

de Véricourt, Francis; Gurkan, Huseyin

Working Paper

Is your machine better than you? You may never know

ESMT Working Paper, No. 22-02 (R1)

Provided in Cooperation with:

ESMT European School of Management and Technology, Berlin

Suggested Citation: de Véricourt, Francis; Gurkan, Huseyin (2022) : Is your machine better than you? You may never know, ESMT Working Paper, No. 22-02 (R1), European School of Management and Technology (ESMT), Berlin,
<https://nbn-resolving.de/urn:nbn:de:101:1-2022121508081072470731>

This Version is available at:

<https://hdl.handle.net/10419/267687>

Standard-Nutzungsbedingungen:

Die Dokumente auf EconStor dürfen zu eigenen wissenschaftlichen Zwecken und zum Privatgebrauch gespeichert und kopiert werden.

Sie dürfen die Dokumente nicht für öffentliche oder kommerzielle Zwecke vervielfältigen, öffentlich ausstellen, öffentlich zugänglich machen, vertreiben oder anderweitig nutzen.

Sofern die Verfasser die Dokumente unter Open-Content-Lizenzen (insbesondere CC-Lizenzen) zur Verfügung gestellt haben sollten, gelten abweichend von diesen Nutzungsbedingungen die in der dort genannten Lizenz gewährten Nutzungsrechte.

Terms of use:

Documents in EconStor may be saved and copied for your personal and scholarly purposes.

You are not to copy documents for public or commercial purposes, to exhibit the documents publicly, to make them publicly available on the internet, or to distribute or otherwise use the documents in public.

If the documents have been made available under an Open Content Licence (especially Creative Commons Licences), you may exercise further usage rights as specified in the indicated licence.

Dec 8, 2022

ESMT Working Paper 22-02 (R1)

Is your machine better than you? You may never know.

Francis de Véricourt, ESMT European School of Management and Technology

Huseyin Gurkan, ESMT European School of Management and Technology

Copyright 2022 by ESMT European School of Management and Technology GmbH, Berlin, Germany, www.esmt.org.

All rights reserved. This document may be distributed for free – electronically or in print – in the same formats as it is available on the website of the ESMT (www.esmt.org) for non-commercial purposes. It is not allowed to produce any derivatives of it without the written permission of ESMT.

Find more ESMT working papers at [ESMT faculty publications](#), [SSRN](#), [RePEc](#), and [EconStor](#).

Is Your Machine Better Than You? You May Never Know.

Francis de Véricourt, Huseyin Gurkan*
ESMT Berlin, {francis.devericourt, huseyin.gurkan}@esmt.org

Artificial intelligence systems are increasingly demonstrating their capacity to make better predictions than human experts. Yet, recent studies suggest that professionals sometimes doubt the quality of these systems and overrule machine-based prescriptions. This paper explores the extent to which a decision maker (DM) supervising a machine to make high-stake decisions can properly assess whether the machine produces better recommendations. To that end, we study a set-up in which a machine performs repeated decision tasks (e.g., whether to perform a biopsy) under the DM’s supervision. Because stakes are high, the DM primarily focuses on making the best choice for the task at hand. Nonetheless, as the DM observes the correctness of the machine’s prescriptions across tasks, she updates her belief about the machine. However, the DM is subject to a so-called verification bias such that the DM verifies the machine’s correctness and updates her belief accordingly only if she ultimately decides to act on the task. In this set-up, we characterize the evolution of the DM’s belief and overruling decisions over time. We identify situations under which the DM hesitates forever whether the machine is better, i.e., she never fully ignores but regularly overrules it. Moreover, the DM sometimes wrongly believes with positive probability that the machine is better. We fully characterize the conditions under which these learning failures occur and explore how mistrusting the machine affects them. These findings provide a novel explanation for human-machine complementarity and suggest guidelines on the decision to fully adopt or reject a machine.

Key words: machine accuracy, decision making, human-in-the-loop, algorithm aversion, dynamic learning

1. Introduction

The adoption of machine learning (ML) algorithms is revolutionizing the delivery of products and services (McKendrick 2021), especially in domains that require human expertise, such as the medical and judiciary sectors. Indeed, artificial intelligence tools have demonstrated a capability to produce higher quality predictions than human judgment for many decision tasks (Grady 2019, Reardon 2019). The deployment of these tools in practice, however, has been limited (Wiens et al. 2019) and challenged by the tendency of decision makers to override—sometimes wrongly—algorithmic prescriptions. For instance, Sun et al. (2021) find in warehouse operations that employees significantly deviated from the recommendations of an algorithm. Lebovitz et al. (2022) also

*The authors are grateful to Santiago R. Balseiro, Denis Gromb, Jean Pauphilet and the seminar attendees at Yale University, Dartmouth College, HEC Paris, The Catholic University of Portugal, ML Approaches for Finance and Management conference at Humboldt University of Berlin, Bilkent University and the European Decision Science seminar for their valuable comments.

report how a team of radiologists in a large US-based hospital abandoned different ML algorithms after using them for several months.

This tendency to override algorithms is typically attributed to an intrinsic mistrust of machine-based predictions, often referred to as an *algorithm aversion* (Dietvorst et al. 2015, Gaube et al. 2021). This bias, however, may not be the sole reason for inappropriately and systematically overriding an algorithm. Indeed, the very context in which a human decision maker (DM) works can also prevent the DM from learning whether a machine produces better prescriptions.

In this paper, we explore the conditions under which making high-stake decisions hampers the DM’s ability to properly learn whether a machine is superior to human expertise. Importantly, we characterize the nature of the inappropriate overriding decisions that these learning failures give rise to, without relying on any mistrust bias. To that end, we analyze a set-up, in which a DM performs repeated decision tasks using the prescriptions of a machine. Each task consists in deciding whether or not to take a specific action. This corresponds, for instance, to deciding on a biopsy in a medical context. To make this choice, the machine produces a recommendation that the DM may overrule based on her own expertise. Crucially, the DM is uncertain about whether the machine makes better or worse decisions than she does, but as the DM verifies the correctness of the different machine’s prescriptions, she forms a belief about the machine’s true accuracy.

Our focus is thus on the DM’s learning behavior once the algorithm has been deployed, that is, after it has been properly trained and evaluated on representative data sets (see Kubat 2017) and possibly shown better-than-human accuracy levels. These datasets, however, never fully capture the ground truth, and the issue of empirical generalizability remains (see, e.g., Lebovitz et al. 2021). Hence, an expert may continue to observe and adjust her belief about the machine after adopting it. Yet, because the machine is deployed and makes prescriptions with real consequences, learning can be impaired in ways that do not exist during the training phase of the algorithm.

In particular, we consider situations in which the DM observes the correctness of the machine’s prediction and updates her belief accordingly only if the action is actually taken (e.g., when a biopsy is performed). In other words, the DM is subject to the so-called *verification bias* (see, e.g., Pepe 2003, p. 169), such that the accuracy parameters of a diagnostic test are learned only when a test result is verified by follow-up work (e.g., a biopsy reveals the presence or absence of the disease). This limitation can also stem from a form of salience bias or inattentional blindness (Taylor and Thompson 1982, Bordalo et al. 2012 and Tiefenbeck et al. 2018), which are especially prevalent for high-stakes decisions (see Lee et al. 2018). In this context, the DM focuses her limited attention on making the decisions at hand but is triggered to reassess the machine’s quality by the salient observation of a verified success or failure. (See Camacho et al. 2011 for an example of similar salient effects in the context of new drug prescriptions.)

Further, the DM only decides what is best for the task at hand and thus never acts for the purpose of verifying the machine’s accuracy. In this sense, the DM’s decisions are *exploration-free*. This restriction may be for legal or ethical concerns, which are often warranted when the stakes are high as in the medical and judiciary sectors (Bastani, Bayati and Khosravi 2021).

In this paper, we mainly examine the case where the machine and the DM are *substitutes* in that the DM’s accuracy is either better or worse than the machine’s. We focus on substitution for two reasons. First, our goal is to study inappropriate overriding decisions without relying on algorithm aversion, the studies of which assume substitution (Dietvorst et al. 2015, Sun et al. 2021). Second, and more importantly, we seek to determine if a complementarity between the DM and the machine might emerge from the DM’s inability to learn the nature of the machine. Assuming substitution enables us to disentangle this learning effect from an intrinsic complementarity between the DM and the machine. Nonetheless, we also explore situations where the machine and DM complement one another (see Section 8).

Our approach thus consists in analytically studying the evolution of the DM’s belief and overruling decisions over time. This enables identification of situations in which the DM properly learns whether the machine makes better predictions than she does. The asymptotic behavior of the DM’s belief further characterizes the different ways in which the DM fails to learn the true nature of the machine.

Following this approach, we find that the DM always properly learns whether the machine is better or worse in the absence of human-machine interactions, i.e., when the machine’s prescriptions never influence the DM’s decisions. Indeed, in this case, the DM verifies the machine independently of its prediction.¹ Hence, inappropriate overriding decisions may occur only if the machine has some influence on the DM’s choices.

When this influence is maximal, i.e., the machine’s prescriptions fully determine the DM’s choices, we find that the DM’s ability to learn depends on her prior about the task. Specifically, the DM properly learns that the machine is better (resp., worse), if the prior probability that an action is required for the task at hand is above (resp., below) a certain threshold. Hence, the DM can end up believing that the machine makes worse predictions than she does, even though the machine is actually better. This occurs when the action is not too frequently required for the tasks. Conversely, the DM learns that the machine is better even though it is actually worse when the action is frequently required.

In these two benchmarks, the DM’s belief about the machine has no effect on her choices. This contrasts with our main set-up, in which the DM’s decision to act, and hence her ability to learn,

¹ To be more precise, the verification event in the no-interaction benchmark is due to the DM’s own judgment, and the event that the machine prescribes to act is independent of the DM’s judgment conditional on the task’s type.

are endogenously determined by her current belief about the machine’s quality. Specifically, in this setting, the DM overrules the machine when the DM’s judgment contradicts the machine’s prescription and the DM sufficiently believes that the machine is worse.

When this is the case, we again find that the prior about the task determines when the DM fails to learn. However, the DM’s overriding decisions fundamentally change the nature of mislearning. Indeed, when the machine is actually better than the DM and the prior about the task is low, the DM’s belief always oscillates over time. In other words, the DM permanently remains unsure about whether the machine is better or not and constantly alternates between following and overriding its prescriptions. Further, and perhaps more interestingly, the DM sometimes treats the machine as if its prediction complements her own judgment, while in fact, the two are full substitutes. In contrast, when the machine is worse and the prior about the task is high, the belief converges to a Bernoulli random variable: the DM properly learns that the machine is worse with a given probability but wrongly learns that it is better with the remaining probability. In other words, the DM randomly ends up incorrectly believing that the machine is better.

Taken together, these results identify two different forms of mislearning—persistent hesitation and random inference—that can occur when a DM works with a machine to make high-stakes decisions. These findings also highlight the key role that the DM’s prior about the task plays in her ability to learn the true nature of the machine. Additionally they uncover a novel rationale—the uncertainty about the machine’s true performance—for why human experts may co-produce their decisions with a machine. These mislearning behaviors do not depend on the DM’s initial belief about the machine and thus hold even when the DM sufficiently believed in the machine’s performance to deploy it in the first place.

These results further suggest guidelines on the decision to fully adopt or reject a machine after it has been deployed. The question is whether the machine should make all the decisions henceforth or be abandoned for good at some point after working with it. Our results indicate that the longer the DM believes the machine is worse, the more likely she is correct in her assessment and hence should abandon the machine. The same is not true, however, if the DM increasingly believes that the machine is better. In this case, our findings suggest to rely on multiple DMs. If a consensus exists among the team that the machine is better, then the larger the team is, the more likely it is that the machine should be adopted. (See Section 6.)

Importantly, the mislearning behaviors we characterize in this paper do not stem from an intrinsic algorithm aversion but rather from certain contexts in which DMs make high-stakes decisions, as captured with the verification bias and the exploration-free condition. Yet, a DM who faces situations such as these may also be subject to mistrust biases against the machine, which can interact with our findings.

Indeed, mistrusting the machine affects the DM’s ability to learn in at least two ways. First, the DM may downplay the machine’s prescription when deciding to act (consistently with the decision-making literature, see, e.g., Soll and Mannes 2011), which alters the DM’s ability to observe the correctness of the machine’s predictions. Second, and in line with the algorithm aversion reported by Dietvorst et al. (2015), the DM’s belief in the machine may disproportionately drop upon observing a machine’s prediction error. We explore how these effects interact with our results (see Section 7) and find that our results hold in the former case but not in the later. When mistrust introduces a negativity bias in the DM’s learning process, the DM does not always properly learn that the machine is better if the prior is sufficiently high, as in the main set-up. Instead, the DM can wrongly learn with a positive probability that the machine is worse. In this sense, algorithm aversion sometimes interacts with our setting to randomize the DM’s ability to learn.

Finally, our results are robust to a partial relaxation of the verification bias, which is legitimate, for instance, when the bias stems from the DM’s limited attention. In this context, the DM also learns from unverified cases. Our results continue to hold as long as unverified cases are sufficiently less salient than verified ones, for which the true state of the world is revealed.

After reviewing the literature in Section 2, we present the model in Section 3. In Section 4, we analyze the no-interaction and no-overriding benchmarks and then focus on the main set-up in Section 5. We highlight the implications of our findings in Section 6 and study the effects of mistrust biases on our findings in Section 7. Further, we explore the settings where the machine and the DM complement one another in Section 8, and the verification bias is partially relaxed in Section 9. Finally, we discuss future research directions in the conclusion.

2. Literature Review

Our study is related to the recent and growing literature on the interaction between human decision-makers and data-driven algorithms. This research explores the extent to which co-production of decisions by a machine and a DM may improve performance. For instance, Boyaci et al. (2020) demonstrate in a rational inattention framework that human-machine interaction improves the overall accuracy of decisions, but sometimes at the cost of higher cognitive effort (see Boyaci et al. 2020 for additional references on formal models of machine-human interactions). Machine learning algorithms have also been proposed to provide interpretable cues to help decision makers improve their decisions (see Bastani, Bastani and Sinchaisri 2021, for instance). This stream of research further explores how to use human judgment to train or improve an algorithm (Van Donselaar et al. 2010, Ibrahim et al. 2021, Cowgill 2019).

We contribute to this literature by providing a novel rationale for why a DM may treat the machine’s prescriptions as a complement to her judgment. In fact, this stream of research typically

assumes that the machine’s accuracy is known and complements the DM’s judgement. In contrast, the DM and the machine are substitutes, and the machine’s accuracy is unknown in our setting.

In this sense, our study is closely related to the literature on overriding decisions and, more generally, trust in algorithmic prescriptions. In particular, Lebovitz et al. (2021) document over several months, how a team of radiologists lost trust in the quality of a machine learning algorithm that helped analyze medical images. Dietvorst et al. (2015) also found in an experimental set-up that their participants overrode a machine’s prescriptions, even after seeing that the machine’s algorithm performed better than the human did on average. This tendency to wrongly override machine-based prescriptions is further supported by empirical evidence in the field. For instance, Sun et al. (2021) observed that packing workers at the warehouses of the Alibaba Group regularly deviated from algorithmic prescriptions, which reduced operational efficiency. Several approaches have been explored to reduce deviations such as these, either with field experiments (Sun et al. 2021) or in the lab (Dietvorst et al. 2018).

In contrast to this stream of papers, our study proposes an alternative explanation for inappropriately overriding decisions such as these, which mostly stems from the context in which the decisions are made. Specifically, we trace these errors to four fundamentals (exploration-free, verification bias, informativeness and substitution), which capture some essential features of high-stakes decision making using machine-based predictions.

Recent studies also suggest that humans follow the principles of Bayesian inference when observing the correctness of machine-based decisions. For instance, Wang et al. (2018) and Guo et al. (2020) analyze in an experimental set-up how observers dynamically update their trust in the machine as they observe the failures and successes of its predictions (without overriding the machine, as in the benchmark of Section 4.2). These studies find that assuming Bayesian observers can explain the empirical level of human trust in the machine over time. The key difference with our set-up, however, is that the DM is not subject to verification bias and thus always observes the correctness of the machine’s prediction in their settings.

Verification bias is a form of selection bias that was first introduced by Ransohoff and Feinstein (1978) to describe situations where the accuracy of a diagnostic test is learned only with the verified cases, i.e., when follow-up actions are taken to confirm a test result. This bias towards verified cases permeates the medical field (e.g., Hujoel et al. 2021, Whiting et al. 2013, Petscavage et al. 2011, Bates et al. 1993, and Greenes and Begg 1985) and has been found in studies evaluating ML algorithms for medical applications (see, e.g., Tschandl et al. 2019). Most of this research has focused on developing estimators of accuracy based on the maximum likelihood to correct the bias.

The conditions required to avoid this bias in the frequentist literature, however, do not hold in our set-up.²

Learning with selective observations as in verification bias can also stem from a form of salience bias or inattentive blindness (Taylor and Thompson 1982, Bordalo et al. 2012 and Tiefenbeck et al. 2018). The behavioral science literature has studied biases such as these, (see, e.g., Kahneman 1973, Chapter 7, for the effects of focused attention on information filtering), which are due to the DM's limited cognitive capacity (Simon 1955). In this sense, our study also contributes to the growing stream of operations and economics literature, which deviates from standard Bayesian learning to account for limited cognitive capacity (see, e.g., Allon et al. 2021, Boyaci et al. 2020 and the references therein). In particular, our set-up is consistent with the notion of selective Bayesian updating (see the seminal work of Schwartzstein 2014), in that the DM only selects the actual success or failure of the machine to update her belief about the machine's type.

Finally, our work is related to the vast literature on learning problems, which have been extensively studied in management science and operations management. For instance, studies have considered price experimentation to learn demand curves by focusing on the tradeoffs between learning and earning, and design heuristic policies achieving good regret performance (Besbes and Zeevi 2009, Boyacı and Özer 2010, Cheung et al. 2017, Keskin and Birge 2019). In this stream of papers, the DM experiments (explores) with different prices in the beginning of the time horizon to earn (exploit) more in the remaining periods. Because of this ability to explore, the DM can, in principle, properly uncover the true demand curve in the limit. The objective of these papers is then to learn sufficiently fast so as to maximize profit. In contrast, we consider situations where exploring is not possible. Thus, the DM optimizes within each period and mislearning may emerge in our set-up.

In this sense, our approach resembles Harrison et al. (2012) which analyzes myopic pricing policies (see Section 4 in particular). In their set-up, demand functions are the focus of learning, whereas we consider unknown accuracy parameters. Therefore, the type of incomplete learning that may occur differs radically in each setting. In particular, incomplete learning takes the form of confounding beliefs in Harrison et al. (2012), such that the myopic policy charges an uninformative price, which prevents Bayesian updating from producing a different posterior. As a result, the DM becomes stuck in the same belief over time. In contrast, mislearning can take the form of belief oscillation in our set-up, which cannot occur in Harrison et al. (2012) per Proposition 2.

² For instance, our main set-up does not satisfy the missing-at-random assumption used by Begg and Greenes (1983) or the restrictions imposed on the data generating process proposed by Zhou (1993). In addition, our no-overriding benchmark corresponds to so-called extreme verification bias (Pepe 2003, p. 180), for which the estimation of accuracy parameters is impossible (Broemeling 2011).

Learning problems such as these are also extensively studied in economics (see for instance Smith and Sørensen 2000, Acemoglu et al. 2011, and references therein), with a particular focus on equilibrium learning dynamics shaped by multiple strategic agents. In this stream of research, Herrera and Hörner (2013) analyzes a set-up with short-lived myopic investors, in which only investing decisions are observable. Although this may resemble our set-up, their payoff, signal and learning structures differ, which yields a different type of mislearning. In particular, the belief converges to an interior point in their set-up (see Propositions 1 and 4 in Herrera and Hörner 2013), while it may not converge in ours.

3. Model Description

We consider a decision maker (DM) who faces a series of independent decision tasks over a discrete time infinite horizon. A machine further assists the DM by producing a recommendation about which decision to take for each task. The DM, however, does not know if the machine's accuracy is superior to her own judgment. As the accuracy of the machine's predictions is revealed over time, the DM forms a belief about whether or not she should override the machine's prediction. Next, we introduce the single decision task problem that the DM performs in each period. We then consider the whole time horizon.

3.1. Single Decision Task

A task consists in deciding whether or not a specific action (e.g., a biopsy) is required. We denote as $\Theta \in \{A, NA\}$ the type of task such that the action is required if $\Theta = A$ and is not required if $\Theta = NA$. The DM does not know the task's type but has a prior belief $p \triangleq \mathbb{P}(\Theta = A)$ that she should act.

To perform this task, the DM applies her expertise and elicits imperfect signal $S^H \in \{+, -\}$, such that $S^H = +$ (resp., $S^H = -$) indicates that $\Theta = A$ (resp., $\Theta = NA$). We denote the sensitivity (true positive rate) and specificity (true negative rate) of the signal by α^H and β^H , respectively. The DM is further assisted by a machine learning algorithm, which makes an independent prediction about type Θ . This prediction corresponds to a second signal, $S^M \in \{+, -\}$, with sensitivity and specificity equal to (α^M, β^M) .

Importantly, the DM is uncertain about whether the machine's accuracy is better than her own. Specifically, we denote the machine's type as $\Gamma \in \{B, W\}$. When $\Gamma = B$ (resp., $\Gamma = W$), signal S^M is *better* (resp., *worse*) than signal S^H , and the sensitivity and specificity of the signal are equal to (α^B, β^B) (resp., (α^W, β^W)). The machine is better (resp., worse) in the sense that the DM never (resp., always) overrules the machine when its type is perfectly known. This corresponds to the notion of substitution, which we introduce and formalize later in this section (see equations 4 and

5). To exclude degenerated cases, we further focus our analysis on situations where $\alpha^B > \alpha^W$ and $\beta^B > \beta^W$.³ This is only for the sake of clarity, as all of our results extend to the more general case.

Probability $b \triangleq \mathbb{P}(\Gamma = B)$ denotes then the DM's belief that the machine outperforms her ability to decide. In effect, these two types of machine induce two different probability measures $\mathbb{P}^B\{\cdot\}$ and $\mathbb{P}^W\{\cdot\}$ on the sample space of the machine's signals, such that $\mathbb{P}^\Gamma(S^M = +, \Theta = A) = \alpha^\Gamma p$ and $\mathbb{P}^\Gamma(S^M = -, \Theta = NA) = \beta^\Gamma \bar{p}$ for $\Gamma \in \{B, W\}$ (with $\bar{x} = 1 - x$ for $x \in [0, 1]$).

