminer.eu: Web framework for interpretable machine learning based on rules and frequent itemsets. Knowledge-Based Systems 150, 111–115.
- [159] Škrabal, R., Šimůnek, M., Vojíř, S., Hazucha, A., Marek, T., Chudán, D., Kliegr, T., 2012. Association rule mining following the web search paradigm. In: Flach, P. A., Bie, T., Cristianini, N. (Eds.), Proceedings of the European Conference on Machine Learning and Knowledge Discovery in Databases (ECML-PKDD 2012). Springer Berlin Heidelberg, pp. 808–811.

- [160] Wang, T., Rudin, C., Doshi-Velez, F., Liu, Y., Klampfl, E., MacNeille, P., 2017. A bayesian framework for learning rule sets for interpretable classification. The Journal of Machine Learning Research 18 (1), 2357– 2393.
- [161] Wason, P. C., 1960. On the failure to eliminate hypotheses in a conceptual task. Quarterly Journal of Experimental Psychology 12 (3), 129–140.
- [162] Webb, G. I., 1994. Recent progress in learning decision lists by prepending inferred rules. In: Proceedings of the 2nd Singapore International Conference on Intelligent Systems. pp. B280–B285.
- [163] Werner, C., Hanea, A. M., Morales-Nápoles, O., 2018. Eliciting multivariate uncertainty from experts: Considerations and approaches along the expert judgement process. In: Elicitation. Springer, pp. 171–210.
- [164] Wilke, A., Mata, R., 2012. Cognitive bias. In: Ramachandran, V. (Ed.), Encyclopedia of Human Behavior (Second Edition), second edition Edition. Academic Press, San Diego, pp. 531 – 535.
- [165] Wille, R., 1982. Restructuring lattice theory: An approach based on hierarchies of concepts. In: Rival, I. (Ed.), Ordered Sets. Reidel, Dordrecht-Boston, pp. 445–470.
- [166] Willis, S., 2010. Standards for the formulation of evaluative forensic science expert opinion association of forensic science providers. Science & Justice 50 (1), 49.
- [167] Winkielman, P., Schwarz, N., Fazendeiro, T., Reber, R., et al., 2003. The hedonic marking of processing fluency: Implications for evaluative judgment. The psychology of evaluation: Affective processes in cognition and emotion, 189–217.
- [168] Winman, A., Juslin, P., Lindskog, M., Nilsson, H., Kerimi, N., 2014. The role of ANS acuity and numeracy for the calibration and the coherence of subjective probability judgments. Frontiers in Psychology 5, 851.
- [169] Wolfe, C. R., Britt, M. A., 2008. The locus of the myside bias in written argumentation. Thinking & Reasoning 14 (1), 1–27.
- [170] Zajonc, R. B., 1968. Attitudinal effects of mere exposure. Journal of Personality and Social Psychology 9 (2, Pt. 2), 1.
- [171] Zhang, C., Zhang, S., 2002. Association Rule Mining: Models and Algorithms. Springer-Verlag.
- [172] Zizzo, D. J., Stolarz-Fantino, S., Wen, J., Fantino, E., 2000. A violation of the monotonicity axiom: Experimental evidence on the conjunction fallacy. Journal of Economic Behavior & Organization 41 (3), 263–276.