Based on realizations s^H and s^M of signals S^H and S^M , respectively, and her belief b about the machine, the DM updates her prior p that an action is required using Bayes' rule. The corresponding posterior probability is thus $\mathbb{P}(\Theta = A | S^H = s^H, S^M = s^M, b)$ (with a slight abuse of notation).⁴

The DM then decides to act if and only if the posterior is above a positive threshold r , i.e., $\mathbb{P}(\Theta = A | S^H = s^H, S^M = s^M, b) \geq r$; the DM does not act otherwise. This decision rule is optimal, for instance, when the DM seeks to maximize the expected value associated with correctly identifying the task's type. In this case, threshold r accounts for the false positive and false negative costs associated with the decision.⁵

Informativeness: In the following, we assume that the signals produced by both the DM and the machine are informative, in the sense that each signal provides sufficient information for the DM to decide. Formally, this corresponds to:

$$\mathbb{P}(\Theta = A | S^H = +) \geq r \text{ and } \mathbb{P}(\Theta = A | S^H = -) < r, \quad (1)$$

$$\mathbb{P}^B(\Theta = A | S^M = +) \geq r \text{ and } \mathbb{P}^B(\Theta = A | S^M = -) < r, \quad (2)$$

$$\mathbb{P}^W(\Theta = A | S^M = +) \geq r \text{ and } \mathbb{P}^W(\Theta = A | S^M = -) < r. \quad (3)$$

In other words, the sole realization of either S^H or S^M , whether the machine is of type B or W, fully determines whether or not the posterior is larger than threshold r , i.e., the DM takes the action. These conditions further imply that considering both signals S^H and S^M together is redundant when their realizations are aligned, i.e., when $s^H = s^M$. One signal is then sufficient for the DM to decide since the DM acts if $S^H = S^M = +$ and does not act if $S^H = S^M = -$. If the realizations are misaligned with $s^H \neq s^M$, however, the DM and the machine may override one another. In this case, we consider situations where the machine and the DM are full substitutes in the following sense.

³ This assumption guarantees that the DM's belief in a better machine decreases (resp., increases) upon observing an incorrect (resp., correct) machine prediction. In contrast, assuming $\alpha^W > \alpha^B$ (resp., $\beta^W > \beta^B$) implies that the DM's belief that the machine is better actually increases after observing a false negative (resp., false positive) error.

⁴ In particular, we have $\mathbb{P}(\Theta = A | S^H = s^H, S^M = s^M, 1) = \mathbb{P}^B(\Theta = A | S^H = s^H, S^M = s^M)$ and $\mathbb{P}(\Theta = A | S^H = s^H, S^M = s^M, 0) = \mathbb{P}^W(\Theta = A | S^H = s^H, S^M = s^M)$.

⁵ See, Alizamir et al. (2013) for instance, for a micro foundation of threshold r .

Substitution: We assume that a type B machine always overrides the DM’s judgment, while the DM always overrides the prescription of a type W machine. Formally, this corresponds to:

$$\mathbb{P}^B(\Theta = A | S^H = +, S^M = -) < r \text{ and } \mathbb{P}^B(\Theta = A | S^H = -, S^M = +) \geq r \quad (4)$$

$$\mathbb{P}^W(\Theta = A | S^H = +, S^M = -) \geq r \text{ and } \mathbb{P}^W(\Theta = A | S^H = -, S^M = +) < r \quad (5)$$

Thus, if the signals of the DM and a type B machine contradict one another, signal S^M alone determines whether or not the posterior probability is larger than the threshold (per equation (4)). Along with the Informativeness property, this means that a type B machine always determines whether the DM should act, independently of the DM’s judgment. In contrast, the DM decides alone and can ignore the prescription of a type W machine (per equation (5)). Hence, if the machine’s type is fully known, the DM and the machine never collaborate to make a decision. In this sense, the DM and the machine are substitutes for the task.

In essence, Informativeness and Substitution are conditions on the DM’s posterior probability about the task’s type, which in turn depends on the signals’ sensitivities and specificities, as well as prior p and threshold r .

3.2. Repeated Tasks and Learning

We now consider the situation where the DM faces a series of decision tasks over a discrete time infinite horizon. Task types Θ_t , $t \in \mathbb{N}$, are independent and identically distributed with probability p . (In the following, we use subscript t to denote the parameters associated with the task of period t .) At the beginning of period $t > 0$, the DM’s belief about the machine’s type is given by b_{t-1} , where b_0 is the prior belief at the beginning of the time horizon. The DM then obtains signals S_t^H , S_t^M and decides whether to act.

Exploration-Free: In making this choice, the DM considers only the task at hand. More formally, the DM acts if $\mathbb{P}(\Theta_t = A | S_t^H, S_t^M, b_{t-1}) \geq r$ and does nothing otherwise. In particular, the DM does not act for the sole purpose of uncovering the true task’s type and thus learning the machine’s. Instead, the DM decides what she thinks is best for the current task and is thus myopic with respect to learning the machine’s type.

At the end of the period, the DM updates her belief b_{t-1} to posterior b_t according to Bayes’ rule, if the DM observes type Θ_t .

Verification Bias: The DM, however, observes the task’s type and updates her belief accordingly if and only if an action is taken. Because decisions are exploration-free, the verification bias

implies that the DM updates her belief if and only if $\mathbb{P}(\Theta_t = \text{A} | S_t^{\text{H}}, S_t^{\text{M}}, b_{t-1}) \geq r$, in which case we assume that the DM follows Bayes' rule. Thus, we have

$$b_t = \begin{cases} b_{t-1} & \text{if } \mathbb{P}(\Theta_t = \text{A} | S_t^{\text{H}} = s^{\text{H}}, S_t^{\text{M}} = s^{\text{M}}, b_{t-1}) < r \\ \left[1 + \frac{\bar{b}_{t-1} \mathbb{P}^{\text{W}}(S_t^{\text{M}} = s^{\text{M}} | \Theta_t = \theta)}{b_{t-1} \mathbb{P}^{\text{B}}(S_t^{\text{M}} = s^{\text{M}} | \Theta_t = \theta)} \right]^{-1} & \text{if } \mathbb{P}(\Theta_t = \text{A} | S_t^{\text{H}} = s^{\text{H}}, S_t^{\text{M}} = s^{\text{M}}, b_{t-1}) \geq r, \end{cases} \quad (6)$$

where $\theta \in \{\text{A}, \text{NA}\}$ is the observed value of Θ_t .

Equation (6) highlights two mechanisms by which the DM's belief about the machine's type is endogenously determined over time. The first corresponds to the Bayesian updating of prior b_{t-1} when the DM observes type Θ_t . The second corresponds to the DM's ability to verify type Θ_t in the first place, that is, whether posterior belief $\mathbb{P}(\Theta_t = \text{A} | S_t^{\text{H}}, S_t^{\text{M}}, b_{t-1})$ is sufficiently large. This, in turn, depends on belief b_{t-1} . Equation (6) further implies that when the DM acts, she increases (resp., decreases) her belief that the machine is better if the machine's prescription turns out to be correct (resp., wrong).

This learning mechanism also corresponds to a selective Bayesian updating set-up (Schwartzstein 2014) in which a DM focuses her limited attention on making diagnostic decisions, instead of evaluating the machine that assists her. The salient observation of an actual success or failure of the machine, however, redirects the DM's attention to reassess her belief about the machine's type. This mechanism resembles the two-stage learning process of Allon et al. (2021), in which agents allocate their attention to different tasks (screening and belief updating in their setting) in each stage. (We relax this limitation on the DM's attention in Section 9.)

When the DM acts, the DM updates belief b_t in part based on signal S_t^{M} . The machine's type, however, determines the probability distribution, $\mathbb{P}^{\text{B}}\{\cdot\}$ or $\mathbb{P}^{\text{W}}\{\cdot\}$, of this signal. Hence, belief $(b_t)_{t \in \mathbb{N}}$ can follow two different stochastic processes depending on machine type Γ . The asymptotic behavior of belief b_t thus captures the DM's ability to learn whether the machine makes better predictions. Indeed, the DM properly learns the machine's type if her belief converges over time to 1 ($b_t \xrightarrow{\text{a.s.}} 1$) when the machine is better ($\Gamma = \text{B}$) and converges to 0 ($b_t \xrightarrow{\text{a.s.}} 0$) when the machine is worse ($\Gamma = \text{W}$). (Notation $\xrightarrow{\text{a.s.}}$ indicates almost-sure convergence.) In contrast, the DM mislearns the machine's type when $b_t \xrightarrow{\text{a.s.}} 0$ (resp., 1) and $\Gamma = \text{B}$ (resp., $\Gamma = \text{W}$). Learning may even be inconclusive when belief b_t converges to an interior point in $(0, 1)$ or oscillates. More formally, a stochastic process Y_t is said to be oscillating and recurrent if recurrent interval \mathcal{I} exists such that for any $\tau > 0$, $\mathbb{P}(Y_t \in \mathcal{I} \text{ for some } t > \tau | Y_\tau \in \mathcal{I}) = 1$ (see Definition 8.1 in Gut 2009 for instance).

Our objective, therefore, is to study the asymptotic behavior of b_t and characterize the resulting learning behavior of the DM.

4. Benchmarks

We first study two settings, in which the DM does not account for her belief about the machine when deciding to act. In the first *no-interaction* setting, the DM always ignores the machine's prescription and bases her choice solely on her own judgment S_t^H . In this sense, the DM and the machine do not interact when deciding on tasks. In the second *no-overriding* setting, the machine fully determines the DM's choice so that the DM never overrides its prediction S_t^M . Importantly, belief b_{t-1} does not determine whether an action is taken in both benchmarks and thus whether a machine's prediction is verified ex post. As a result, the second mechanism by which belief b_{t-1} influences posterior b_t in equation (6) is mute. This belief affects learning through only the first mechanism, i.e., the application of Bayes' rule when the DM acts.

More specifically, the DM acts if and only if $S_t^H = +$, regardless of the machine's signal S_t^M in the no-interaction benchmark, and if and only if $S_t^M = +$, regardless of her own judgment S_t^H in the no-overriding benchmark. The condition for acting, $\mathbb{P}(\Theta_t = A | S_t^H, S_t^M, b_{t-1}) \geq r$, thus reduces to $\mathbb{P}(\Theta_t = A | S_t^H) \geq r$ in the first benchmark and to $\mathbb{P}(\Theta_t = A | S_t^M, b_{t-1}) \geq r$ in the second one, which are, respectively, equivalent to $S_t^H = +$ and to $S_t^M = +$ for any b_{t-1} due to Informativeness (1)-(3). In both benchmarks, equation (6) then becomes

$$b_t = \begin{cases} b_{t-1} & \text{if } S_t^\sigma = - \\ \left[1 + \frac{\bar{b}_{t-1} \mathbb{P}^W(S_t^M | \Theta_t = \theta)}{b_{t-1} \mathbb{P}^B(S_t^M | \Theta_t = \theta)} \right]^{-1} & \text{if } S_t^\sigma = +. \end{cases} \quad (7)$$

where $\sigma = H$ and $\sigma = M$ in the no-interaction and no-overriding benchmark, respectively. In particular, the realization of S_t^σ is independent of belief b_{t-1} , which is in contrast to the condition in equation (6).

To study the asymptotic behavior of b_t , we consider instead the log-likelihood ratio L_t of the probability that $\Gamma = B$ in period t . Formally, L_t is a monotone continuous transformation of b_t given by $L_t \triangleq \log\left(\frac{b_t}{1-b_t}\right)$, such that

$$L_t = L_{t-1} + R_t,$$

where $(R_t)_{t \in \mathbb{N}}$ are i.i.d. random jumps. In particular, the log-likelihood ratio is increasing in the DM's belief, and the asymptotic behavior of L_t fully determines the asymptotic behavior of b_t . Indeed, we have $b_t \xrightarrow{\text{a.s.}} 1$ (and $b_t \xrightarrow{\text{a.s.}} 0$) if and only if $L_t \xrightarrow{\text{a.s.}} +\infty$ (and resp. $L_t \xrightarrow{\text{a.s.}} -\infty$) per the continuous mapping theorem.

Log-likelihood ratio L_t is a random walk governed by jumps $(R_t)_{t \in \mathbb{N}}$, which capture the magnitude and direction of the belief's updates. These random jumps take three possible values: a positive (resp., negative) value when the DM observes that the machine's prediction is correct (resp., wrong), i.e., $S_t^M = +$ and $\Theta_t = A$ (resp., $\Theta_t = NA$), or zero when the task's type is not observed, i.e.,

when $S_t^H = -$ in the no-interaction benchmark and $S_t^M = -$ in the no-overriding benchmark. The asymptotic behavior of L_t is then fully determined by the sign of the mean $\mathbb{E}^\Gamma[R_t]$. If $\mathbb{E}^\Gamma[R_t] > 0$ (resp., < 0), then log-likelihood ratio L_t increases in expectation and converges almost surely to $+\infty$ (resp., $-\infty$),⁶ while L_t does not converge when $\mathbb{E}^\Gamma[R_t] = 0$.

4.1. No-Interaction Benchmark

First, we characterize the DM’s ability to learn the machine’s type in the absence of DM-machine interactions, i.e., when the DM’s decisions always ignore the machine’s prescriptions. Specifically, when equation (7) holds with $\sigma = H$, we have

THEOREM 1 (Learning with No-Interaction). *When the machine is better $\Gamma = B$ (is worse $\Gamma = W$), $b_t \xrightarrow{a.s.} 1$ (and resp., $b_t \xrightarrow{a.s.} 0$).*

To prove this result, we first establish that $\mathbb{E}^\Gamma[R_t] > 0$ (resp., < 0) if $\Gamma = B$ (resp., $\Gamma = W$) and then apply the continuous mapping theorem. (All proofs are in the appendix.)

Theorem 1 states that the DM correctly learns the machine’s type when she decides to act solely based on her own judgment. In particular, verification bias does not prevent learning in the limit. This is because the DM’s sampling of the machine’s correct and wrong predictions is not biased in this case. Indeed, the DM acts and thus verifies the machine when $\{S_t^H = +\}$, regardless of the realization of S_t^M and hence, the probability to verify is independent of the machine’s prescription (conditional on the task’s type). This, in effect, relaxes the exploration-free condition by randomly enabling learning (with probability $\mathbb{P}(S_t^H = +)$), even when the machine’s prescription would have induced the DM not to act (i.e., when $\mathbb{P}(\Theta_t = A | S_t^H = +, S_t^M = -, b_{t-1}) < r$) in equation (6).

Overall, the theorem reveals that inappropriate overriding decisions may occur only when the machine biases the DM’s decisions since full learning occurs in the absence of DM-machine interactions. Next, we explore the case where the machine fully biases these decisions and hence the sampling of observations.

4.2. No-Overriding Benchmark

In our next result, we characterize the DM’s ability to learn the machine’s type when the DM’s choices are fully determined by the machine’s prescriptions. In this case, the bias of the machine on the DM’s decision is extreme. Specifically, we consider the DM’s asymptotic learning behavior when equation (7) holds with $\sigma = M$. We then have

THEOREM 2 (Learning with No-Overriding). *Unique thresholds p^B and p^W exist such that $p^B < p^W$ and,*

⁶ The divergence of L_t is due to the strong law large numbers; see Gut (Theorem 8.3 in p. 68 2009) for more details.

- when the machine is better ($\Gamma = B$), $b_t \xrightarrow{a.s.} 0$ if $p < p^B$, $b_t \xrightarrow{a.s.} 1$ if $p > p^B$; and b_t is recurrent and oscillates if $p = p^B$,
- when the machine is worse ($\Gamma = W$), $b_t \xrightarrow{a.s.} 0$ if $p < p^W$, $b_t \xrightarrow{a.s.} 1$ if $p > p^W$; and b_t is recurrent and oscillates if $p = p^W$.

Further, we have $p^B \triangleq \frac{\bar{\beta}^B \log\left(\frac{\bar{\beta}^W}{\bar{\beta}^B}\right)}{\bar{\beta}^B \log\left(\frac{\bar{\beta}^W}{\bar{\beta}^B}\right) + \alpha^B \log\left(\frac{\alpha^B}{\alpha^W}\right)}$ and $p^W \triangleq \frac{\bar{\beta}^W \log\left(\frac{\bar{\beta}^W}{\bar{\beta}^B}\right)}{\bar{\beta}^W \log\left(\frac{\bar{\beta}^W}{\bar{\beta}^B}\right) + \alpha^W \log\left(\frac{\alpha^B}{\alpha^W}\right)}$.

In essence, Theorem 2 states that the DM's ability to learn depends on her prior about the task as well as the machine's type. The DM learns that the machine is worse (resp., better) when her prior is below (resp., above) p^Γ for type $\Gamma \in \{B, W\}$. Importantly, this means that the DM may not properly learn whether the machine is better than her. Indeed, when prior p is low ($p < p^B$), the DM learns that the machine is worse ($b_t \xrightarrow{a.s.} 0$), while the machine is actually better ($\Gamma = B$). Similarly, when prior p is high ($p > p^W$), the DM learns that the machine is better ($b_t \xrightarrow{a.s.} 1$), while the machine is actually worse ($\Gamma = W$).

Theorem 2 stems from the fact that, in this benchmark, the DM acts and observes the machine's correctness only if the machine's signal is positive. In other words, the DM's observations are sampled solely from true and false positive predictions and never from true or false negative ones. Indeed, recall first that the DM's belief increases (resp., decreases) when the DM observes a correct (resp., incorrect) machine recommendation. Because the DM is able to observe this only when the machine's signal is positive, the DM increases her belief if and only if the machine correctly prescribes to act ($S_t^M = +$ and $\Theta_t = A$) and decreases her belief if and only if the machine wrongly prescribes to act ($S_t^M = +$ and $\Theta_t = NA$). But whether the task truly requires an action (i.e., $\Theta_t = A$) is determined by prior p . Hence, the DM increases her belief more frequently when the task is more likely to require an action, i.e., prior p takes higher values. The DM's belief will converge to one (resp., to zero) when prior p is sufficiently high (resp. low) such that the number of observed correct predictions is relatively higher than the number of incorrect ones.

This effect of prior p is absent from the no-interaction benchmark because the DM also observes the correctness of the machine's predictions when the machine recommends not to act (i.e., when $S_t^M = -$ and $\Theta_t = A$, or $S_t^M = -$ and $\Theta_t = NA$).

Note finally that magnitudes of these changes in beliefs do not depend on prior p but are determined by the accuracy parameters of the machine. Threshold p^Γ thus corresponds to the break-even value of prior p such that the expected increase in belief compensates for the expected decrease when the machine is of type Γ . When $p > p^\Gamma$, the expected increase dominates the expected decrease and the belief converges to one. When $p < p^\Gamma$, the opposite is true, and the belief converges to zero.

5. Main Set-up: Accounting for the DM's Belief about the Machine

In our main set-up, and in contrast to the previous benchmarks, the DM's belief about the machine influences the DM's decision to act, and hence her ability to verify the machine's predictions. As a result, learning is endogenously determined by the DM's current belief about the machine. This fundamentally changes the nature of mislearning.

More specifically, recall that due to Informativeness, the DM always decides according to her signal when it is consistent with the machine's signal with $S_t^H = S_t^M$. When these signals differ with $S_t^H \neq S_t^M$, the DM may override the machine when her current belief b_{t-1} that the machine is better is sufficiently low. Hence, belief b_{t-1} influences posterior b_t via the two mechanisms captured by equation (6). The following result determines when such overriding decisions occur. (The result follows from Substitution (4)-(5) and the continuity of the posterior probabilities in b_{t-1} ; see Appendix B.)

LEMMA 1. *Unique thresholds $b^- \in (0,1)$ and $b^+ \in (0,1)$ exist such that*

$$\mathbb{P}(\Theta_t = A | S_t^H = +, S_t^M = -, b_{t-1}) \geq r \Leftrightarrow b_{t-1} \leq b^-, \quad (8)$$

$$\mathbb{P}(\Theta_t = A | S_t^H = -, S_t^M = +, b_{t-1}) \leq r \Leftrightarrow b_{t-1} \leq b^+. \quad (9)$$

Lemma 1 states that when the DM's judgment contradicts the machine's prescription, i.e., $S_t^H \neq S_t^M$, the DM overrides the machine if and only if her belief in a better machine is sufficiently low, i.e., below a threshold. However, the DM can override the machine in two different ways, depending on whether the machine prescribes to act or not. This yields two different thresholds b^- and b^+ , which correspond to an overriding decision for a negative and positive machine signal, respectively.

These thresholds actually correspond to the value of belief b that makes the DM indifferent between acting and not acting when $S_t^H = -, S_t^M = +$ and $S_t^H = +, S_t^M = -$, respectively. Note also that the ranking between b^- and b^+ depends on the problem's parameters, and we define the minimum and maximum of these two thresholds as $b^H \triangleq \min(b^-, b^+)$ and $b^M \triangleq \max(b^-, b^+)$, respectively (where b^H and b^M can be equal).

Thus, when belief b_{t-1} is sufficiently large with $b_{t-1} > b^M$, the DM has sufficient confidence in the machine to always follow its prescriptions; this corresponds to the no-overriding benchmark. However, when the belief is sufficiently low with $b_{t-1} < b^H$, the DM always overrides the machine and decides solely based on her judgment, which corresponds to the no-interaction benchmark. Overall, these two cases are consistent with Substitution, which stipulates that the machine is either better or worse than the DM.

Interestingly, Lemma 1 further reveals that the DM may treat the machine's prescription as complementing—instead of substituting—her expertise. This occurs when the DM is sufficiently

unsure about the machine's type with $b_{t-1} \in (b^H, b^M)$. In this case, the DM and the machine complement one another in two possible ways, depending on whether threshold b^- is larger or smaller than threshold b^+ . If $b^+ < b^-$: the DM overrules the machine when its signal is negative but follows the machine's prescription when it is positive. In other words, the DM assumes that she makes better positive but worse negative decisions than the machine. In this sense, the DM and the machine collaborate on the task. As a result, the DM acts if and only if either the DM or the machine find evidence to do so ($S_t^H = +$ or $S_t^M = +$). If $b^- < b^+$, however, the DM overrules a positive machine's signal but follows a negative machine's signal and thus acts if and only if the DM and the machine agree that an action is required ($S_t^H = +$ and $S_t^M = +$).

Overall, Lemma 1 indicates that the DM's ability to learn the true type of task depends on her current belief about the machine. This means, in particular, that the random jumps of the corresponding log-likelihood ratio also depend on the current ratio. Formally, we have

$$L_t = L_{t-1} + R_t^{\text{HM}}(L_{t-1})$$

when the DM can override the machine. In contrast to the no-overriding benchmark, the random jumps R_t^{HM} are no longer i.i.d., as their distributions now depend on the magnitude of L_{t-1} . Thus, the sign of the expected jump, which determines the asymptotic behavior of belief b_t , is path-dependent. Next, we explore how this dependency affects the ability of the DM to learn the true machine type.

5.1. Learning When the Machine is Better

We first study the DM's ability to properly learn the machine's type when the machine is in fact better. Our next result characterizes the situations in which mislearning occurs in this case.

THEOREM 3. *When the machine is better, i.e. $\Gamma = \mathcal{B}$, if $p \leq p^{\mathcal{B}}$, then b_t oscillates and is recurrent; otherwise $b_t \xrightarrow{\text{a.s.}} 1$.*

Thus the DM's ability to learn the machine's type continues to depend on whether her prior about the task is sufficiently high. In fact, the threshold characterizing when proper learning occurs is the exact same as the one without overriding (defined in Theorem 2). Specifically, the DM learns that the machine is indeed better ($b_t \xrightarrow{\text{a.s.}} 1$) if and only if prior p is sufficiently high with $p > p^{\mathcal{B}}$. Figure 1b illustrates this case and exhibits asymptotic behavior.

The DM's ability to override the machine, however, fundamentally changes the nature of mislearning. Indeed, when prior p is such that $p \leq p^{\mathcal{B}}$, the DM wrongly learns that the machine is worse in the no-overriding benchmark. In the main set-up, however, the belief oscillates as illustrated in Figure 1a. Additionally, because the belief is also recurrent, the DM constantly switches among

overruling the machine ($b_t < b^H$), collaborating with it ($b_t \in (b^H, b^M)$) or letting the machine decide ($b_t > b^M$), as stated by the following corollary.

COROLLARY 1. *When the machine is better, i.e., $\Gamma = B$ and $p \leq p^B$, intervals $(0, b^H]$, (b^H, b^M) and $[b^M, 1)$ are recurrent for belief b_t .*

Hence, when the DM sufficiently believes that the machine is better, she never overrules it and we retrieve the dynamics of the no-overriding benchmark. That is, when $b_t > b^M$, learning is entirely driven by whether or not a machine's prescription to act is correct. Additionally, because prior p is low, the frequency of these correct predictions is also low, so the belief is decreasing in expectation.

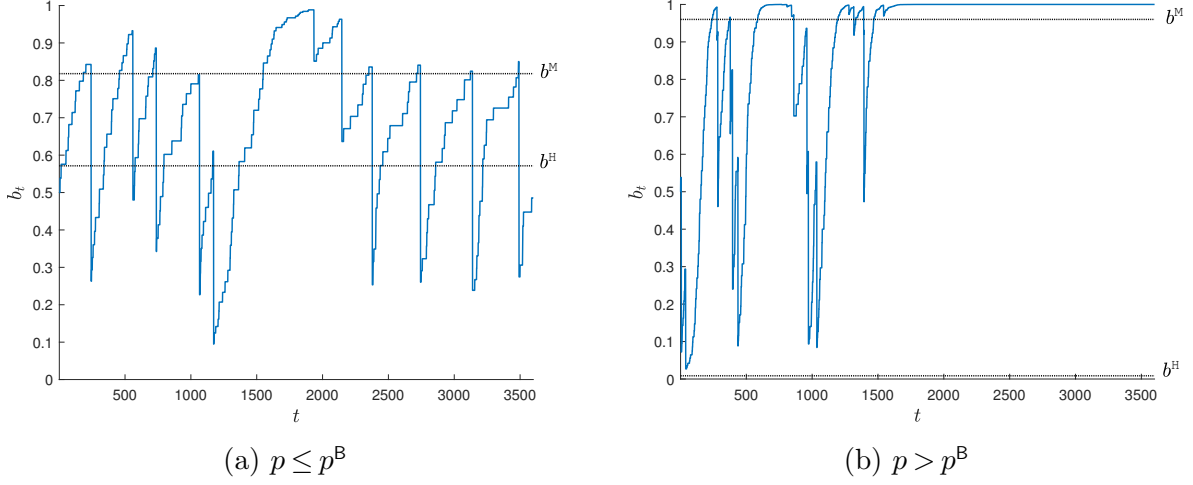
In contrast, when the DM sufficiently believes that the machine is worse with $b_t < b^H$, she always overrides the machine. In this case, the DM sometimes observes the machine's accuracy even when it prescribes not to act. This occurs when the DM's signal is positive and overrules a machine's negative signal. In this case, learning is driven by the true machine type, and because the machine is actually better, the belief increases in expectation.

Overall, the result holds because the DM's belief in the machine's type negatively reinforces her sampling bias: When the belief biases the sample of observations, the resulting updated belief tends to debias the sampling—and vice versa. As a result, belief b_t is pushed back downward when it reaches high values ($b_t > b^M$) and is pushed upward when it takes low values ($b_t < b^H$). Hence, the DM never fully learns that the machine is better, but due to overriding, never mislearns that it is worse either. In this sense, the DM always remains in perpetual uncertainty about whether or not to disregard the machine.

Interestingly, this long-run uncertainty induces the DM to sometimes treat the machine's prescription as a complement to her judgment. This happens when the belief reaches $b_t \in (b^H, b^M)$, which is a recurrent event. In this case, the DM and the machine co-produce the decision per Lemma 1 (and the explanations that follow). Because the machine and DM are actually substitutes, the emergence of this complementarity is driven only by the DM's inability to learn the true machine type.

5.2. Learning When the Machine is Worse

Per Theorems 2 and 3, the DM properly learns that the machine is good when prior p takes high values (i.e., $p > p^B$), whether the DM can override the machine or not. In this case, overriding essentially prevents the DM from wrongly learning that the machine is worse, which creates a perpetual state of uncertainty. In contrast, when the machine is indeed worse and the DM can overrule it, the DM may learn its true type for any prior p . This, however, occurs only randomly when prior p takes low values, as stated by the following result.

Figure 1 The DM's belief b_t when the machine is better, $\Gamma = B$ 

Note. $\alpha^H = \beta^H = 0.95$, $\alpha^B = \beta^B = 0.99$, $\alpha^W = \beta^W = 0.85$, $p^B = 0.15$ and $r = 0.07$, (a) $p = 0.05$, $b^H = 0.57$, $b^M = 0.81$, (b) $p = 0.2$, $b^H = 0.01$, $b^M = 0.96$.

THEOREM 4. *When the machine is worse, i.e., $\Gamma = W$, if $p \leq p^W$, then $b_t \xrightarrow{a.s.} 0$; otherwise, $b_t \xrightarrow{a.s.} X$ where X is a Bernoulli random variable.*

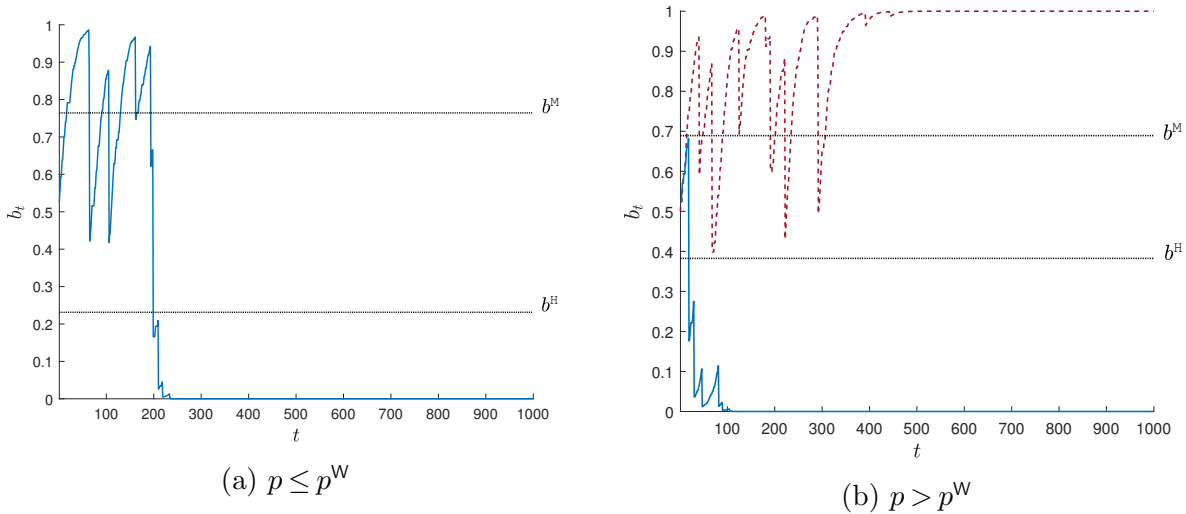
Theorem 4 indicates that the DM's ability to learn hinges again on prior p . As in the no-overriding benchmark, the DM can properly learn that the machine is worse ($b \xrightarrow{a.s.} 0$) if prior p takes low values ($p < p^W$, where threshold p^W is, again, the same as that in the no-overriding benchmark). Figure 2a illustrates this point, and depicts a sample path of b_t .

When the prior is high ($p > p^W$), however, the belief converges to a Bernoulli random variable. That is, the sample paths of belief b_t converge to zero with a certain probability and to one with the complement probability. In particular, the belief never oscillates nor converges to an interior point in the limit. Thus, the DM's ability to properly learn the machine's type is random in this case. In particular, the DM may sometimes wrongly learn that the machine is better, while it is actually worse. Figure 2b illustrates this point and depicts examples of the two possible sample paths for b_t , one (dashed line) converging to one and the other (solid line) to zero.

Similar to the better machine case, learning is driven by prior p as in the no-overriding benchmark when the belief is high ($b_t > b^M$), and by the true type of the machine when the belief is low ($b_t < b^H$). In the latter case, the belief decreases in expectation since the machine is worse.

Thus, for low prior $p < p^W$, the belief moves downward in expectation when it takes sufficiently high or low values and hence converges to zero in the long run. The DM then properly learns that the machine is worse.

For high prior $p > p^W$, however, the belief increases in expectation when the belief is already high and decreases when it is already low. In the long-run, the belief is thus pushed close to either

Figure 2 The DM's belief b_t when the machine is worse, $\Gamma = W$ 

Note. $\alpha^H = \beta^H = 0.95$, $\alpha^B = \beta^B = 0.99$, $\alpha^W = \beta^W = 0.9$, $p^W = 0.72$ and $r = 0.8$, (a) $p = 0.71$, $b^H = 0.23$, $b^M = 0.76$, (b) $p = 0.75$, $b^H = 0.38$, $b^M = 0.68$.

one or zero. This is because, in contrast to the case where the machine is better, the DM's belief positively reinforces the sampling bias: When the belief biases (resp., does not bias) the sample, the resulting updated belief tends to perpetuate the bias (resp., remain unbiased). Whether the belief reaches high or low values is then determined by realizations of the different signals and the task types and is thus random. Note that when the belief takes intermediary values ($b_t \in (b^H, b^M)$) it can either decrease or increase in expectation depending on the problem parameters. However, this region is transient since the belief is pushed away from the region when the belief is more extreme ($b_t \notin (b^H, b^M)$).

6. Implications

6.1. Learning and Mislearning

Taken together, these results provide theoretical limits on our ability to learn whether a machine makes better decisions than an expert. Interestingly, this inability to learn sometimes induces the DM to treat the machine's prescription as a complement to her own judgment, even though the two are actually substitutes. For instance, the DM may believe that her predictions have better sensitivity but worse specificity than those of the machine, while in fact, the machine is better in terms of both accuracy metrics. In this sense, the DM's uncertainty about the machine provides a novel rationale for why experts and machines may collaborate in practice.

Our results further identify the uncertainty surrounding the decision task as the key factor for mislearning. In fact, the DM fails to learn when she is most certain a priori about whether an action is required for a task (i.e., when prior p takes more extreme values with $p \leq p^B$ or $p > p^W$).

Conversely, the DM always properly learns the machine’s type when she is most uncertain about whether or not to act (i.e., prior p takes moderate values), as stated by the next corollary of Theorems 3 and 4.

COROLLARY 2. *The DM always correctly learns the type of the machine if and only if $p \in (p^B, p^W]$.*

6.2. Learning with Anticipation

In our set-up, as in the literature, the DM updates her belief using the past history of the observed accuracy of the machine’s predictions. Nonetheless, our results characterize the asymptotic behavior of this learning process and, as such, provide guidelines for DMs who anticipate the future behavior of their own belief. In particular, the nature of a learning failure is indicative of the machine’s type in our results. The DM may thus leverage this information to determine whether the machine is better.

Indeed, the DM’s belief may oscillate only when the machine is better (see Figure 1a), and always converges when it is worse (see Figure 2). Thus, the longer the DM remains uncertain (in the sense of Theorem 3), the more likely the machine is actually better. Similarly, the DM’s belief can converge to zero only if the machine is worse. Indeed, the belief either oscillates or converges to one when the machine is better. Hence, the longer the DM believes that the machine is worse, the more likely she is correct in her assessment.

Assessing if the DM is correct when she increasingly believes that the machine is better appears to be more challenging. Indeed, the DM’s belief can converge to one whether the machine is better (see Figure 1b) or worse (see Figure 2b). To circumvent this issue, one approach consists in relying on more than one decision maker. To see how, consider several identical decision makers who independently handle a series of tasks that are randomly drawn from the same sample and use the same machine. If this machine is better, all DMs should have the same limiting behavior: they either all remain uncertain or all learn that the machine is indeed better (per Theorem 3). However, if the machine is worse, the convergence to either zero or one is random (per Theorem 4). Thus, if a single DM in the team believes over time that the machine is worse, then the machine is indeed likely to be worse—even if the rest of the team believes it to be better. In contrast, if there is consensus in the team that the machine is better, then the larger the team is, the more likely it is that the machine is better.

In short, long-term uncertainty or a unanimous belief among large teams that the machine is better is indicative of a better machine. In contrast, persistently overruling the machine is indicative of a worse one.

6.3. Adoption or Rejection

Our study also sheds lights on the decision to fully adopt or reject the machine. Indeed, after observing and at times overriding the machine’s prescriptions, the DM’s belief may reach extreme levels. In these cases, the DM decides either to let the machine make all the decisions (as in Section 4.2), or to abandon the machine altogether, depending on whether the belief is sufficiently high or low, respectively. Once a machine is abandoned, however, the DM cannot learn about it anymore.

If the DM decides to fully adopt the machine—but continues to observe its performance—our results indicate that the DM will become increasingly confident about her adoption decision over time when prior p about the task is high ($p > p^B$ for a better machine, and $p > p^W$, for a worse one per Theorems 3-4). This occurs even when the machine is actually worse and should be abandoned.

In contrast, when the prior about the task is low ($p < p^B$ or $p < p^W$, depending on the true machine type), the DM increasingly doubts her adoption decision over time. This is because the DM’s belief in a better machine decreases in expectation over time and always approaches 0 in the limit in this case (per Theorems 3-4). This happens even when the machine is actually better and should be adopted.

Recall, finally, that when the machine is better and the prior about the task is low, the DM’s belief in the main set-up oscillates and is recurrent in $(0,1)$ (per Theorem 3 and Corollary 1). Therefore, the DM’s belief eventually reaches any low level with probability one. In other words, the DM always ends up abandoning the machine in the long run, even though the machine is actually better.

7. Mistrust Biases against the Machine

A key feature of the mislearning behavior in Theorems 3-4 is that they do not stem from an inherent mistrust against the machine. Instead, they stem from four fundamentals (verification bias, exploration-free decisions, informativeness and substitution), which characterize the set-up in which the DM works with the machine. Nonetheless, the DM may also be subject to mistrust biases against the machine in situations where these fundamentals are at play. In this section, we explore to which extent biases such as these interact with our four fundamentals to affect the DM’s learning behavior.

In our main set-up, mistrusting the machine affects the DM’s ability to learn in at least two ways. First, the DM may downplay the machine’s prescription when deciding to act, which alters the DM’s ability to observe the correctness of the machine’s predictions. Second, and in line with the algorithm aversion reported by Dietvorst et al. (2015), the DM’s belief in the machine may disproportionately drop upon observing the failure of a machine’s prediction. In the following, we inspect these different biases in turn.

7.1. Mistrusting the Machine When Deciding

The DM's mistrust in the machine may affect the way she weights the machine's prescription when deciding to act. This is consistent with the decision-making literature, which finds that individuals tend to discount information coming from external sources and overweight their own opinions (see for instance, Soll and Mannes 2011). To account for this possibility, we follow Stone (1961), who proposes a non-Bayesian approach to represent the aggregation of different opinions. This approach is commonly used to model mistrust bias (Özer and Zheng 2018), in particular when a human makes decisions based on the input of a data-driven algorithm (Ahsen et al. 2019, Boyaci et al. 2020). Specifically, the DM's updated belief that an action is required given the DM's and machine's signals is a linear combination between the updated belief of the human and that of the machine. More formally, the DM's posterior belief about the task is defined as

$$\tilde{\mathbb{P}}_\lambda(s^H, s^M, b_{t-1}) \triangleq \lambda \mathbb{P}(\Theta_t = A | S^H = s^H) + (1 - \lambda) \mathbb{P}(\Theta_t = A | S^M = s^M, b_{t-1}) \quad (10)$$

where $\lambda \in (0, 1)$ represents the DM's mistrust bias against the machine's signal. The higher the value of λ is, the more the DM mistrusts the machine. Belief $\tilde{\mathbb{P}}_\lambda(S^H, S^M, b_{t-1})$ corresponds then to the posterior probability $\mathbb{P}(\Theta = A | S^H, S^M, b_{t-1})$ derived from Bayes' rule in the main set-up. In particular, the DM decides to act if and only if $\tilde{\mathbb{P}}_\lambda(s^H, s^M, b_{t-1}) \geq r$.

In this set-up, the DM always overrules a better machine (never overrules a worse machine) if the bias is too high (resp., too low). In these extreme situations, the DM and the machine no longer substitute one another. We thus restrict our analysis to moderate values of mistrust parameter λ , as formalized by the following lemma, where $\tilde{\mathbb{P}}_\lambda^B(s^H, s^M) \triangleq \tilde{\mathbb{P}}_\lambda(s^H, s^M, 1)$ and $\tilde{\mathbb{P}}_\lambda^W(s^H, s^M) \triangleq \tilde{\mathbb{P}}_\lambda(s^H, s^M, 0)$ represent the DM's beliefs about the task's type when the machine's type is known.

LEMMA 2. *Two thresholds λ_{min} and λ_{max} exist such that if $\lambda \in (\lambda_{min}, \lambda_{max})$, then*

$$\tilde{\mathbb{P}}_\lambda^B(S_t^H = +, S_t^M = -) < r \text{ and } \tilde{\mathbb{P}}_\lambda^B(S_t^H = -, S_t^M = +) \geq r, \quad (11)$$

$$\tilde{\mathbb{P}}_\lambda^W(S_t^H = +, S_t^M = -) \geq r \text{ and } \tilde{\mathbb{P}}_\lambda^W(S_t^H = -, S_t^M = +) < r. \quad (12)$$

In the following, we focus on $\lambda \in (\lambda_{min}, \lambda_{max})$ to ensure that Substitution persists in the presence of mistrust bias. The next theorem shows that under these conditions, the structure of our main results continue to hold.

THEOREM 5. *Assume $\lambda \in (\lambda_{min}, \lambda_{max})$.*

- *When $\Gamma = B$, if $p \leq p^B$, then b_t oscillates and is recurrent; otherwise, $b_t \xrightarrow{a.s.} 1$.*
- *When $\Gamma = W$, if $p \leq p^W$, then $b_t \xrightarrow{a.s.} 0$; otherwise, $b_t \xrightarrow{a.s.} X$ where X is a Bernoulli random variable.*

Thus, the DM's learning behavior characterized in Theorems 3-4 does not change overall if the DM is biased against the machine's prescription when making a decision. Note also that Corollary 1 continues to hold in this case but with thresholds b^- and b^+ depending on λ (see Lemma 5 in the appendix).

7.2. Mistrusting the Machine When Updating Belief

The DM's mistrust against the machine can also affect the way the DM updates her belief about the machine. Dietvorst et al. (2015), for instance, experimentally show that individuals are more likely to ignore algorithm-based predictions after observing these algorithms err. More generally, the observation of negative outcomes, such as a prediction failure, more strongly impact the formation of an individual impression than do positive ones—a phenomenon referred to as negativity bias in the literature (Baumeister et al. 2001).

To account for this bias, we follow the literature (see for instance Coutts 2019, Möbius et al. 2022) and allow updated belief b_t to drop significantly upon observing an incorrect machine prediction. Specifically, the DM updates her belief following Bayes' rule when the machine is correct but magnifies the decrease in belief when the machine is wrong. More formally, we have

$$b_t = \begin{cases} b_{t-1} & \text{if } \mathbb{P}(\Theta_t = A \mid S_t^H = s^H, S_t^M = s^M, b_{t-1}) < r \\ \left[1 + \frac{\bar{b}_{t-1}}{b_{t-1}} \left[\frac{\mathbb{P}^W(S_t^M = s^M \mid \Theta_t = \theta)}{\mathbb{P}^B(S_t^M = s^M \mid \Theta_t = \theta)} \right] \phi(s^M, \theta) \right]^{-1} & \text{if } \mathbb{P}(\Theta_t = A \mid S_t^H = s^H, S_t^M = s^M, b_{t-1}) \geq r, \end{cases}$$

where function $\phi(s^M, \theta) = \mu > 1$ if the machine is incorrect (i.e., for $s^M = +$ and $\theta = \text{NA}$ or $s^M = -$ and $\theta = A$) and $\phi(s^M, \theta) = 1$ otherwise. Because ratio $\mathbb{P}^W(S_t^M \mid \Theta_t) / \mathbb{P}^B(S_t^M \mid \Theta_t) > 1$ when the machine is incorrect, the higher the value of mistrust parameter μ is, the lower belief b_t becomes. In particular, the main set-up corresponds to $\mu = 1$, which coincides with Bayes' rule.

The next result characterizes the asymptotic behavior of the DM's belief in the presence of this negativity bias.

THEOREM 6 (Learning with Negativity Bias). *Unique thresholds μ^B, μ^W, μ^H exist such that*

- *when the machine is better ($\Gamma = B$),*
 - *if $\mu \geq \mu^B$ and $\mu > \mu^H$, then $b_t \xrightarrow{a.s.} 0$.*
 - *if $\mu^B > \mu > \mu^H$, then $b_t \xrightarrow{a.s.} X$ where X is a Bernoulli random variable.*
 - *if $\mu^H \geq \mu \geq \mu^B$, then b_t is recurrent and oscillates.*
 - *if $\mu^B > \mu$ and $\mu^H \geq \mu$, then $b_t \xrightarrow{a.s.} 1$.*
- *when the machine is worse ($\Gamma = W$),*
 - *if $\mu \geq \mu^W$, then $b_t \xrightarrow{a.s.} 0$.*
 - *if $\mu^W > \mu$, then $b_t \xrightarrow{a.s.} X$ where X is a Bernoulli random variable.*

Further, we have

$$\mu^B \triangleq \frac{p\alpha^B \log\left(\frac{\alpha^B}{\alpha^W}\right)}{\bar{p}\bar{\beta}^B \log\left(\frac{\bar{\beta}^W}{\bar{\beta}^B}\right)}, \mu^W \triangleq \frac{p\alpha^W \log\left(\frac{\alpha^B}{\alpha^W}\right)}{\bar{p}\bar{\beta}^W \log\left(\frac{\bar{\beta}^W}{\bar{\beta}^B}\right)} \text{ and } \mu^H \triangleq \frac{p\alpha^H \alpha^B \log\left(\frac{\alpha^B}{\alpha^W}\right) + \bar{p}\bar{\beta}^H \beta^B \log\left(\frac{\beta^B}{\beta^W}\right)}{p\alpha^H \bar{\alpha}^B \log\left(\frac{\alpha^W}{\alpha^B}\right) + \bar{p}\bar{\beta}^H \bar{\beta}^B \log\left(\frac{\bar{\beta}^W}{\bar{\beta}^B}\right)} > 1.$$

Note that thresholds μ^B and μ^W actually play the same role as p^B and p^W in Theorems 3 and 4, respectively. Indeed, thresholds p_μ^B and p_μ^W exist such that $\mu > \mu^\Gamma \Leftrightarrow p < p_\mu^\Gamma$, for $\Gamma = \{B, W\}$. In particular, the structure of the results when the machine is worse (see Theorem 4) does not change in the presence of negativity bias. In this case, mistrust parameter μ affects the learning only through the value of threshold μ^W and, hence, p_μ^W .

When the machine is better, however, the presence of mistrust changes the structure of the results. In particular, under moderate mistrust such that $\mu^B > \mu > \mu^H$, the DM learns that the machine is better only with some probability (second point in Theorem 4). This is in contrast to the main set-up without mistrust, in which the DM always learns that the machine is better if $p > p^B$. In fact, the DM's belief can converge to a Bernoulli random variable in our main set-up only when the machine is worse. If the mistrust in the machine is too strong with $\mu > \max(\mu^H, \mu^B)$ (first point of the theorem), however, the DM always wrongly learns that the machine is worse. Otherwise, the bias does not alter the learning behavior. In fact, Theorem 6 reduces to Theorems 3-4 when $\mu = 1$. In this case, we have $\mu^H > \mu = 1$ and the belief either converges to one or oscillates depending on whether $\mu^B \leq \mu = 1$ or not, which is equivalent to $p^B \geq p$.

Overall, mistrust in the form of a negativity bias interacts with our fundamentals in a meaningful way only when the level of mistrust is moderate and the machine is actually better. In this case, whether the DM learns the true nature of the machine becomes random—while the DM always properly learns that the machine is better when she is not biased against the machine.

8. Complementarity

Thus far, we have focused on settings in which the machine and the DM are substitutes. Nonetheless, our framework also applies to the case of complementarity, which can take different forms. In this section, we explore the learning behavior of a DM who uncovers how a machine may complement her own judgment.

In our set-up, only two possible ways actually exist by which the machine and the DM complement one another. Indeed, a machine that complements the DM is superior in only one of the two dimensions of a judgment—positive or negative signals—and inferior in the other. Thus, the first form of complementarity corresponds to a DM who always overrides the machine when her judgment indicates that an action is required but always follows the machine's prescription if she finds

that she should not act. The converse holds for the second form: the DM overrides the machine when her judgment indicates not to act, but always follows the machine otherwise.

We denote these two machine types as C^+ and C^- , respectively, and their sensitivity and specificity are α^Γ and β^Γ for $\Gamma \in \{C^+, C^-\}$. In this section, we study the DM's learning behavior in the same settings as our base model, except for the substitution (4)-(5), which we replace by the following complementarity conditions.

Complementarity:

$$\mathbb{P}^{C^+}(\Theta = A | S^H = +, S^M = -) < r \text{ and } \mathbb{P}^{C^+}(\Theta = A | S^H = -, S^M = +) < r \quad (13)$$

$$\mathbb{P}^{C^-}(\Theta = A | S^H = +, S^M = -) \geq r \text{ and } \mathbb{P}^{C^-}(\Theta = A | S^H = -, S^M = +) \geq r \quad (14)$$

where $\mathbb{P}^{C^+}\{\cdot\}$ and $\mathbb{P}^{C^-}\{\cdot\}$ denote the probability measures induced by the two types of the machine.

The DM does not know the machine's type a priori. However, she forms a belief over time about which type of complementarity she is facing. With a slight abuse of notation, we refer to b_t as the DM's prior belief that $\Gamma = C^+$. The next result then characterizes the DM's ability to learn how the machine complements her judgment.

THEOREM 7. *We have,*

- *when $\Gamma = C^+$, then $b_t \xrightarrow{a.s.} 1$,*
- *when $\Gamma = C^-$, a unique threshold p^C exists such that $b_t \xrightarrow{a.s.} 1$ if $p \leq p^C$; otherwise $b_t \xrightarrow{a.s.} X$ where X is a Bernoulli random variable. (Threshold p^C is defined in (95) in Appendix D.)*

Thus, the DM always properly learns the machine's type when the actual form of complementarity is C^+ . In contrast, the DM can mislearn how the machine complements her judgment when the true type is C^- and the prior about the task is high (i.e., $p > p^C$). In this case, learning is random and the DM wrongly learns with positive probability that the machine is of type C^+ .

This result is akin to Theorem 4 when the machine and DM are substitutes. In contrast to Theorem 3, however, the DM's belief never oscillates and always converges to either zero or one. Thus, in the limit, the DM is always certain of the form of complementarity that the machine provides. In particular, the DM never behaves as if the machine and DM were substitutes. Again, this is in contrast to our main set-up, in which the DM's decisions sometimes exhibit complementarity, while in fact, the DM and the machine are substitutes (see Corollary 1).

Note finally that conditions (13)-(14) correspond to a complementarity between the machine's and the DM's signals. Other forms of complementarity, however, exist. In particular, the DM may seek to uncover for which decision tasks the machine is better and for which ones the DM is. In the context of biopsies, for instance, this corresponds to understanding for what kinds of patients the

machine does better and for what kinds of patients it does worse. A possible approach to study this case is to consider our main set-up but with more than one type of decision task. Denote this type as \mathbb{T} , with threshold $r^{\mathbb{T}}$ and prior $p^{\mathbb{T}}$ for $\mathbb{T} \in \{\mathbb{T}_1, \mathbb{T}_2, \dots\}$. These different types may correspond to different kinds of patients, for instance. The DM then forms different beliefs $b_t^{\mathbb{T}}$ over time so that the findings of Section 5 independently apply to each task's type $\mathbb{T} \in \{\mathbb{T}_1, \mathbb{T}_2, \dots\}$. With multiple task types, these findings characterize when the DM wrongly learns the decision tasks for which the machine's predictions are superior, and the ones for which her own judgment is better.

9. Partial Relaxation of the Verification Bias

In this section, we explore the robustness of our results when the verification bias is partially relaxed, which is legitimate when the bias stems from the DM's limited attention. In this context, the DM also learns from unverified cases and updates her belief based on the machine's (and her own) signal when she does not act. Because of salience effects and inattentive blindness, however, the DM assigns relatively less weight to this unverified information compared to information based on a verified case, for which the true state is revealed.

Formally, we consider inattentive blindness parameter $\varepsilon \in [0, 1]$, such that the DM's belief b_{t-1} is updated to b_t , as follows

$$b_t = \begin{cases} \left[1 + \frac{\bar{b}_{t-1}}{b_{t-1}} \left(\frac{\mathbb{P}^W(S_t^M = s^M, S_t^H = s^H)}{\mathbb{P}^B(S_t^M = s^M, S_t^H = s^H)} \right)^\varepsilon \right]^{-1} & \text{if } \mathbb{P}(\Theta_t = A | S_t^H = s^H, S_t^M = s^M, b_{t-1}) < r \\ \left[1 + \frac{\bar{b}_{t-1}}{b_{t-1}} \frac{\mathbb{P}^W(S_t^M = s^M | \Theta_t = \theta)}{\mathbb{P}^B(S_t^M = s^M | \Theta_t = \theta)} \right]^{-1} & \text{if } \mathbb{P}(\Theta_t = A | S_t^H = s^H, S_t^M = s^M, b_{t-1}) \geq r. \end{cases} \quad (15)$$

Here, ε represents how less salient unverified information is compared to verified information.⁷ The higher the value of ε is, the more sensitive the DM is to the informativeness of the machine's signal for an unverified case compared to a verified one. The verification bias is fully relaxed and proper learning occurs when $\varepsilon = 1$ per Proposition 1 in Appendix E.⁸ By contrast, our main set-up corresponds to $\varepsilon = 0$. The next theorem shows that our results continue to hold when ε is positive but sufficiently low.

THEOREM 8. *Unique thresholds ε^B and ε^W exist such that*

- *when the machine is better ($\Gamma = B$), if $\varepsilon \leq \varepsilon^B$ and $p < p^B$, then b_t oscillates and is recurrent; otherwise, $b_t \xrightarrow{a.s.} 1$.*
- *when the machine is worse ($\Gamma = W$), if $\varepsilon < \varepsilon^W$ and $p > p^W$, then $b_t \xrightarrow{a.s.} X$ where X is a Bernoulli random variable; otherwise, $b_t \xrightarrow{a.s.} 0$.*

⁷ To see this, consider a set-up with two different absolute weights for the verified and unverified cases, say ω_v and ω_u , respectively. This set-up is equivalent to the one in Section 9 by taking $\varepsilon = \omega_u / \omega_v$.

⁸ This proposition is consistent with the frequentist consistency of Bayesian updating (see, e.g., Diaconis and Freedman 1986), which implies perfect learning when the verification bias is fully relaxed with $\varepsilon = 1$.

Thresholds ε^B and ε^W are, respectively, defined in (98) and (99) in Appendix E, and p^B and p^W are in Theorem 2.

Theorem 8 corresponds to Theorems 3-4 with the additional condition that ε is less than ε^Γ for $\Gamma \in \{B, W\}$, respectively. In particular, when the unverified cases are sufficiently less salient than the verified ones with $\varepsilon < \min(\varepsilon^B, \varepsilon^W)$, our main results always hold.

10. Conclusion

This paper proposes a framework in which a machine performs repeated decision tasks under the supervision of a DM. In this set-up, we fully characterize the evolution of the DM’s belief about the machine and overruling decisions over time. We find that mislearning can take two radically different forms: a constant change of mind (oscillation of the DM’s belief per Theorem 3) and a chance of being persuaded that the machine has the wrong accuracy levels (convergence of the belief to a Bernoulli variable per Theorem 4). This contrasts with the convergence of the DM’s belief to an interior point in $(0, 1)$, which is often found in the dynamic learning literature (see e.g., confounding beliefs in Harrison et al. 2012). This analysis also provides a novel explanation for the joint production of decisions by machines and experts and suggests several guidelines for adopting or abandoning a machine.

The different forms of mislearning we uncover in this paper stem from the interaction between the DM’s belief in the machine and her decision to act, which in turn determines her sampling of correct and incorrect machine predictions. The belief and resulting sampling bias interact in a negative feed-back loop when the machine is better, while the feed-back loop is positive when the machine is worse.

These learning failures do not arise from an intrinsic mistrust bias against machine-based predictions, such as algorithmic aversion. Rather, they stem from the problem of learning about a machine while actually using its predictions to make high-stake decisions. We capture the key features of this problem with four fundamentals: informativeness, substitution, verification bias and exploration-free decisions.

Of these four, the last two conditions are crucial for our findings. Indeed, the DM always properly learns the true nature of the machine when the verification bias is sufficiently relaxed (per Theorem 8). Similarly, our no-interaction benchmark corresponds to a partial relaxation of the exploration-free condition, which also induces proper learning (see Theorem 1). In contrast, we find that the DM sometimes randomly fails to learn the machine’s accuracy when its predictions complement the DM’s judgment (see Theorem 7). We further expect mislearning to occur even when some of the signals are not informative, although the problem can become degenerate in this case (when none of the signals are informative, for instance).

We also restrict our analysis to two possible machine types, mostly for simplicity but our framework can be extended to account for more, possibly continuous types. Our results should not change overall as long as the previous fundamentals hold. Indeed, the DM’s belief that the machine outperforms her expertise is what fundamentally matters when deciding to override the machine. This, in essence, divides the different possible machine types into two distinct partitions depending on whether or not the type is better than the DM. In this sense, we retrieve a setup with two—albeit more convoluted—machine types.

Even though we assume them away, a DM may nonetheless be subject to mistrust biases against the machine in our set-up. Our results indicate that these biases can interact with our results in a significant way. In particular, the presence of mistrust bias akin to algorithm aversion sometimes randomizes the DM’s ability to properly learn the true nature of the machine. These results also provide novel hypotheses that future experimental research can test.

We focus on mistrust biases in this paper, but our framework can potentially accommodate other types of biases such as overconfidence and loss aversion (Benjamin 2019). Further, our framework can potentially account for situations in which the DM does not perfectly know her own accuracy, or has a misspecified representation of the machine (Fudenberg et al. 2017). Alternatively, the machine may provide partial explanations for the machine’s prescription, which may help the DM to learn the true machine accuracy (see, e.g., Puranam and Tsetlin 2021, for a way to model explainability).

Note finally that our framework may also be applied to situations where an expert supervises another expert instead of a machine. Doing so, however, requires assuming that experts learn the level of expertise solely by observing the ex post accuracy of someone’s judgments. While this precise setting may exist, experts such as radiologists typically provide a rationale or causal explanation to justify their prescriptions. These explanations are also indicative of someone’s knowledge and expertise. In other words, a human expert can more directly, and a priori, assess the quality of someone’s judgment in a way that is difficult with an ML algorithm (see, e.g., Cukier et al. 2021 for more on the difference between machine-based predictions and human cognition). In this sense, our framework is better suited for and offers a fruitful approach to exploring the issue of learning whether human expertise should overrule machine-based prescriptions.

References

- Acemoglu, D., Dahleh, M. A., Lobel, I. and Ozdaglar, A. (2011), ‘Bayesian learning in social networks’, *The Review of Economic Studies* **78**(4), 1201–1236.
- Ahsen, M. E., Ayvaci, M. U. S. and Raghunathan, S. (2019), ‘When algorithmic predictions use human-generated data: A bias-aware classification algorithm for breast cancer diagnosis’, *Information Systems Research* **30**(1), 97–116.

- Alizamir, S., de Véricourt, F. and Sun, P. (2013), ‘Diagnostic accuracy under congestion’, *Management Sci.* **59**(1), 157–171.
- Allon, G., Drakopoulos, K. and Manshadi, V. (2021), ‘Information inundation on platforms and implications’, *Operations Research* .
- Bastani, H., Bastani, O. and Sinchaisri, W. P. (2021), ‘Learning best practices: Can machine learning improve human decision-making?’.
- Bastani, H., Bayati, M. and Khosravi, K. (2021), ‘Mostly exploration-free algorithms for contextual bandits’, *Management Sci.* **67**(3), 1329–1349.
- Bates, A. S., Margolis, P. A. and Evans, A. T. (1993), ‘Verification bias in pediatric studies evaluating diagnostic tests’, *The Journal of pediatrics* **122**(4), 585–590.
- Baumeister, R. F., Bratslavsky, E., Finkenauer, C. and Vohs, K. D. (2001), ‘Bad is stronger than good’, *Review of general psychology* **5**(4), 323–370.
- Begg, C. B. and Greenes, R. A. (1983), ‘Assessment of diagnostic tests when disease verification is subject to selection bias’, *Biometrics* pp. 207–215.
- Benjamin, D. J. (2019), ‘Errors in probabilistic reasoning and judgment biases’, *Handbook of Behavioral Economics: Applications and Foundations 1* **2**, 69–186.
- Besbes, O. and Zeevi, A. (2009), ‘Dynamic pricing without knowing the demand function: Risk bounds and near-optimal algorithms’, *Operations Research* **57**(6), 1407–1420.
- Bordalo, P., Gennaioli, N. and Shleifer, A. (2012), ‘Salience theory of choice under risk’, *The Quarterly journal of economics* **127**(3), 1243–1285.
- Boyacı, T., Canyakmaz, C. and de Véricourt, F. (2020), ‘Human and machine: The impact of machine input on decision-making under cognitive limitations’.
- Boyacı, T. and Özer, Ö. (2010), ‘Information acquisition for capacity planning via pricing and advance selling: When to stop and act?’ , *Operations Research* **58**(5), 1328–1349.
- Broemeling, L. D. (2011), ‘Bayesian estimation of combined accuracy for tests with verification bias’, *Diagnostics* **1**(1), 53–76.
- Camacho, N., Donkers, B. and Stremersch, S. (2011), ‘Predictably non-bayesian: Quantifying salience effects in physician learning about drug quality’, *Marketing Science* **30**(2), 305–320.
- Cheung, W. C., Simchi-Levi, D. and Wang, H. (2017), ‘Dynamic pricing and demand learning with limited price experimentation’, *Operations Research* **65**(6), 1722–1731.
- Chung, K. L. and Zhong, K. (2001), *A course in probability theory*, Academic press.
- Coutts, A. (2019), ‘Good news and bad news are still news: Experimental evidence on belief updating’, *Experimental Economics* **22**(2), 369–395.

- Cowgill, B. (2019), ‘Bias and productivity in humans and machines’, *Columbia Business School Research Paper Forthcoming*.
- Csiszár, I. (1975), ‘I-divergence geometry of probability distributions and minimization problems’, *The annals of probability* pp. 146–158.
- Cukier, K., Mayer-Schönberger, V. and de Véricourt, F. (2021), *Framers: Human advantage in an age of technology and turmoil*, Penguin.
- Diaconis, P. and Freedman, D. (1986), ‘On the consistency of bayes estimates’, *The Annals of Statistics* pp. 1–26.
- Dietvorst, B. J., Simmons, J. P. and Massey, C. (2015), ‘Algorithm aversion: People erroneously avoid algorithms after seeing them err.’, *Journal of Experimental Psychology: General* **144**(1), 114.
- Dietvorst, B. J., Simmons, J. P. and Massey, C. (2018), ‘Overcoming algorithm aversion: People will use imperfect algorithms if they can (even slightly) modify them’, *Management Sci.* **64**(3), 1155–1170.
- Fudenberg, D., Romanyuk, G. and Strack, P. (2017), ‘Active learning with a misspecified prior’, *Theoretical Economics* **12**(3), 1155–1189.
- Gaube, S., Suresh, H., Raue, M., Merritt, A., Berkowitz, S. J., Lerner, E., Coughlin, J. F., Gutttag, J. V., Colak, E. and Ghassemi, M. (2021), ‘Do as ai say: susceptibility in deployment of clinical decision-aids’, *NPJ digital medicine* **4**(1), 1–8.
- Grady, D. (2019), ‘Ai took a test to detect lung cancer. it got an a’, *The New York Times* **20**.
- Greenes, R. A. and Begg, C. B. (1985), ‘Assessment of diagnostic technologies. methodology for unbiased estimation from samples of selectively verified patients.’, *Investigative radiology* **20**(7), 751–756.
- Guo, Y., Zhang, C. and Yang, X. J. (2020), Modeling trust dynamics in human-robot teaming: A bayesian inference approach, in ‘Extended Abstracts of the 2020 CHI Conference on Human Factors in Computing Systems’, pp. 1–7.
- Gut, A. (2009), *Stopped random walks*, Springer.
- Harrison, J. M., Keskin, N. B. and Zeevi, A. (2012), ‘Bayesian dynamic pricing policies: Learning and earning under a binary prior distribution’, *Management Sci.* **58**(3), 570–586.
- Herrera, H. and Hörner, J. (2013), ‘Biased social learning’, *Games and Economic Behavior* **80**, 131–146.
- Hujoel, I. A., Jansson-Knodell, C. L., Hujoel, P. P., Hujoel, M. L., Murray, J. A., Rubio-Tapia, A. et al. (2021), ‘Estimating the impact of verification bias on celiac disease testing’, *Journal of clinical gastroenterology* **55**(4), 327.
- Ibrahim, R., Kim, S.-H. and Tong, J. (2021), ‘Eliciting human judgment for prediction algorithms’, *Management Sci.* **67**(4), 2314–2325.
- Kahneman, D. (1973), *Attention and effort*, Vol. 1063, Citeseer.

- Kemperman, J. (1974), ‘The oscillating random walk’, *Stochastic Processes and their applications* **2**(1), 1–29.
- Keskin, N. B. and Birge, J. R. (2019), ‘Dynamic selling mechanisms for product differentiation and learning’, *Operations research* **67**(4), 1069–1089.
- Kubat, M. (2017), *An introduction to machine learning*, Vol. 2, Springer.
- Lebovitz, S., Levina, N. and Lifshitz-Assaf, H. (2021), ‘Is ai ground truth really “true”? the dangers of training and evaluating ai tools based on experts’ know-what’, *Management Information Systems Quarterly* .
- Lebovitz, S., Lifshitz-Assaf, H. and Levina, N. (2022), ‘To engage or not to engage with ai for critical judgments: How professionals deal with opacity when using ai for medical diagnosis’, *Organization Science* .
- Lee, H. C. B., Ba, S., Li, X. and Stallaert, J. (2018), ‘Salience bias in crowdsourcing contests’, *Information Systems Research* **29**(2), 401–418.
- MacKay, D. J., Mac Kay, D. J. et al. (2003), *Information theory, inference and learning algorithms*, Cambridge university press.
- McKendrick, J. (2021), ‘Ai adoption skyrocketed over the last 18 months’, <https://hbr.org/2021/09/ai-adoption-skyrocketed-over-the-last-18-months> [Online; accessed 18-February-2022].
- Möbius, M. M., Niederle, M., Niehaus, P. and Rosenblat, T. S. (2022), ‘Managing self-confidence: Theory and experimental evidence’, *Management Science* p. Forthcoming.
- Özer, Ö. and Zheng, Y. (2018), ‘Trust and trustworthiness’, *The Handbook of Behavioral Operations* pp. 489–523.
- Pepe, M. S. (2003), *The statistical evaluation of medical tests for classification and prediction*, Oxford University Press, USA.
- Petscavage, J. M., Richardson, M. L. and Carr, R. B. (2011), ‘Verification bias: an underrecognized source of error in assessing the efficacy of medical imaging’, *Academic Radiology* **18**(3), 343–346.
- Puranam, P. and Tsetlin, I. (2021), ‘Explainability as an optimal stopping problem: Implications for human-ai interaction’.
- Ransohoff, D. F. and Feinstein, A. R. (1978), ‘Problems of spectrum and bias in evaluating the efficacy of diagnostic tests’, *New England Journal of Medicine* **299**(17), 926–930.
- Reardon, S. (2019), ‘Rise of robot radiologists’, *Nature* **576**(7787), S54–S54.
- Schwartzstein, J. (2014), ‘Selective attention and learning’, *Journal of the European Economic Association* **12**(6), 1423–1452.
- Simon, H. A. (1955), ‘A behavioral model of rational choice’, *The quarterly journal of economics* **69**(1), 99–118.
- Smith, L. and Sørensen, P. (2000), ‘Pathological outcomes of observational learning’, *Econometrica* **68**(2), 371–398.

- Soll, J. B. and Mannes, A. E. (2011), ‘Judgmental aggregation strategies depend on whether the self is involved’, *International Journal of Forecasting* **27**(1), 81–102.
- Stone, M. (1961), ‘The opinion pool’, *The Annals of Mathematical Statistics* pp. 1339–1342.
- Sun, J., Zhang, D. J., Hu, H. and Van Mieghem, J. A. (2021), ‘Predicting human discretion to adjust algorithmic prescription: A large-scale field experiment in warehouse operations’, *Management Sci.* **68**(2), 846–865.
- Taylor, S. E. and Thompson, S. C. (1982), ‘Stalking the elusive” vividness” effect.’, *Psychological review* **89**(2), 155.
- Tiefenbeck, V., Goette, L., Degen, K., Tasic, V., Fleisch, E., Lalive, R. and Staake, T. (2018), ‘Overcoming salience bias: How real-time feedback fosters resource conservation’, *Management science* **64**(3), 1458–1476.
- Tschandl, P., Codella, N., Akay, B. N., Argenziano, G., Braun, R. P., Cabo, H., Gutman, D., Halpern, A., Helba, B., Hofmann-Wellenhof, R. et al. (2019), ‘Comparison of the accuracy of human readers versus machine-learning algorithms for pigmented skin lesion classification: an open, web-based, international, diagnostic study’, *The lancet oncology* **20**(7), 938–947.
- Van Donselaar, K. H., Gaur, V., Van Woensel, T., Broekmeulen, R. A. and Fransoo, J. C. (2010), ‘Ordering behavior in retail stores and implications for automated replenishment’, *Management Sci.* **56**(5), 766–784.
- Vatutin, V. A. and Wachtel, V. (2009), ‘Local probabilities for random walks conditioned to stay positive’, *Probability Theory and Related Fields* **143**(1), 177–217.
- Wang, C., Zhang, C. and Yang, X. J. (2018), Automation reliability and trust: A bayesian inference approach, in ‘Proceedings of the Human Factors and Ergonomics Society Annual Meeting’, Vol. 62, SAGE Publications Sage CA: Los Angeles, CA, pp. 202–206.
- Whiting, P. F., Rutjes, A. W., Westwood, M. E., Mallett, S., Group, Q.-. S. et al. (2013), ‘A systematic review classifies sources of bias and variation in diagnostic test accuracy studies’, *Journal of clinical epidemiology* **66**(10), 1093–1104.
- Wiens, J., Saria, S., Sendak, M., Ghassemi, M., Liu, V. X., Doshi-Velez, F., Jung, K., Heller, K., Kale, D., Saeed, M. et al. (2019), ‘Do no harm: a roadmap for responsible machine learning for health care’, *Nature medicine* **25**(9), 1337–1340.
- Zhou, X.-h. (1993), ‘Maximum likelihood estimators of sensitivity and specificity corrected for verification bias’, *Communications in Statistics-Theory and Methods*, **22**(11), 3177–3198.

Appendix A: No-Interaction and No-Overriding Benchmarks

We prove Theorems 1-2 by first deriving a recursive expression (in terms of b_{t-1}) for belief b_t using Bayes' rule. Then, we focus on the log-likelihood ratio process L_t defined by $L_t = \log(b_t/(1-b_t))$. Observe that when $L_t \rightarrow \infty$ (and/or $L_t \rightarrow -\infty$) almost surely, then it immediately follows that $b_t \rightarrow 1$ (and resp., $b_t \rightarrow 0$) due to the continuous mapping theorem since L_t is a continuous monotone transformation of b_t .

Proof of Theorem 1. In the no-interaction benchmark, the DM acts only if $S_t^H = +$, so it follows that

$$b_t = \left[1 + \frac{\bar{b}_{t-1}}{b_{t-1}} \zeta \right]^{-1} \quad (16)$$

where

$$\zeta \triangleq \left\{ \left[\left(\frac{\bar{\alpha}^W}{\bar{\alpha}^B} \right)^{1_{\{S_t^M=-\}}} \left(\frac{\alpha^W}{\alpha^B} \right)^{1_{\{S_t^M=+\}}} \right]^{1_{\{\Theta_t=A\}}} \left[\left(\frac{\beta^W}{\bar{\beta}^B} \right)^{1_{\{S_t^M=-\}}} \left(\frac{\bar{\beta}^W}{\bar{\beta}^B} \right)^{1_{\{S_t^M=+\}}} \right]^{1_{\{\Theta_t=NA\}}} \right\}^{1_{\{S_t^H=+\}}} \quad (17)$$

Here, $1_{\{\cdot\}}$ is the indicator variable. Using the definition of L_t , we obtain $L_t = L_{t-1} + R_t$, where

$$\begin{aligned} R_t \triangleq & 1_{\{S_t^H=+, S_t^M=+, \Theta_t=A\}} \log \left(\frac{\alpha^B}{\alpha^W} \right) + 1_{\{S_t^H=+, S_t^M=-, \Theta_t=A\}} \log \left(\frac{\bar{\alpha}^B}{\bar{\alpha}^W} \right) \\ & + 1_{\{S_t^H=+, S_t^M=+, \Theta_t=NA\}} \log \left(\frac{\bar{\beta}^B}{\bar{\beta}^W} \right) + 1_{\{S_t^H=+, S_t^M=-, \Theta_t=NA\}} \log \left(\frac{\beta^B}{\beta^W} \right). \end{aligned} \quad (18)$$

Therefore, L_t is a random walk with i.i.d. random jumps R_t . When the machine's type is $\Gamma \in \{B, W\}$, the mean of the random jump is

$$\mathbb{E}^\Gamma [R_t] = p\alpha^H \left[\alpha^\Gamma \log \left(\frac{\alpha^B}{\alpha^W} \right) + \bar{\alpha}^\Gamma \log \left(\frac{\alpha^B}{\alpha^W} \right) \right] + \bar{p}\bar{\beta}^H \left[\bar{\beta}^\Gamma \log \left(\frac{\bar{\beta}^B}{\bar{\beta}^W} \right) + \beta^\Gamma \log \left(\frac{\beta^B}{\beta^W} \right) \right]. \quad (19)$$

The mean $\mathbb{E}^\Gamma [R_t]$ is positive (negative) when $\Gamma = B$ (and resp., $\Gamma = W$). This is because, the terms inside the square brackets are the Kullback-Leibler (KL) divergence (Csiszár 1975) when $\Gamma = B$, i.e., $D_{KL}(\mathbb{P}^B(\cdot | \Theta = A) || \mathbb{P}^W(\cdot | \Theta = A))$ inside the first square brackets and $D_{KL}(\mathbb{P}^B(\cdot | \Theta = NA) || \mathbb{P}^W(\cdot | \Theta = NA))$ inside the first square brackets. Similarly, $-D_{KL}(\mathbb{P}^W(\cdot | \Theta = A) || \mathbb{P}^B(\cdot | \Theta = A))$ is inside the first square brackets and $-D_{KL}(\mathbb{P}^W(\cdot | \Theta = NA) || \mathbb{P}^B(\cdot | \Theta = NA))$ is inside the second square brackets when $\Gamma = W$. Gibbs' inequality (see Section 2.6 in MacKay et al. 2003) implies that KL divergence is always positive when the distributions are not the same, which implies the result. Q.E.D.

Proof of Theorem 2. In the no-overriding benchmark, the DM acts only if $S_t^M = +$, so it follows that

$$b_t = \frac{b_{t-1} (\alpha^B)^{1_{\{S_t^M=+, \Theta_t=A\}}} (\bar{\beta}^B)^{1_{\{S_t^M=+, \Theta_t=NA\}}}}{b_{t-1} (\alpha^B)^{1_{\{S_t^M=+, \Theta_t=A\}}} (\bar{\beta}^B)^{1_{\{S_t^M=+, \Theta_t=NA\}}} + \bar{b}_{t-1} (\alpha^W)^{1_{\{S_t^M=+, \Theta_t=A\}}} (\bar{\beta}^W)^{1_{\{S_t^M=+, \Theta_t=NA\}}}}. \quad (20)$$

Using the definition of L_t , we obtain $L_t = L_{t-1} + R_t$, where

$$R_t \triangleq 1_{\{S_t^M=+, \Theta_t=A\}} \log \left(\frac{\alpha^B}{\alpha^W} \right) + 1_{\{S_t^M=+, \Theta_t=NA\}} \log \left(\frac{\bar{\beta}^B}{\bar{\beta}^W} \right). \quad (21)$$

Therefore, L_t is a random walk with i.i.d. random jumps R_t .

When the machine's type is $\Gamma \in \{B, W\}$, then the mean of the random jump is

$$\mathbb{E}^\Gamma[R_t] = p\alpha^\Gamma \log \left(\frac{\alpha^B}{\alpha^W} \right) + \bar{p}\bar{\beta}^\Gamma \log \left(\frac{\bar{\beta}^B}{\bar{\beta}^W} \right). \quad (22)$$

If $p < p^\Gamma$, it follows that the mean $\mathbb{E}^\Gamma[R_t]$ and, hence, the drift of the random walk L_t is negative so $L_t \rightarrow -\infty$ (see Gut 2009, Theorem 9.1). The reverse condition (with strict inequality) implies the divergence to ∞ . If the mean $\mathbb{E}^\Gamma[R_t]$ equals 0 ($p = p^\Gamma$), then L_t is a martingale; hence, b_t oscillates (see Theorem 8.3.4 in Chung and Zhong 2001).

Finally, p^B and p^W are such that $p^B < p^W$ because $\alpha^B/\bar{\beta}^B > \alpha^W/\bar{\beta}^W$, which is implied by Substitution (4)-(5). Q.E.D.

Appendix B: Main Set-up

Proof of Lemma 1. This lemma follows from the fact that posterior probabilities are continuous and monotone in b_{t-1} and the boundary values 1 and 0 are on different sides of r due to Substitution (4)-(5). In particular, the posterior probabilities are

$$\begin{aligned} \mathbb{P}(\Theta_t = A \mid S_t^H = +, S_t^M = -, b_{t-1}) &= \frac{\alpha^H(b_{t-1}\bar{\alpha}^B + \bar{b}_{t-1}\bar{\alpha}^W)p}{\alpha^H(b_{t-1}\bar{\alpha}^B + \bar{b}_{t-1}\bar{\alpha}^W)p + \bar{\beta}^H(b_{t-1}\beta^B + \bar{b}_{t-1}\beta^W)\bar{p}}, \\ \mathbb{P}(\Theta_t = A \mid S_t^H = -, S_t^M = +, b_{t-1}) &= \frac{\bar{\alpha}^H(b_{t-1}\alpha^B + \bar{b}_{t-1}\alpha^W)p}{\bar{\alpha}^H(b_{t-1}\alpha^B + \bar{b}_{t-1}\alpha^W)p + \beta^H(b_{t-1}\bar{\beta}^B + \bar{b}_{t-1}\bar{\beta}^W)\bar{p}}. \end{aligned}$$

By solving the following equations for b_{t-1} ,

$$\mathbb{P}(\Theta_t = A \mid S_t^H = +, S_t^M = -, b_{t-1}) = r \text{ and } \mathbb{P}(\Theta_t = A \mid S_t^H = -, S_t^M = +, b_{t-1}) = r \quad (23)$$

we obtain the thresholds

$$b^- = \frac{\left(\frac{\bar{r}p\alpha^H}{r\bar{p}\bar{\beta}^H}\right)\bar{\alpha}^W - \beta^W}{\left(\frac{\bar{r}p\alpha^H}{r\bar{p}\bar{\beta}^H}\right)[\alpha^B - \alpha^W] + \beta^B - \beta^W} \text{ and } b^+ = \frac{\bar{\beta}^W - \left(\frac{\bar{r}p\bar{\alpha}^H}{r\bar{p}\beta^H}\right)\alpha^W}{\left(\frac{\bar{r}p\bar{\alpha}^H}{r\bar{p}\beta^H}\right)[\alpha^B - \alpha^W] + \beta^B - \beta^W}. \quad (24)$$

Thus, the result follows. Q.E.D.

Before proving Theorems 3 and 4, we first provide two constructive lemmas and their proofs.

LEMMA 3. *In the main set-up, the log-likelihood ratio process L_t is as follows.*

Case 1: If $b^+ > b^-$, then we have $L_t = L_{t-1} + R_t^{HM}(L_{t-1})$, where

$$R_t^{HM}(L_{t-1}) \triangleq \begin{cases} 1_{\{S_t^M=+\}} \left[1_{\{\Theta_t=A\}} \log \left(\frac{\alpha^B}{\alpha^W} \right) + 1_{\{\Theta_t=NA\}} \log \left(\frac{\bar{\beta}^B}{\bar{\beta}^W} \right) \right] & \text{if } L_{t-1} \geq L_h \\ 1_{\{S_t^H=+, S_t^M=+\}} \left[1_{\{\Theta_t=A\}} \log \left(\frac{\alpha^B}{\alpha^W} \right) + 1_{\{\Theta_t=NA\}} \log \left(\frac{\bar{\beta}^B}{\bar{\beta}^W} \right) \right] & \text{if } L_h > L_{t-1} > L_\ell \\ 1_{\{S_t^H=+\}} [v_1 + v_2] & \text{if } L_\ell \geq L_{t-1} \end{cases}$$

with $L_h \triangleq \log\left(\frac{b^+}{b^-}\right)$, $L_\ell \triangleq \log\left(\frac{b^-}{b^+}\right)$ and

$$\begin{aligned} v_1 &\triangleq 1_{\{\Theta_t=A\}} \left(1_{\{S_t^M=-\}} \log\left(\frac{\bar{\alpha}^B}{\bar{\alpha}^W}\right) + 1_{\{S_t^M=+\}} \log\left(\frac{\alpha^B}{\alpha^W}\right) \right), \\ v_2 &\triangleq 1_{\{\Theta_t=NA\}} \left(1_{\{S_t^M=-\}} \log\left(\frac{\beta^B}{\beta^W}\right) + 1_{\{S_t^M=+\}} \log\left(\frac{\bar{\beta}^B}{\bar{\beta}^W}\right) \right). \end{aligned}$$

Case 2: If $b^+ \leq b^-$, then we have $L_t = L_{t-1} + R_t^{HM}(L_{t-1})$ where

$$R_t^{HM}(L_{t-1}) \triangleq \begin{cases} 1_{\{S_t^M=+\}} \left[1_{\{\Theta_t=A\}} \log\left(\frac{\alpha^B}{\alpha^W}\right) + 1_{\{\Theta_t=NA\}} \log\left(\frac{\bar{\beta}^B}{\bar{\beta}^W}\right) \right] & \text{if } L_{t-1} > L_h \\ 1_{\{\Theta_t=A\}} \nu_1 + 1_{\{\Theta_t=NA\}} \nu_2 & \text{if } L_h \geq L_{t-1} \geq L_\ell \\ 1_{\{S_t^H=+\}} [\nu_1 + \nu_2] & \text{if } L_\ell > L_{t-1} \end{cases} \quad (25)$$

with $L_h \triangleq \log\left(\frac{b^-}{b^+}\right)$, $L_\ell \triangleq \log\left(\frac{b^+}{b^-}\right)$ and

$$\begin{aligned} \nu_1 &\triangleq 1_{\{S_t^M=+\}} \log\left(\frac{\alpha^B}{\alpha^W}\right) + 1_{\{S_t^M=-\}} 1_{\{S_t^H=+\}} \log\left(\frac{\bar{\alpha}^B}{\bar{\alpha}^W}\right), \\ \nu_2 &\triangleq 1_{\{S_t^M=+\}} \log\left(\frac{\bar{\beta}^B}{\bar{\beta}^W}\right) + 1_{\{S_t^M=-\}} 1_{\{S_t^H=+\}} \log\left(\frac{\beta^B}{\beta^W}\right). \end{aligned}$$

Proof of Lemma 3. As in the proof of Theorem 2, we derive a recursive expression for b_t in terms of b_{t-1} using Bayes' rule but this time using the decision-making procedure characterized in Lemma 1.

Case 1. When $b^+ > b^-$, we have the following three regimes

- $b_{t-1} \geq b^+$:

$$\mathbb{P}(\Theta_t = A \mid S_t^H = -, S_t^M = +, b_{t-1}) \geq r \quad (26)$$

$$\mathbb{P}(\Theta_t = A \mid S_t^H = +, S_t^M = -, b_{t-1}) < r \quad (27)$$

In this case, $S_t^M = +$ is sufficient and necessary to act, which implies the machine overrides the human's signal $S_t^H = -$. Further, $S_t^H = +$ is also overruled by $S_t^M = -$.

- $b^+ > b_{t-1} > b^-$:

$$\mathbb{P}(\Theta_t = A \mid S_t^H = -, S_t^M = +, b_{t-1}) < r \quad (28)$$

$$\mathbb{P}(\Theta_t = A \mid S_t^H = +, S_t^M = -, b_{t-1}) < r \quad (29)$$

In this case, $S_t^M = S_t^H = +$ is the only condition for acting.

- $b^- \geq b_{t-1}$:

$$\mathbb{P}(\Theta_t = A \mid S_t^H = -, S_t^M = +, b_{t-1}) < r \quad (30)$$

$$\mathbb{P}(\Theta_t = A \mid S_t^H = +, S_t^M = -, b_{t-1}) \geq r \quad (31)$$

In this case, $S_t^H = +$ is sufficient and necessary to act. The signal of the human overrides the machine's signal in both conflicting cases.

Finally, the belief update when $b^+ > b^-$ is as follows.

$$b_t = \begin{cases} \left[1 + \frac{\bar{b}_{t-1}}{b_{t-1}} \left(\frac{\alpha^W}{\alpha^B} \right)^{1_{\{\Theta_t=A, S_t^M=+\}}} \left(\frac{\bar{\beta}^W}{\bar{\beta}^B} \right)^{1_{\{\Theta_t=NA, S_t^M=+\}}} \right]^{-1} & \text{if } b_{t-1} \geq b^+ \\ \left[1 + \frac{\bar{b}_{t-1}}{b_{t-1}} \left(\frac{\alpha^W}{\alpha^B} \right)^{1_{\{\Theta_t=A, S_t^M=+, S_t^H=+\}}} \left(\frac{\bar{\beta}^W}{\bar{\beta}^B} \right)^{1_{\{\Theta_t=NA, S_t^M=+, S_t^H=+\}}} \right]^{-1} & \text{if } b^+ > b_{t-1} > b^- \\ \left[1 + \frac{\bar{b}_{t-1}}{b_{t-1}} \zeta \right]^{-1} & \text{if } b^- \geq b_{t-1} \end{cases} \quad (32)$$

where ζ is defined in (17).

Case 2. When $b^+ \leq b^-$ we have the following three regimes

- $b_{t-1} > b^-$:

$$\mathbb{P}(\Theta_t = A \mid S_t^H = -, S_t^M = +, b_{t-1}) \geq r \quad (33)$$

$$\mathbb{P}(\Theta_t = A \mid S_t^H = +, S_t^M = -, b_{t-1}) < r \quad (34)$$

In this case, $S_t^M = +$ is sufficient and necessary to act, which implies the machine overrides the human's signal $S_t^H = -$. Further, $S_t^H = +$ is also overruled by $S_t^M = -$.

- $b^- \geq b_{t-1} \geq b^+$:

$$\mathbb{P}(\Theta_t = A \mid S_t^H = -, S_t^M = +, b_{t-1}) \geq r \quad (35)$$

$$\mathbb{P}(\Theta_t = A \mid S_t^H = +, S_t^M = -, b_{t-1}) \geq r \quad (36)$$

In this case, $S_t^M = +$ or $S_t^H = +$ is sufficient for acting.

- $b^+ > b_{t-1}$:

$$\mathbb{P}(\Theta_t = A \mid S_t^H = -, S_t^M = +, b_{t-1}) < r \quad (37)$$

$$\mathbb{P}(\Theta_t = A \mid S_t^H = +, S_t^M = -, b_{t-1}) \geq r \quad (38)$$

In this case, $S_t^H = +$ is sufficient and necessary to act. The signal of the human overrides the machine's signal in both conflicting cases.

Finally the belief update if $b^- \leq b^+$ is as follows.

$$b_t = \begin{cases} \left[1 + \frac{\bar{b}_{t-1}}{b_{t-1}} \left(\frac{\alpha^W}{\alpha^B} \right)^{1_{\{\Theta_t=A, S_t^M=+\}}} \left(\frac{\bar{\beta}^W}{\bar{\beta}^B} \right)^{1_{\{\Theta_t=NA, S_t^M=+\}}} \right]^{-1} & \text{if } b_{t-1} > b^- \\ \left[1 + \frac{\bar{b}_{t-1}}{b_{t-1}} l \right]^{-1} & \text{if } b^- \geq b_{t-1} \geq b^+ \\ \left[1 + \frac{\bar{b}_{t-1}}{b_{t-1}} \zeta \right]^{-1} & \text{if } b^+ > b_{t-1} \end{cases} \quad (39)$$

where ζ is defined in (17) and

$$l \triangleq \left[\left(\frac{\alpha^W}{\alpha^B} \right)^{1_{\{S_t^M=+\}}} \left(\frac{\bar{\alpha}^W}{\bar{\alpha}^B} \right)^{1_{\{S_t^M=-, S_t^H=+\}}} \right]^{1_{\{\Theta_t=A\}}} \left[\left(\frac{\bar{\beta}^W}{\bar{\beta}^B} \right)^{1_{\{S_t^M=+\}}} \left(\frac{\beta^W}{\beta^B} \right)^{1_{\{S_t^M=-, S_t^H=+\}}} \right]^{1_{\{\Theta_t=NA\}}}$$

The log-likelihood ratio process L_t is then obtained by $L_t = \log(b_t/\bar{b}_t)$ in both cases. Q.E.D.

LEMMA 4. Consider the following sequences of random variables

- *i.i.d.* $Y_{1,t}$ with $\mathbb{E}[Y_{1,t}] > 0$ and $|Y_{1,t}| \leq Y_{1,h}$ where $\infty > Y_{1,h} > 0$ for $t = 1, \dots$,
- *i.i.d.* $Y_{2,t}$ with $\mathbb{P}(Y_{2,t} > 0) > 0$ and $|Y_{2,t}| \leq Y_{2,h}$ where $\infty > Y_{2,h} > 0$ for $t = 1, \dots$,
- and, *i.i.d.* $Y_{3,t}$ with $\mathbb{E}[Y_{3,t}] \geq 0$ and $|Y_{3,t}| \leq Y_{3,h}$ where $\infty > Y_{3,h} > 0$ for $t = 1, \dots$

Let Z_t be a discrete stochastic process governed by $Y_{i,t}$ and let two thresholds Z_ℓ and Z_h be such that $Z_h > Z_\ell$ as follows.

$$Z_{t+1} = Z_t + Y_{1,t}1_{\{Z_t \geq Z_h\}} + Y_{2,t}1_{\{Z_t \in (Z_\ell, Z_h)\}} + Y_{3,t}1_{\{Z_t \leq Z_\ell\}}. \quad (40)$$

Then, $Z_t \xrightarrow{a.s.} \infty$.

Proof of Lemma 4. Note that the divergence immediately follows when the mean of $Z_{t+1} - Z_t$ is positive for any Z_t because in that case in all three regions process Z_t is driven by random walks drifting to ∞ . Thus, we focus on the most extreme parameter regime in terms of divergence to ∞ , in particular, when the mean of $Z_{t+1} - Z_t$ is nonpositive in (Z_ℓ, Z_h) and 0 for $Z_t \leq Z_\ell$. To prove this result, we show that there exists a finite (but random) period τ such that process Z_t remains above $Z_t \geq Z_h$ for all $t > \tau$. This is sufficient for divergence to ∞ because process Z_t is governed by a random walk drifting to ∞ above Z_h . We prove this in two steps, but first, we define the following stopping times recursively by assuming, without loss of generality, that $Z_0 > Z_h$.

$$T_1 = \inf\{t : Z_t < Z_h\}, \quad (41)$$

$$\tilde{T}_1 = \inf\{t > T_1 : Z_t \geq Z_h\}, \quad (42)$$

$$T_i = \inf\{t > \tilde{T}_{i-1} : Z_t < Z_h\} \quad \forall i > 1, \quad (43)$$

$$\tilde{T}_i = \inf\{t > T_i : Z_t \geq Z_h\} \quad \forall i > 1. \quad (44)$$

Here, if one of the sets above is empty, i.e., no such t exists, we assign $T_i = \infty$ (and respectively $\tilde{T}_i = \infty$) for the corresponding set.

Step 1. In the first step, we prove that $\mathbb{P}(\tilde{T}_i < \infty | T_i < \infty) = 1$ for all i . In particular, if process Z_t goes below Z_h once, then with probability one, it will cross up Z_h in finite steps. This result also implies that the sequence of infinite stopping times (if it exists) is started by $T_j = \infty$ (but not $\tilde{T}_j = \infty$) for some j . To prove this result, first fix i and assume that T_i is finite; then we define a new sequence of stopping times for process Z_t .

$$V_1 = \inf\{t > T_i : Z_t \leq Z_\ell\} \quad (45)$$

$$\tilde{V}_1 = \inf\{t > V_1 : Z_t > Z_\ell\} \quad (46)$$

$$V_i = \inf\{t > \tilde{V}_{i-1} : Z_t \leq Z_\ell\} \quad \forall i > 1, \quad (47)$$

$$\tilde{V}_i = \inf\{t > V_i : Z_t > Z_\ell\} \quad \forall i > 1. \quad (48)$$

First, note that V_i and \tilde{V}_i for $i \geq 1$ are proper random variables, i.e., they take finite values with probability one. This is because

$$\mathbb{P}(V_i < \infty \mid \tilde{V}_{i-1} < \infty) = 1 \quad (49)$$

$$\mathbb{P}(\tilde{V}_i < \infty \mid V_i < \infty) = 1. \quad (50)$$

The first equality holds because as discussed at the beginning of this proof, we focus on the case where the mean of $Z_{t+1} - Z_t$ in (Z_ℓ, Z_h) is either negative or 0. If negative, Z_t is governed by a random walk drifting to $-\infty$ in (Z_ℓ, Z_h) ; thus it crosses Z_ℓ with probability one in finite steps (see, Theorem 9.1 in Gut 2009, p. 70). If 0, Z_t in (Z_ℓ, Z_h) is a martingale and oscillates (see, Theorem 8.3 in Gut 2009, p. 68). The second equality holds because Z_t is governed again by an oscillating random walk when $Z_t < Z_\ell$ thus it crosses Z_ℓ in finite steps. Furthermore, $T_i < \infty$; thus, V_1 is also finite. Therefore, it follows that $V_i < \infty$ and $\tilde{V}_i < \infty$ for $i \geq 1$.

To prove that \tilde{T}_i is finite, we will show that there exists a finite j such that $\tilde{T}_i < V_j$. To do so, define events $A_t = \{\tilde{T}_i \geq V_t\}$. Then, the Borel-Cantelli lemma implies the following,

$$\sum_{t=T_i}^{\infty} \mathbb{P}(A_t) < \infty \Rightarrow \mathbb{P}(\limsup_{i \rightarrow \infty} A_t) = 0 \quad (51)$$

First, consider $\mathbb{P}(A_1) = \mathbb{P}(\tilde{T}_i \geq V_1)$, this probability is bounded, i.e., $\mathbb{P}(A_1) < \delta < 1$ because $Z_h - Z_\ell$ is bounded and Z_t has positive size jumps with positive probability in (Z_ℓ, Z_h) . Proceeding similarly, we obtain the following

$$\mathbb{P}(\tilde{T}_i \geq V_t) = \mathbb{P}(\tilde{T}_i \geq V_{t-1})\mathbb{P}(\tilde{T}_i \geq V_t \mid \tilde{T}_i \geq V_{t-1}) < \delta^t \quad (52)$$

Therefore, it follows that $\mathbb{P}(\limsup_{i \rightarrow \infty} A_t) = 0$, which implies that there exists a finite j such that $\tilde{T}_i < V_j$. As discussed, V_j is finite; thus, it follows that $\mathbb{P}(\tilde{T}_i < \infty \mid T_i < \infty) = 1$.

Step 2. In this step, we show that there exists a finite j such that $T_j = \infty$ and $\tilde{T}_k < \infty$ for all $k < j$. Define the following events $E_t = \{T_t < \infty\}$ for $t > 1$. Then, the Borel-Cantelli lemma implies the following

$$\sum_{t=1}^{\infty} \mathbb{P}(E_t) < \infty \Rightarrow \mathbb{P}(\limsup_{t \rightarrow \infty} E_t) = 0. \quad (53)$$

In particular, if the summation condition is satisfied, then events E_t cannot occur infinitely many times, i.e., there exists a finite j such that $T_j = \infty$. To show that the summation condition is indeed satisfied, we construct a finite upper bound for it. Above Z_h , process Z_t is driven by a random walk drifting to ∞ . Thus, stopping time T_1 is defective, i.e., $\mathbb{P}(T_1 < \infty) < \varphi < 1$ for some φ (see Theorem 9.1 in Gut 2009, p.70). Moreover, we bound $\mathbb{P}(T_2 < \infty)$ as follows:

$$\mathbb{P}(T_2 < \infty) = \mathbb{P}(T_2 < \infty \mid \tilde{T}_1 < \infty)\mathbb{P}(\tilde{T}_1 < \infty \mid T_1 < \infty)\mathbb{P}(T_1 < \infty) < \varphi^2. \quad (54)$$

In Step 1 of this proof, we show that $\mathbb{P}(\tilde{T}_1 < \infty | T_1 < \infty) = 1$. Moreover, we have $\mathbb{P}(T_2 < \infty | \tilde{T}_1 < \infty) < \varphi < 1$ because $\tilde{T}_1 < \infty$ implies that process Z_t crosses Z_h up in finite steps, and after crossing Z_h , process Z_t is again driven by the same random walk drifting to ∞ . Proceeding similarly, it follows that $\mathbb{P}(E_t) < \varphi^t$. Hence, there exists a finite (but random) j such that $T_j = \infty$, and Step 1 implies that $\tilde{T}_k < \infty$ for $k < j$ because $T_k < \infty$. Therefore, after sufficiently large t , process Z_t always remains above Z_h and diverges to ∞ . Q.E.D.

Proof of Theorem 3. In this proof, we focus on the log-likelihood ratio process L_t because as also discussed in the proof of Theorem 2, the continuous mapping theorem implies that the limit of L_t characterizes the limit of b_t . Our proof is based on analyzing the mean of $L_{t+1} - L_t$ when L_t takes different values in comparison to L_h and L_ℓ for Case 1 and Case 2 characterized in Lemma 3. In both cases, process L_t is governed by different random walks with random jumps whose means and size change depending on the previous state L_{t-1} .⁹ Trivial cases arise when the mean of $L_{t+1} - L_t$ always remains positive regardless of L_t . In particular, L_t diverges to ∞ when the mean of $L_{t+1} - L_t$ is always positive despite changing values of L_t due to the strong law of large numbers because three random walks (for $L_t \geq L_h$, $L_t \in (L_h, L_\ell)$ and $L_t \leq L_\ell$ under Case 1) that establish the trajectory of process L_t all diverge to ∞ (see Theorem 8.3 in Gut 2009, p. 68). As a result, b_t converges to 0 as L_t diverges to ∞ . We consider the different values of p in the statement of the theorem separately.

Step 1. Let $p \leq p^B$. To prove that L_t and hence b_t oscillate, we need to show that there exists a number that process L_t crosses infinitely often (see, for instance, Vatutin and Wachtel 2009 for a mathematical definition of oscillation). This property (oscillation) holds for a random walk with a noise term whose mean is 0 (see Theorem 8.2 in Gut 2009, p. 68). However, process L_t does not satisfy this property because the mean of $L_{t+1} - L_t$ is always positive when $L_t \leq L_\ell$ in Case 1 ($L_t < L_\ell$ in Case 2) as discussed. Nevertheless, we can show that process L_t oscillates because the mean of $L_{t+1} - L_t$ is either negative or zero when $p \leq p^B$ for $L_t \geq L_h$ in Case 1 ($L_t > L_h$ in Case 2). Therefore, process L_t returns to interval (L_ℓ, L_h) in finite steps after lying outside it. Oscillation is regardless of the sign of the mean of $L_{t+1} - L_t$ in (L_ℓ, L_h) because process L_t goes out of (L_ℓ, L_h) in finite steps. Specifically, in finite steps, process L_t i) crosses L_h up if the mean of $L_{t+1} - L_t$ in (L_ℓ, L_h) is positive, ii) crosses L_ℓ down if it is negative, and iii) goes out of (L_ℓ, L_h) if it is zero.

To illustrate this process, we focus on Case 1 and the setting where the mean of $L_{t+1} - L_t$ in (L_ℓ, L_h) is negative. Define the following stopping times

$$J = \inf\{t \geq 1 : L_t \leq L_\ell\} \tag{55}$$

$$\tilde{J} = \inf\{t \geq 1 : L_t > L_\ell\} \tag{56}$$

⁹ The log-likelihood ratio process L_t is similar to the *oscillating random walk* defined in Kemperman (1974) such that the partition of \mathbb{R} consists of $(-\infty, L_\ell]$, (L_ℓ, L_h) and $[L_h, \infty)$.

If J and \tilde{J} are proper random variables (i.e., $\mathbb{P}(J < \infty) = \mathbb{P}(\tilde{J} < \infty) = 1$), then process L_t oscillates. Assume $L_0 > L_\ell$ without loss of generality; then J is a proper (almost surely finite) random variable because process L_t is driven by a random walk drifting to $-\infty$ for $L_t > L_\ell$. Thus, in finite steps it will cross L_ℓ down. After J steps, process L_t is driven by a random walk drifting to ∞ ,¹⁰ thus it crosses L_ℓ up in finite steps, so $\tilde{J} - J$ is a proper random variable. Since J is also proper, $\tilde{J} = \tilde{J} - J + J$ is also proper. Since L_t crosses L_ℓ infinitely often, b_t crosses $1/(1 + \exp(L_\ell))$ infinitely often and, hence, oscillates. The same approach can be repeated for the remaining setting where the mean of $L_{t+1} - L_t$ is positive or zero in (L_ℓ, L_h) .

The oscillation property also implies that the process L_t and hence b_t is recurrent because the mean of the random jumps is bounded. With probability 1, process L_t will revisit interval (L_ℓ, L_h) for Case 1 in finite steps interval $[L_\ell, L_h]$ for Case 2 in finite steps. Furthermore, intervals that can be reached from (L_ℓ, L_h) with positive probability are also recurrent, which implies Corollary 1.

Step 2. Let $p > p^B$. Then, the mean of $L_{t+1} - L_t$ when $L_t \leq L_\ell$ for Case 1 and $L_t < L_\ell$ for Case 2 is positive (see Footnote 10). Further, we know from the proof of Theorem 2 that the mean of $L_{t+1} - L_t$ when $L_t \geq L_h$ for Case 1 and $L_t > L_h$ for $p > p^B$ is also positive. In Case 1, the mean of $L_{t+1} - L_t$ for $L_t \in (L_\ell, L_h)$ is

$$\alpha^H \alpha^B p \log \left[\frac{\alpha^B}{\alpha^W} \right] + \bar{\beta}^H \bar{\beta}^B \bar{p} \log \left[\frac{\bar{\beta}^B}{\bar{\beta}^W} \right] = \bar{\beta}^H \left[\alpha^B p \log \left[\frac{\alpha^B}{\alpha^W} \right] + \bar{\beta}^B \bar{p} \log \left[\frac{\bar{\beta}^B}{\bar{\beta}^W} \right] \right] + \bar{\beta}^H \alpha^B \left[\frac{\alpha^H}{\bar{\beta}^H} - 1 \right] \log \left[\frac{\alpha^B}{\alpha^W} \right].$$

Note that the term above is positive for $p > p^B$. Hence, as discussed L_t diverges to ∞ and b_t converges to 1 for Case 1 when p is larger than p^B . In Case 2, the mean of $L_{t+1} - L_t$ for $L_t \in [L_\ell, L_h]$ is not necessarily positive. When it is positive, we immediately obtain the same result. Nevertheless, we obtain the same divergence despite having negative or zero mean of $L_{t+1} - L_t$ for $L_t \in [L_\ell, L_h]$ as long as the means of $L_{t+1} - L_t$ for $L_t < L_\ell$ and $L_t > L_h$ are positive, which is true as proven in Lemma 4. Thus, using Lemma 4, we capture all possible values of L_t with respect to L_h and L_ℓ in Cases 1 and 2 for $p > p^B$, and show that L_t diverges to ∞ , which implies b_t converges to 1. Hence, we conclude the proof. Q.E.D.

Proof of Theorem 4. The first part of the result in this theorem when $p \leq p^W$ follows from Lemma 4 by considering the reflection of the stochastic process. In particular, updating the conditions in the statement of Lemma 4 as $\mathbb{E}[Y_{1,t}] \leq 0$, $\mathbb{P}(Y_{2,t} < 0) > 0$ and $\mathbb{E}[Y_{3,t}] < 0$ would imply $Z_t \xrightarrow{\text{a.s.}} -\infty$ in that lemma. This is because the mean of $L_{t+1} - L_t$ is nonpositive when $L_t \in [L_h, \infty)$ in Case 1 (and $L_t \in (L_h, \infty)$ in Case 2); and is negative when $L_t \in (-\infty, L_\ell]$ in Case 1 (and $L_t \in (-\infty, L_\ell)$ in Case 2).¹¹ Thus, the case for $p \leq p^W$ follows from Lemma 4 with a slight modification.

¹⁰ This claim can be proved by invoking Lemma A.2 in Harrison et al. (2012) to show that $\mathbb{E}^B[R_t^{\text{HM}}(L_{t-1})] > 0$ for $L_{t-1} \leq L_\ell$ in Case 1, and $L_{t-1} < L_\ell$ in Case 2.

¹¹ The second part of this claim can be proved by invoking Lemma A.2 in Harrison et al. (2012) to show that $\mathbb{E}^B[R_t^{\text{HM}}(L_{t-1})] < 0$ for $L_{t-1} \leq L_\ell$ in Case 1, and $L_{t-1} < L_\ell$ in Case 2.

To prove the second part of this result $p > p^w$, we focus on Case 1 (Case 2 can be addressed by adjusting weak and strict inequalities in the same way) by assuming the mean of $L_{t+1} - L_t$ is negative when $L_t \in (L_\ell, L_h)$ and define the following stopping times by assuming $L_0 > L_h$ without loss of generality.

$$T_1 = \inf\{t : L_t < L_h\} \tag{57}$$

$$\tilde{T}_1 = \inf\{t > T_1 : L_t \geq L_h\} \tag{58}$$

$$T_i = \inf\{t > \tilde{T}_{i-1} : L_t < L_h\} \quad \text{for } i \geq 2 \tag{59}$$

$$\tilde{T}_i = \inf\{t > T_i : L_t \geq L_h\} \quad \text{for } i \geq 2 \tag{60}$$

Here, if a set is empty, then the stopping times takes the value of ∞ . Therefore, if one of the stopping times T_i or \tilde{T}_i is not finite for some i , then all the following stopping times for $j > i$ also are ∞ . We first show that $\sum_{i=1}^{\infty} \mathbb{P}(T_i < \infty) < \infty$ and then use the Borel-Cantelli lemma to deduce that there exists a finite j such that $\mathbb{P}(T_j = \infty \text{ or } \tilde{T}_j = \infty) = 1$. If $T_j = \infty$, then $L_t \rightarrow \infty$; otherwise, ($\tilde{T}_j = \infty$) then $L_t \rightarrow -\infty$. Thus, $L_t \rightarrow L$, where $L \in \{-\infty, \infty\}$; hence b_t converges to a Bernoulli random variable.

First, consider T_1 ; it follows that $\mathbb{P}(T_1 < \infty) < \xi < 1$ (i.e., T_1 is a defective random variable) because process L_t is driven by a random walk drifting to ∞ when $L_t > L_h$ and $p > p^w$. Next, considering \tilde{T}_1 and T_2 ; we obtain that

$$\mathbb{P}(\tilde{T}_1 < \infty) = \mathbb{P}(\tilde{T}_1 < \infty | T_1 < \infty) \mathbb{P}(T_1 < \infty) < \xi^2 \tag{61}$$

$$\mathbb{P}(T_2 < \infty) = \mathbb{P}(T_2 < \infty | \tilde{T}_1 < \infty) \mathbb{P}(\tilde{T}_1 < \infty) < \xi^3 \tag{62}$$

Here, the first term is strictly less than 1, i.e., $\mathbb{P}(\tilde{T}_1 < \infty | T_1 < \infty) < \xi < 1$ because process L_t is driven by random walks drifting to $-\infty$ after T_1 . The inequality in the second line follows because L_t is driven by a random walk drifting to ∞ . Proceeding similarly, it follows that $\mathbb{P}(T_i < \infty) < \xi^{(2i-1)}$ and hence $\sum_{i=1}^{\infty} \mathbb{P}(T_i < \infty)$ is finite.

Note that we assume at the beginning that the mean of $L_{t+1} - L_t$ is negative when $L_t \in (L_\ell, L_h)$. If it is positive, the same steps can be followed by redefining stopping times using L_ℓ as the threshold instead of L_h . If it is zero, then stopping times T_i is defined by considering the time when process L_t enters (L_ℓ, L_h) and \tilde{T}_i is the time when L_t exits (L_ℓ, L_h) in any direction. The approach follows because outside (L_ℓ, L_h) , process L_t diverges (either to ∞ or $-\infty$ depending on being above L_h or below L_ℓ) and it oscillates in (L_ℓ, L_h) , which guarantees it remains in (L_ℓ, L_h) for finite steps.

More specifically, the main driver of this result is that process L_t may diverge to ∞ when above L_h and to $-\infty$ when below L_ℓ . The sign of the mean between L_ℓ and L_h does not affect this characteristic in the limit. Thus, we conclude the proof. Q.E.D.

Appendix C: Mistrust Bias against the Machine

Proof of Lemma 2. Note that Informativeness (1)-(3) imply that the posterior after + (after -) is larger than or equal to r (resp. lower than r). Thus, there exist thresholds $\lambda_{\Gamma}^{-}, \lambda_{\Gamma}^{+} \in (0, 1)$ for $\Gamma \in \{\mathbf{B}, \mathbf{W}\}$ given by

$$\lambda_{\Gamma}^{-} \triangleq \frac{r - \mathbb{P}^{\Gamma}(\Theta_t = \mathbf{A} | S^{\mathbf{M}} = -)}{\mathbb{P}(\Theta_t = \mathbf{A} | S^{\mathbf{H}} = +) - \mathbb{P}^{\Gamma}(\Theta_t = \mathbf{A} | S^{\mathbf{M}} = -)}, \quad (63)$$

$$\lambda_{\Gamma}^{+} \triangleq \frac{\mathbb{P}^{\Gamma}(\Theta_t = \mathbf{A} | S^{\mathbf{M}} = +) - r}{\mathbb{P}^{\Gamma}(\Theta_t = \mathbf{A} | S^{\mathbf{M}} = +) - \mathbb{P}(\Theta_t = \mathbf{A} | S^{\mathbf{H}} = -)}. \quad (64)$$

These equations imply that

$$\lambda \mathbb{P}(\Theta_t = \mathbf{A} | S^{\mathbf{H}} = +) + \bar{\lambda} \mathbb{P}^{\mathbf{B}}(\Theta_t = \mathbf{A} | S^{\mathbf{M}} = -) \geq r \Leftrightarrow \lambda \geq \lambda_{\mathbf{B}}^{-}, \quad (65)$$

$$\lambda \mathbb{P}(\Theta_t = \mathbf{A} | S^{\mathbf{H}} = -) + \bar{\lambda} \mathbb{P}^{\mathbf{B}}(\Theta_t = \mathbf{A} | S^{\mathbf{M}} = +) \leq r \Leftrightarrow \lambda \geq \lambda_{\mathbf{B}}^{+}, \quad (66)$$

$$\lambda \mathbb{P}(\Theta_t = \mathbf{A} | S^{\mathbf{H}} = +) + \bar{\lambda} \mathbb{P}^{\mathbf{W}}(\Theta_t = \mathbf{A} | S^{\mathbf{M}} = -) \leq r \Leftrightarrow \lambda \leq \lambda_{\mathbf{W}}^{-}, \quad (67)$$

$$\lambda \mathbb{P}(\Theta_t = \mathbf{A} | S^{\mathbf{H}} = -) + \bar{\lambda} \mathbb{P}^{\mathbf{W}}(\Theta_t = \mathbf{A} | S^{\mathbf{M}} = +) \geq r \Leftrightarrow \lambda \leq \lambda_{\mathbf{W}}^{+}, \quad (68)$$

because the left-hand sides of the first and third inequalities are increasing and the second and the fourth ones are decreasing in λ .

First, note that λ_{Γ}^{+} is increasing in $\mathbb{P}^{\Gamma}(\Theta_t = \mathbf{A} | S^{\mathbf{M}} = +)$, and λ_{Γ}^{-} is decreasing in $\mathbb{P}^{\Gamma}(\Theta_t = \mathbf{A} | S^{\mathbf{M}} = -)$. We also know from Substitution (4)-(5) that $\mathbb{P}^{\mathbf{B}}(\Theta_t = \mathbf{A} | S^{\mathbf{M}} = +) > \mathbb{P}^{\mathbf{W}}(\Theta_t = \mathbf{A} | S^{\mathbf{M}} = +)$, and $\mathbb{P}^{\mathbf{W}}(\Theta_t = \mathbf{A} | S^{\mathbf{M}} = -) > \mathbb{P}^{\mathbf{B}}(\Theta_t = \mathbf{A} | S^{\mathbf{M}} = -)$. We conclude the proof by defining $\lambda_{min} \triangleq \max(\lambda_{\mathbf{W}}^{-}, \lambda_{\mathbf{W}}^{+})$ and $\lambda_{max} = \min(\lambda_{\mathbf{B}}^{-}, \lambda_{\mathbf{B}}^{+})$. Q.E.D.

Before proving Theorem 5, we first provide a constructive lemma and its proof.

LEMMA 5. *For $\lambda \in (\lambda_{min}, \lambda_{max})$, unique thresholds $b_{\lambda}^{-} \in (0, 1)$ and $b_{\lambda}^{+} \in (0, 1)$ exist such that*

$$\lambda \mathbb{P}(\Theta_t = \mathbf{A} | S^{\mathbf{H}} = +) + (1 - \lambda) \mathbb{P}(\Theta_t = \mathbf{A} | S^{\mathbf{M}} = -, b_{t-1}) \geq r \Leftrightarrow b_{t-1} \leq b_{\lambda}^{-}, \quad (69)$$

$$\lambda \mathbb{P}(\Theta_t = \mathbf{A} | S^{\mathbf{H}} = -) + (1 - \lambda) \mathbb{P}(\Theta_t = \mathbf{A} | S^{\mathbf{M}} = +, b_{t-1}) \geq r \Leftrightarrow b_{t-1} \leq b_{\lambda}^{+}. \quad (70)$$

Proof of Lemma 5. We obtain the thresholds by solving the following equations.

$$\lambda \mathbb{P}(\Theta_t = \mathbf{A} | S^{\mathbf{H}} = +) + (1 - \lambda) \mathbb{P}(\Theta_t = \mathbf{A} | S^{\mathbf{M}} = -, b_{\lambda}^{-}) = r, \quad (71)$$

$$\lambda \mathbb{P}(\Theta_t = \mathbf{A} | S^{\mathbf{H}} = -) + (1 - \lambda) \mathbb{P}(\Theta_t = \mathbf{A} | S^{\mathbf{M}} = +, b_{\lambda}^{+}) = r. \quad (72)$$

The closed-form expressions are

$$b_{\lambda}^{-} = \frac{\varphi_{\lambda}^{-}(1 - \alpha^{\mathbf{W}}) - \beta^{\mathbf{W}}}{\beta^{\mathbf{B}} - \beta^{\mathbf{W}} + \varphi_{\lambda}^{-}(\alpha^{\mathbf{B}} - \alpha^{\mathbf{W}})} \quad \text{and} \quad b_{\lambda}^{+} = \frac{\bar{\beta}^{\mathbf{W}} - \varphi_{\lambda}^{+} \alpha^{\mathbf{W}}}{\varphi_{\lambda}^{+}(\alpha^{\mathbf{B}} - \alpha^{\mathbf{W}}) + \beta^{\mathbf{B}} - \beta^{\mathbf{W}}}$$

where

$$\varphi_{\lambda}^{-} = \frac{p}{\bar{p}} \left[\frac{(1 - \lambda) - r + \lambda \mathbb{P}(\Theta_t = \mathbf{A} | S_t^{\mathbf{H}} = +)}{r - \lambda \mathbb{P}(\Theta_t = \mathbf{A} | S_t^{\mathbf{H}} = +)} \right] \quad \text{and} \quad \varphi_{\lambda}^{+} = \frac{p}{\bar{p}} \left[\frac{(1 - \lambda) - r + \lambda \mathbb{P}(\Theta_t = \mathbf{A} | S_t^{\mathbf{H}} = -)}{r - \lambda \mathbb{P}(\Theta_t = \mathbf{A} | S_t^{\mathbf{H}} = -)} \right].$$

Finally, we can similarly define $b_{\lambda}^{\mathbf{H}} = \min(b_{\lambda}^{-}, b_{\lambda}^{+})$ and $b_{\lambda}^{\mathbf{M}} = \max(b_{\lambda}^{-}, b_{\lambda}^{+})$ as in Corollary 1. Q.E.D.

Proof of Theorem 5. Note that Lemma 5 shows that there exist thresholds $b_\lambda^-, b_\lambda^+ \in (0, 1)$ for $\lambda \in (\lambda_{min}, \lambda_{max})$ as in the case of Lemma 1. Thus, the exact same steps at which thresholds $b^-, b^+ \in (0, 1)$ are replaced by $b_\lambda^-, b_\lambda^+ \in (0, 1)$ in the proofs of Lemma 3, Theorems 3-4 will imply this result because the results in those theorems do not depend on the exact values of thresholds b^- and b^+ , as long as they are interior to $(0, 1)$. Q.E.D.

Proof of Theorem 6. Note that Lemma 1 also characterizes the DM's decision rule for the biased case. Nevertheless, the belief updating (6) is replaced with the one in Section 7.2 with a bias term μ when the machine's prediction is not correct. Thus, we need to modify Lemma 3 slightly to incorporate this change. In particular, the updated log-likelihood ratio process \tilde{L}_t is as follows.

Case 1: If $b^+ > b^-$, we have $\tilde{L}_t = \tilde{L}_{t-1} + \tilde{R}_t^{\text{HM}}(\tilde{L}_{t-1})$ where

$$\tilde{R}_t^{\text{HM}}(\tilde{L}_{t-1}) \triangleq \begin{cases} 1_{\{S_t^{\text{M}}=+\}} \left[1_{\{\Theta_t=\text{A}\}} \log\left(\frac{\alpha^{\text{B}}}{\alpha^{\text{W}}}\right) + 1_{\{\Theta_t=\text{NA}\}} \mu \log\left(\frac{\bar{\beta}^{\text{B}}}{\bar{\beta}^{\text{W}}}\right) \right] & \text{if } \tilde{L}_{t-1} \geq L_h \\ 1_{\{S_t^{\text{H}}=+, S_t^{\text{M}}=+\}} \left[1_{\{\Theta_t=\text{A}\}} \log\left(\frac{\alpha^{\text{B}}}{\alpha^{\text{W}}}\right) + 1_{\{\Theta_t=\text{NA}\}} \mu \log\left(\frac{\bar{\beta}^{\text{B}}}{\bar{\beta}^{\text{W}}}\right) \right] & \text{if } L_h > \tilde{L}_{t-1} > L_\ell \\ 1_{\{S_t^{\text{H}}=+\}} [\tilde{v}_1 + \tilde{v}_2] & \text{if } L_\ell \geq \tilde{L}_{t-1} \end{cases}$$

with $L_h \triangleq \log\left(\frac{b^+}{b^+}\right)$, $L_\ell \triangleq \log\left(\frac{b^-}{b^-}\right)$ and

$$\begin{aligned} \tilde{v}_1 &\triangleq 1_{\{\Theta_t=\text{A}\}} \left(1_{\{S_t^{\text{M}}=-\}} \mu \log\left(\frac{\bar{\alpha}^{\text{B}}}{\bar{\alpha}^{\text{W}}}\right) + 1_{\{S_t^{\text{M}}=+\}} \log\left(\frac{\alpha^{\text{B}}}{\alpha^{\text{W}}}\right) \right), \\ \tilde{v}_2 &\triangleq 1_{\{\Theta_t=\text{NA}\}} \left(1_{\{S_t^{\text{M}}=-\}} \log\left(\frac{\beta^{\text{B}}}{\beta^{\text{W}}}\right) + 1_{\{S_t^{\text{M}}=+\}} \mu \log\left(\frac{\bar{\beta}^{\text{B}}}{\bar{\beta}^{\text{W}}}\right) \right). \end{aligned}$$

Case 2: If $b^+ \leq b^-$, we have $\tilde{L}_t = \tilde{L}_{t-1} + \tilde{R}_t^{\text{HM}}(\tilde{L}_{t-1})$ where

$$\tilde{R}_t^{\text{HM}}(\tilde{L}_{t-1}) \triangleq \begin{cases} 1_{\{S_t^{\text{M}}=+\}} \left[1_{\{\Theta_t=\text{A}\}} \log\left(\frac{\alpha^{\text{B}}}{\alpha^{\text{W}}}\right) + 1_{\{\Theta_t=\text{NA}\}} \mu \log\left(\frac{\bar{\beta}^{\text{B}}}{\bar{\beta}^{\text{W}}}\right) \right] & \text{if } \tilde{L}_{t-1} > L_h \\ 1_{\{\Theta_t=\text{A}\}} \tilde{v}_1 + 1_{\{\Theta_t=\text{NA}\}} \tilde{v}_2 & \text{if } L_h \geq \tilde{L}_{t-1} \geq L_\ell \\ 1_{\{S_t^{\text{H}}=+\}} [\tilde{v}_1 + \tilde{v}_2] & \text{if } L_\ell > \tilde{L}_{t-1} \end{cases} \quad (73)$$

with $L_h \triangleq \log\left(\frac{b^-}{b^-}\right)$, $L_\ell \triangleq \log\left(\frac{b^+}{b^+}\right)$ and

$$\begin{aligned} \tilde{v}_1 &\triangleq 1_{\{S_t^{\text{M}}=+\}} \log\left(\frac{\alpha^{\text{B}}}{\alpha^{\text{W}}}\right) + 1_{\{S_t^{\text{M}}=-\}} 1_{\{S_t^{\text{H}}=+\}} \mu \log\left(\frac{\bar{\alpha}^{\text{B}}}{\bar{\alpha}^{\text{W}}}\right), \\ \tilde{v}_2 &\triangleq 1_{\{S_t^{\text{M}}=+\}} \mu \log\left(\frac{\bar{\beta}^{\text{B}}}{\bar{\beta}^{\text{W}}}\right) + 1_{\{S_t^{\text{M}}=-\}} 1_{\{S_t^{\text{H}}=+\}} \log\left(\frac{\beta^{\text{B}}}{\beta^{\text{W}}}\right). \end{aligned}$$

Note that the thresholds in the theorem are determined as the break-even points of the following equations.

$$p\alpha^{\text{B}} \log\left(\frac{\alpha^{\text{B}}}{\alpha^{\text{W}}}\right) + \bar{p}\bar{\beta}^{\text{B}} \mu^{\text{B}} \log\left(\frac{\bar{\beta}^{\text{B}}}{\bar{\beta}^{\text{W}}}\right) = 0 \quad (74)$$

$$p\alpha^{\text{W}} \log\left(\frac{\alpha^{\text{B}}}{\alpha^{\text{W}}}\right) + \bar{p}\bar{\beta}^{\text{W}} \mu^{\text{W}} \log\left(\frac{\bar{\beta}^{\text{B}}}{\bar{\beta}^{\text{W}}}\right) = 0 \quad (75)$$

$$p\alpha^{\text{H}} \left(\bar{\alpha}^{\text{B}} \mu^{\text{H}} \log\left(\frac{\bar{\alpha}^{\text{B}}}{\bar{\alpha}^{\text{W}}}\right) + \alpha^{\text{B}} \log\left(\frac{\alpha^{\text{B}}}{\alpha^{\text{W}}}\right) \right) + \bar{p}\bar{\beta}^{\text{H}} \left(\beta^{\text{B}} \log\left(\frac{\beta^{\text{B}}}{\beta^{\text{W}}}\right) + \bar{\beta}^{\text{B}} \mu^{\text{H}} \log\left(\frac{\bar{\beta}^{\text{B}}}{\bar{\beta}^{\text{W}}}\right) \right) = 0 \quad (76)$$

The left-hand sides of these equations correspond to the mean of $\tilde{L}_{t+1} - \tilde{L}_t$ when evaluated at μ in place of the thresholds at corresponding values of \tilde{L}_t . In the remainder of the proof, we explain the sign of the mean of $\tilde{L}_{t+1} - \tilde{L}_t$ given the machine's type and the value of μ . When the sign of this mean is known, the results follow from the proofs of previous theorems on which we elaborate.

- when $\Gamma = B$,

— $\mu \geq \mu^B$ and $\mu > \mu^H$ implies that the mean of $\tilde{L}_{t+1} - \tilde{L}_t$ for $L_\ell \geq \tilde{L}_t$ for Case 1 ($L_\ell > \tilde{L}_t$ for Case 2) is negative. Further, the mean of $\tilde{L}_{t+1} - \tilde{L}_t$ for $\tilde{L}_t \geq L_h$ for Case 1 ($\tilde{L}_t > L_h$ for Case 2) is nonpositive. Thus, the proof of the first part (for $p \leq p^W$) of Theorem 4 applies to this parameter regime.

— $\mu^B > \mu > \mu^H$ implies that the mean of $\tilde{L}_{t+1} - \tilde{L}_t$ for $L_\ell \geq \tilde{L}_t$ for Case 1 ($L_\ell > \tilde{L}_t$ for Case 2) is positive. Further, the mean of $\tilde{L}_{t+1} - \tilde{L}_t$ for $\tilde{L}_t \geq L_h$ for Case 1 ($\tilde{L}_t > L_h$ for Case 2) is negative. Thus, the proof of the second part (for $p > p^W$) of Theorem 4 applies to this parameter regime.

— $\mu^H \geq \mu \geq \mu^B$ implies that the mean of $\tilde{L}_{t+1} - \tilde{L}_t$ for $L_\ell \geq \tilde{L}_t$ for Case 1 ($L_\ell > \tilde{L}_t$ for Case 2) is nonnegative. Further, the mean of $\tilde{L}_{t+1} - \tilde{L}_t$ for $\tilde{L}_t \geq L_h$ for Case 1 ($\tilde{L}_t > L_h$ for Case 2) is nonpositive. Thus, the proof of the first part (for $p \leq p^B$) of Theorem 3 applies to this parameter regime.

— if $\mu^B > \mu$ and $\mu^H \geq \mu$ implies that the mean of $\tilde{L}_{t+1} - \tilde{L}_t$ for $L_\ell \geq \tilde{L}_t$ for Case 1 ($L_\ell > \tilde{L}_t$ for Case 2) is nonnegative. Further, the mean of $\tilde{L}_{t+1} - \tilde{L}_t$ for $\tilde{L}_t \geq L_h$ for Case 1 ($\tilde{L}_t > L_h$ for Case 2) is positive. Thus, the proof of the second part (for $p > p^B$) of Theorem 3 applies to this parameter regime.

- when $\Gamma = W$,

— $\mu \geq \mu^W$ implies that the mean of $\tilde{L}_{t+1} - \tilde{L}_t$ for $\tilde{L}_t \geq L_h$ for Case 1 ($\tilde{L}_t > L_h$ for Case 2) is nonpositive. Further, the mean of $\tilde{L}_{t+1} - \tilde{L}_t$ for $L_\ell \geq \tilde{L}_t$ for Case 1 ($L_\ell > \tilde{L}_t$ for Case 2) is negative for all $\mu \geq 1$. Thus, the proof of the first part (for $p \leq p^W$) of Theorem 4 applies to this parameter regime.

— $\mu^W > \mu$ implies that the mean of $\tilde{L}_{t+1} - \tilde{L}_t$ for $\tilde{L}_t \geq L_h$ for Case 1 ($\tilde{L}_t > L_h$ for Case 2) is positive. Further, the mean of $\tilde{L}_{t+1} - \tilde{L}_t$ for $L_\ell \geq \tilde{L}_t$ for Case 1 ($L_\ell > \tilde{L}_t$ for Case 2) is negative for all $\mu \geq 1$. Thus, the proof of the second part (for $p > p^W$) of Theorem 4 applies to this parameter regime.

Finally, we show that $\mu^H > 1$. Note that the following term is always positive when evaluated at $\mu^H = 1$; see Footnote 10.

$$p\alpha^H \left(\bar{\alpha}^B \mu^H \log \left(\frac{\bar{\alpha}^B}{\bar{\alpha}^W} \right) + \alpha^B \log \left(\frac{\alpha^B}{\alpha^W} \right) \right) + \bar{p}\bar{\beta}^H \left(\beta^B \log \left(\frac{\beta^B}{\beta^W} \right) + \bar{\beta}^B \mu^H \log \left(\frac{\bar{\beta}^B}{\bar{\beta}^W} \right) \right)$$

To make this equal to 0, μ^H has to be larger than 1 because it is decreasing in μ^H . Hence, we conclude the proof. Q.E.D.

Appendix D: Complementarity

Before providing the proof of Theorem 7, we first provide some constructive lemmas and their proofs. The first lemma is analogous to Lemmas 1 and 2 in terms of characterizing the DM's decision rule as a function of her belief.

LEMMA 6. *Unique thresholds b_C^- and b_C^+ exist such that*

$$\mathbb{P}(\Theta_t = A \mid S_t^H = +, S_t^M = -, b_{t-1}) \geq r \iff b_{t-1} \leq b_C^-, \quad (77)$$

$$\mathbb{P}(\Theta_t = A \mid S_t^H = -, S_t^M = +, b_{t-1}) \geq r \iff b_{t-1} \leq b_C^+. \quad (78)$$

Proof of Lemma 6. As mentioned in the proof of Lemma 1, the posterior probabilities in the statement of this lemma are continuous and monotone in b_{t-1} and the boundary values 1 and 0 are on different sides of r due to Complementarity (13)-(14). We next provide the posterior probabilities.

$$\mathbb{P}(\Theta_t = A \mid S_t^H = +, S_t^M = -, b_{t-1}) = \frac{\alpha^H(b_{t-1}\bar{\alpha}^{C^+} + \bar{b}_{t-1}\bar{\alpha}^{C^-})p}{\alpha^H(b_{t-1}\bar{\alpha}^{C^+} + \bar{b}_{t-1}\bar{\alpha}^{C^-})p + \beta^H(b_{t-1}\beta^{C^+} + \bar{b}_{t-1}\beta^{C^-})\bar{p}} \quad (79)$$

$$\mathbb{P}(\Theta_t = A \mid S_t^H = -, S_t^M = +, b_{t-1}) = \frac{\bar{\alpha}^H(b_{t-1}\alpha^{C^+} + \bar{b}_{t-1}\alpha^{C^-})p}{\bar{\alpha}^H(b_{t-1}\alpha^{C^+} + \bar{b}_{t-1}\alpha^{C^-})p + \beta^H(b_{t-1}\bar{\beta}^{C^+} + \bar{b}_{t-1}\bar{\beta}^{C^-})\bar{p}} \quad (80)$$

We solve the following equations for b_{t-1}

$$\mathbb{P}(\Theta_t = A \mid S_t^H = +, S_t^M = -, b_{t-1}) = r \text{ and } \mathbb{P}(\Theta_t = A \mid S_t^H = -, S_t^M = +, b_{t-1}) = r. \quad (81)$$

We then obtain

$$b_C^- = \frac{\left(\frac{\bar{r}p\alpha^H}{r\bar{p}\beta^H}\right)\bar{\alpha}^{C^-} - \beta^{C^-}}{\left(\frac{\bar{r}p\alpha^H}{r\bar{p}\beta^H}\right)[\alpha^{C^+} - \alpha^{C^-}] + \beta^{C^+} - \beta^{C^-}} \text{ and } b_C^+ = \frac{\bar{\beta}^{C^-} - \left(\frac{\bar{r}p\bar{\alpha}^H}{r\bar{p}\beta^H}\right)\alpha^{C^-}}{\left(\frac{\bar{r}p\bar{\alpha}^H}{r\bar{p}\beta^H}\right)[\alpha^{C^+} - \alpha^{C^-}] + \beta^{C^+} - \beta^{C^-}}. \quad (82)$$

The result follows. Q.E.D.

LEMMA 7. *Assume that the accuracy parameters $(\alpha^\Gamma, \beta^\Gamma)$ for $\Gamma \in \{C^+, C^-\}$ satisfy Complementarity (13)-(14). Then, it follows that*

$$\alpha^{C^+} > \alpha^{C^-} \text{ and } \beta^{C^+} < \beta^{C^-}. \quad (83)$$

Proof of Lemma 7. We prove this result by obtaining two subsets of $[0, 1] \times [0, 1]$ using Complementarity (13)-(14) that contain accuracy parameters $(\alpha^{C^+}, \beta^{C^+})$ and respectively $(\alpha^{C^-}, \beta^{C^-})$. These sets are disjoint besides the lowest sensitivity in the first set is higher than the highest sensitivity in the second set, and vice a versa for specificity.

We first define the following constants.

$$\eta \triangleq \frac{\alpha^H \bar{r} p}{\beta^H r \bar{p}} \text{ and } \tilde{\eta} \triangleq \frac{\bar{\alpha}^H \bar{r} p}{\beta^H r \bar{p}} \quad (84)$$

Note that Informativeness (1) implies that $\eta \geq 1$ and $\tilde{\eta} < 1$. Now, using these and (13)-(14), we obtain that

$$\mathcal{G}^{\text{C}^+} \triangleq \{(\alpha, \beta) \in [0, 1] \times [0, 1] : \beta + \eta\alpha > \eta \text{ and } 1 > \beta + \tilde{\eta}\alpha\} \quad (85)$$

$$\mathcal{G}^{\text{C}^-} \triangleq \{(\alpha, \beta) \in [0, 1] \times [0, 1] : \beta + \eta\alpha \leq \eta \text{ and } 1 \leq \beta + \tilde{\eta}\alpha\} \quad (86)$$

First observe that $\mathcal{G}^{\text{C}^+} \cap \mathcal{G}^{\text{C}^-} = \emptyset$. Next, we define $\hat{\mathcal{G}}^{\text{C}^+} \triangleq \{(\alpha, \beta) \in [0, 1] \times [0, 1] : \beta + \eta\alpha \geq \eta \text{ and } 1 \geq \beta + \tilde{\eta}\alpha\}$. It follows that

$$\min_{(\alpha, \beta) \in \mathcal{G}^{\text{C}^+}} \alpha > \min_{(\alpha, \beta) \in \hat{\mathcal{G}}^{\text{C}^+}} \alpha$$

because $\mathcal{G}^{\text{C}^+} \subset \hat{\mathcal{G}}^{\text{C}^+}$ and the optimal (α, β) pair minimizing α over $\hat{\mathcal{G}}^{\text{C}^+}$ is at the corner solution where both inequalities binding; and that point is not in \mathcal{G}^{C^+} since inequalities are weak. In particular, we obtain solving the system of equations that $\min_{(\alpha, \beta) \in \hat{\mathcal{G}}^{\text{C}^+}} \alpha = (\eta - 1)/(\eta - \tilde{\eta})$.

Next, we consider $\max_{(\alpha, \beta) \in \mathcal{G}^{\text{C}^-}} \alpha$. Again, the optimal solution is at the corner where both inequalities are binding, and thus it follows that $\max_{(\alpha, \beta) \in \mathcal{G}^{\text{C}^-}} \alpha = (\eta - 1)/(\eta - \tilde{\eta})$.

Combining these, we get that $\alpha^{\text{C}^+} > \alpha^{\text{C}^-}$. Following the same steps to minimize and maximize β over these sets yields $\beta^{\text{C}^+} < \beta^{\text{C}^-}$. Hence, we conclude the proof. Q.E.D.

We are now ready to prove Theorem 7.

Proof of Theorem 7. We prove this result in 3 steps. Note that the thresholds b_{C}^- and b_{C}^+ defined in Lemma 6 imply the DM's decision rule. In the first step, we provide the log-likelihood ratio process generated by the DM's decision and signal realizations. In the second, we analyze the mean of the random jumps that govern the log-likelihood ratio process. In the last step, we combine our findings to derive the limit result.

Step 1. Assume that $b_{\text{C}}^+ > b_{\text{C}}^-$.¹² Then, we have the following three regimes

- $b_{t-1} > b_{\text{C}}^+$:

$$\mathbb{P}(\Theta_t = \text{A} \mid S_t^{\text{H}} = -, S_t^{\text{M}} = +, b_{t-1}) < r \quad (87)$$

$$\mathbb{P}(\Theta_t = \text{A} \mid S_t^{\text{H}} = +, S_t^{\text{M}} = -, b_{t-1}) < r \quad (88)$$

In this case, $S_t^{\text{M}} = S_t^{\text{H}} = +$ is sufficient and necessary to act, which implies the DM decides as in type C^+ .

- $b_{\text{C}}^+ \geq b_{t-1} > b_{\text{C}}^-$:

$$\mathbb{P}(\Theta_t = \text{A} \mid S_t^{\text{H}} = -, S_t^{\text{M}} = +, b_{t-1}) \geq r \quad (89)$$

$$\mathbb{P}(\Theta_t = \text{A} \mid S_t^{\text{H}} = +, S_t^{\text{M}} = -, b_{t-1}) < r \quad (90)$$

In this case, $S_t^{\text{M}} = +$ is necessary and sufficient condition for acting.

¹² The case when $b_{\text{C}}^+ \leq b_{\text{C}}^-$ can be analyzed in the same way, and the result still follows as in the case of Theorems 3-4. Thus, we omit it to avoid repetition.

- $b_C^- \geq b_{t-1}$:

$$\mathbb{P}(\Theta_t = \mathbf{A} \mid S_t^{\mathbf{H}} = -, S_t^{\mathbf{M}} = +, b_{t-1}) \geq r \quad (91)$$

$$\mathbb{P}(\Theta_t = \mathbf{A} \mid S_t^{\mathbf{H}} = +, S_t^{\mathbf{M}} = -, b_{t-1}) \geq r \quad (92)$$

In this case, the DM decides to act after $S_t^{\mathbf{H}} = +$ or $S_t^{\mathbf{M}} = +$.

Finally, the belief update is as follows.

$$b_t = \begin{cases} \left[1 + \frac{\bar{b}_{t-1}}{b_{t-1}} \left(\frac{\alpha^{C^-}}{\alpha^{C^+}} \right)^{1_{\{\Theta_t = \mathbf{A}, S_t^{\mathbf{M}} = +, S_t^{\mathbf{H}} = +\}}} \left(\frac{\bar{\beta}^{C^-}}{\bar{\beta}^{C^+}} \right)^{1_{\{\Theta_t = \mathbf{NA}, S_t^{\mathbf{M}} = +, S_t^{\mathbf{H}} = +\}}} \right]^{-1} & \text{if } b_{t-1} > b_C^+ \\ \left[1 + \frac{\bar{b}_{t-1}}{b_{t-1}} \left(\frac{\alpha^{C^-}}{\alpha^{C^+}} \right)^{1_{\{\Theta_t = \mathbf{A}, S_t^{\mathbf{M}} = +\}}} \left(\frac{\bar{\beta}^{C^-}}{\bar{\beta}^{C^+}} \right)^{1_{\{\Theta_t = \mathbf{NA}, S_t^{\mathbf{M}} = +\}}} \right]^{-1} & \text{if } b_C^+ \geq b_{t-1} > b_C^- \\ \left[1 + \frac{\bar{b}_{t-1}}{b_{t-1}} \tilde{\zeta} \right]^{-1} & \text{if } b_C^- \geq b_{t-1} \end{cases} \quad (93)$$

where $\tilde{\zeta}$ is defined as

$$\tilde{\zeta} \triangleq \left(\frac{\alpha^{C^-}}{\alpha^{C^+}} \right)^{1_{\{\Theta_t = \mathbf{A}, S_t^{\mathbf{M}} = +\}}} \left(\frac{\bar{\alpha}^{C^-}}{\bar{\alpha}^{C^+}} \right)^{1_{\{\Theta_t = \mathbf{A}, S_t^{\mathbf{M}} = -, S_t^{\mathbf{H}} = +\}}} \left(\frac{\bar{\beta}^{C^-}}{\bar{\beta}^{C^+}} \right)^{1_{\{\Theta_t = \mathbf{NA}, S_t^{\mathbf{M}} = +\}}} \left(\frac{\beta^{C^-}}{\beta^{C^+}} \right)^{1_{\{\Theta_t = \mathbf{NA}, S_t^{\mathbf{M}} = -, S_t^{\mathbf{H}} = -\}}} \quad (94)$$

Then, the transition of the log-likelihood ratio process follows from its definition using the recursive expression of b_t . As discussed, in the proofs of Theorem 3-4, the transition rule of the log-likelihood ratio process when b_{t-1} is in $(b_C^+, 1)$ and $(0, b_C^-]$ determine the limit of L_t . Thus, in the following steps we analyze the mean of the random jumps that govern L_t when b_{t-1} lies in $(b_C^+, 1)$ and $(0, b_C^-]$.

Step 2. Assume b_t lies in $(b_C^+, 1)$. In this case, the log-likelihood ratio L_t process is governed by the random jumps $R_t^{\mathbf{HM}}$ whose mean is given by

$$\mathbb{E}^\Gamma[R_t^{\mathbf{HM}}] = p\alpha^{\mathbf{H}}\alpha^\Gamma \log \left(\frac{\alpha^{C^+}}{\alpha^{C^-}} \right) + \bar{p}\bar{\beta}^{\mathbf{H}}\bar{\beta}^\Gamma \log \left(\frac{\bar{\beta}^{C^+}}{\bar{\beta}^{C^-}} \right)$$

for $\Gamma \in \{C^+, C^-\}$. Since $\alpha^{C^+} > \alpha^{C^-}$ and $\beta^{C^+} < \beta^{C^-}$ (which implies $\bar{\beta}^{C^+} > \bar{\beta}^{C^-}$) as shown in Lemma 7, it follows that $\mathbb{E}^\Gamma[R_t^{\mathbf{HM}}] > 0$.

Assume now b_{t-1} lies in $(0, b_C^-]$. In this case, we have

$$\begin{aligned} \mathbb{E}^\Gamma[R_t^{\mathbf{HM}}] &= p\alpha^{\mathbf{H}} \underbrace{\left[\alpha^\Gamma \log \left(\frac{\alpha^{C^+}}{\alpha^{C^-}} \right) + \bar{\alpha}^\Gamma \log \left(\frac{\bar{\alpha}^{C^+}}{\bar{\alpha}^{C^-}} \right) \right]}_I + \bar{p}\bar{\beta}^{\mathbf{H}} \underbrace{\left[\bar{\beta}^\Gamma \log \left(\frac{\bar{\beta}^{C^+}}{\bar{\beta}^{C^-}} \right) + \beta^\Gamma \log \left(\frac{\beta^{C^+}}{\beta^{C^-}} \right) \right]}_{II} \\ &\quad + \underbrace{p\bar{\alpha}^{\mathbf{H}}\alpha^\Gamma \log \left(\frac{\alpha^{C^+}}{\alpha^{C^-}} \right) + \bar{p}\bar{\beta}^{\mathbf{H}}\bar{\beta}^\Gamma \log \left(\frac{\bar{\beta}^{C^+}}{\bar{\beta}^{C^-}} \right)}_{III} \end{aligned}$$

Note that the sign of terms I and II are determined by Γ because the terms inside the square brackets are KL-divergence (see proof of Theorem 1 for more). In particular, terms I and II are positive when $\Gamma = C^+$ and negative when $\Gamma = C^-$. Differently, term III is always positive regardless of Γ and p because $\alpha^{C^+} > \alpha^{C^-}$ and $\beta^{C^+} < \beta^{C^-}$ (which implies $\bar{\beta}^{C^+} > \bar{\beta}^{C^-}$) as shown in

Lemma 7. Therefore, when $\Gamma = \mathbf{C}^+$, it follows that $\mathbb{E}^\Gamma[R_t^{\text{HM}}] > 0$. However, when $\Gamma = \mathbf{C}^-$, the value of p determines the sign of $\mathbb{E}^\Gamma[R_t^{\text{HM}}]$. Observe that $\mathbb{E}^{\mathbf{C}^-}[R_t^{\text{HM}}]$ is linear function of p , and it is positive when $p = 0$ since

$$\bar{\beta}^{\mathbf{C}^-} \log\left(\frac{\bar{\beta}^{\mathbf{C}^+}}{\bar{\beta}^{\mathbf{C}^-}}\right) \geq \beta^{\mathbf{C}^-} \log\left(\frac{\beta^{\mathbf{C}^-}}{\beta^{\mathbf{C}^+}}\right) > \bar{\beta}^{\mathbf{H}} \beta^{\mathbf{C}^-} \log\left(\frac{\beta^{\mathbf{C}^-}}{\beta^{\mathbf{C}^+}}\right)$$

Here, the first inequality follows from the fact that KL-divergence is positive, and the second follows from $1 > \bar{\beta}^{\mathbf{H}}$. Hence, we define $p^{\mathbf{C}}$ as follows.

$$p^{\mathbf{C}} \triangleq \begin{cases} \frac{\text{term}_b}{\text{term}_b - \text{term}_a} & \text{if } \text{term}_b > \text{term}_a \\ 1 & \text{otherwise.} \end{cases} \quad (95)$$

where

$$\begin{aligned} \text{term}_a &= \alpha^{\mathbf{C}^-} \log\left(\frac{\alpha^{\mathbf{C}^+}}{\alpha^{\mathbf{C}^-}}\right) + \alpha^{\mathbf{H}} \bar{\alpha}^{\mathbf{C}^-} \log\left(\frac{\bar{\alpha}^{\mathbf{C}^+}}{\bar{\alpha}^{\mathbf{C}^-}}\right) \\ \text{term}_b &= \bar{\beta}^{\mathbf{C}^-} \log\left(\frac{\bar{\beta}^{\mathbf{C}^+}}{\bar{\beta}^{\mathbf{C}^-}}\right) + \bar{\beta}^{\mathbf{H}} \beta^{\mathbf{C}^-} \log\left(\frac{\beta^{\mathbf{C}^+}}{\beta^{\mathbf{C}^-}}\right) \end{aligned}$$

Thus, it follows that $\mathbb{E}^{\mathbf{C}^-}[R_t^{\text{HM}}] \geq 0$ if $p \leq p^{\mathbf{C}}$; and $\mathbb{E}^{\mathbf{C}^-}[R_t^{\text{HM}}] < 0$ if $p > p^{\mathbf{C}}$.

Step 3. We now combine our findings in the previous steps to analyze the limiting behavior of the log-likelihood ratio process. First, we consider $\Gamma = \mathbf{C}^+$. In this case, the sign of the mean $\mathbb{E}^{\mathbf{C}^+}[R_t^{\text{HM}}]$ is always positive for b_{t-1} in $(b_{\mathbf{C}}^+, 1)$ and $(0, b_{\mathbf{C}}^-]$. Thus, L_t converges to infinity and b_t converges to 1. Next, consider $\Gamma = \mathbf{C}^-$. In this case, for $p \leq p^{\mathbf{C}}$, we have the sign of the mean $\mathbb{E}^{\mathbf{C}^-}[R_t^{\text{HM}}]$ is positive for b_{t-1} in $(b_{\mathbf{C}}^+, 1)$ and nonnegative for b_{t-1} in $(0, b_{\mathbf{C}}^-]$. Thus, b_t again converges to 1 because L_t goes to infinity. If $p > p^{\mathbf{C}}$, then $\mathbb{E}^{\mathbf{C}^-}[R_t^{\text{HM}}] < 0$ for b_{t-1} in $(0, b_{\mathbf{C}}^-]$, while $\mathbb{E}^{\mathbf{C}^-}[R_t^{\text{HM}}] > 0$ for b_{t-1} in $(b_{\mathbf{C}}^+, 1)$. Therefore, b_t converges to a Bernoulli random variable.¹³ Q.E.D.

Appendix E: Relaxation of the Verification Bias

In this section of the appendix, we relax the verification bias by changing the belief updating process in two different ways. Specifically, we first consider a setup where the DM updates her belief about the machine's type based on the signal generated by the machine when the correctness of the machine's prediction is not revealed. Second, we allow the DM to observe the correctness of the machine's prediction regardless of the DM's decision to act.

First, we consider that the DM updates her belief when she decides not to act. In this case, the DM accounts for the machine prescription and her own judgment to update her belief. The modified belief updating rule that replaces equation (6) is given by

$$b_t = \begin{cases} \left[1 + \frac{\bar{b}_{t-1} \mathbb{P}^{\mathbf{W}}(S_t^{\mathbf{M}} = s^{\mathbf{M}}, S_t^{\mathbf{H}} = s^{\mathbf{H}})}{\bar{b}_{t-1} \mathbb{P}^{\mathbf{B}}(S_t^{\mathbf{M}} = s^{\mathbf{M}}, S_t^{\mathbf{H}} = s^{\mathbf{H}})} \right]^{-1} & \text{if } \mathbb{P}(\Theta_t = \mathbf{A} \mid S_t^{\mathbf{H}} = s^{\mathbf{H}}, S_t^{\mathbf{M}} = s^{\mathbf{M}}, b_{t-1}) < r \\ \left[1 + \frac{\bar{b}_{t-1} \mathbb{P}^{\mathbf{W}}(S_t^{\mathbf{M}} = s^{\mathbf{M}} \mid \Theta_t = \theta)}{\bar{b}_{t-1} \mathbb{P}^{\mathbf{B}}(S_t^{\mathbf{M}} = s^{\mathbf{M}} \mid \Theta_t = \theta)} \right]^{-1} & \text{if } \mathbb{P}(\Theta_t = \mathbf{A} \mid S_t^{\mathbf{H}} = s^{\mathbf{H}}, S_t^{\mathbf{M}} = s^{\mathbf{M}}, b_{t-1}) \geq r. \end{cases} \quad (96)$$

Thus, we directly focus on the proof when the DM updates her belief

¹³ See proofs of Theorems 3-4 for more about the relation between the sign of the mean of the random jumps and the convergence behavior.

PROPOSITION 1. *If the DM updates her belief using (96), then $b_t \xrightarrow{a.s.} 1_{\{\Gamma=B\}}$.*

Proof of Proposition 1. We prove this result by analyzing the no-interaction and no-overriding benchmarks in the first two steps. Finally, we combine our findings to conclude the proof for the main set-up.

No-Interaction Benchmark. The log-likelihood ratio process L_t becomes a random walk such that the mean of the i.i.d. random jumps is given by

$$\begin{aligned} \mathbb{E}^\Gamma[R_t] = & (p\bar{\alpha}^H\bar{\alpha}^\Gamma + \bar{p}\beta^H\beta^\Gamma) \log\left(\frac{p\bar{\alpha}^H\bar{\alpha}^B + \bar{p}\beta^H\beta^B}{p\bar{\alpha}^H\bar{\alpha}^W + \bar{p}\beta^H\beta^W}\right) + (p\bar{\alpha}^H\alpha^\Gamma + \bar{p}\beta^H\bar{\beta}^\Gamma) \log\left(\frac{p\bar{\alpha}^H\alpha^B + \bar{p}\beta^H\bar{\beta}^B}{p\bar{\alpha}^H\alpha^W + \bar{p}\beta^H\bar{\beta}^W}\right) \\ & + p \left[\alpha^H \left(\bar{\alpha}^\Gamma \log\left(\frac{\bar{\alpha}^B}{\bar{\alpha}^W}\right) + \alpha^\Gamma \log\left(\frac{\alpha^B}{\alpha^W}\right) \right) \right] + \bar{p} \left[\bar{\beta}^H \left(\bar{\beta}^\Gamma \log\left(\frac{\bar{\beta}^B}{\bar{\beta}^W}\right) + \beta^\Gamma \log\left(\frac{\beta^B}{\beta^W}\right) \right) \right] \end{aligned}$$

We next show that $\mathbb{E}^B[R_t] > 0$ and $\mathbb{E}^W[R_t] < 0$. The terms inside the square brackets above are in fact the same as in (19). In the proof of Theorem 1, we prove that the sign of those terms are aligned with the true type of the machine. Thus, we focus on the following term.

$$(p\bar{\alpha}^H\bar{\alpha}^\Gamma + \bar{p}\beta^H\beta^\Gamma) \log\left(\frac{p\bar{\alpha}^H\bar{\alpha}^B + \bar{p}\beta^H\beta^B}{p\bar{\alpha}^H\bar{\alpha}^W + \bar{p}\beta^H\beta^W}\right) + (p\bar{\alpha}^H\alpha^\Gamma + \bar{p}\beta^H\bar{\beta}^\Gamma) \log\left(\frac{p\bar{\alpha}^H\alpha^B + \bar{p}\beta^H\bar{\beta}^B}{p\bar{\alpha}^H\alpha^W + \bar{p}\beta^H\bar{\beta}^W}\right)$$

Assume first $\Gamma = B$. Then, we show that the above term is negative in the following.

$$\begin{aligned} & - (p\bar{\alpha}^H\bar{\alpha}^B + \bar{p}\beta^H\beta^B) \log\left(\frac{p\bar{\alpha}^H\bar{\alpha}^B + \bar{p}\beta^H\beta^B}{p\bar{\alpha}^H\bar{\alpha}^W + \bar{p}\beta^H\beta^W}\right) - (p\bar{\alpha}^H\alpha^B + \bar{p}\beta^H\bar{\beta}^B) \log\left(\frac{p\bar{\alpha}^H\alpha^B + \bar{p}\beta^H\bar{\beta}^B}{p\bar{\alpha}^H\alpha^W + \bar{p}\beta^H\bar{\beta}^W}\right) \\ & = (p\bar{\alpha}^H\bar{\alpha}^B + \bar{p}\beta^H\beta^B) \log\left(\frac{p\bar{\alpha}^H\bar{\alpha}^W + \bar{p}\beta^H\beta^W}{p\bar{\alpha}^H\bar{\alpha}^B + \bar{p}\beta^H\beta^B}\right) + (p\bar{\alpha}^H\alpha^B + \bar{p}\beta^H\bar{\beta}^B) \log\left(\frac{p\bar{\alpha}^H\alpha^W + \bar{p}\beta^H\bar{\beta}^W}{p\bar{\alpha}^H\alpha^B + \bar{p}\beta^H\bar{\beta}^B}\right) \\ & \leq \log[p\bar{\alpha}^H\bar{\alpha}^W + \bar{p}\beta^H\beta^W + p\bar{\alpha}^H\alpha^W + \bar{p}\beta^H\bar{\beta}^W] = \log[p\bar{\alpha}^H + \bar{p}\beta^H] < \log(1) = 0 \end{aligned}$$

Here, the first inequality follows because \log is a concave function. The second inequality follows because $p\bar{\alpha}^H + \bar{p}\beta^H$ is less than one. Thus, we obtain that $\mathbb{E}^B[R_t] > 0$, and following the same steps imply that $\mathbb{E}^W[R_t] < 0$.

No-Overriding Benchmark. In this case, the mean of the random jumps that govern the log-likelihood ratio process is given by

$$\begin{aligned} \mathbb{E}^\Gamma[R_t] = & (p\bar{\alpha}^\Gamma\alpha^H + \bar{p}\beta^\Gamma\bar{\beta}^H) \log\left(\frac{p\bar{\alpha}^B\alpha^H + \bar{p}\beta^B\bar{\beta}^H}{p\bar{\alpha}^W\alpha^H + \bar{p}\beta^W\bar{\beta}^H}\right) + (p\bar{\alpha}^\Gamma\bar{\alpha}^H + \bar{p}\beta^\Gamma\beta^H) \log\left(\frac{p\bar{\alpha}^B\bar{\alpha}^H + \bar{p}\beta^B\beta^H}{p\bar{\alpha}^W\bar{\alpha}^H + \bar{p}\beta^W\beta^H}\right) \\ & + p\alpha^\Gamma \log\left(\frac{\alpha^B}{\alpha^W}\right) + \bar{p}\bar{\beta}^\Gamma \log\left(\frac{\bar{\beta}^B}{\bar{\beta}^W}\right). \end{aligned}$$

Both of these terms are KL-divergence when $\Gamma = B$ and $-$ KL-divergence when $\Gamma = W$. Thus it is respectively positive and negative which implies that the log-likelihood ratio process converges to infinity and respectively minus infinity. Hence, we obtain $b_t \xrightarrow{a.s.} 1_{\Gamma=B}$.

Main Set-up. Note that the decision rule characterized in Lemma 1 implies that the DM follows no-interaction and no-overriding decision rules at low and respectively high b_t . As shown in the previous steps, the log-likelihood ratio process under those decision rules converge to ∞ for $\Gamma = \text{B}$, and $-\infty$ for $\Gamma = \text{W}$. Therefore, the DM correctly learns the machine's type. Q.E.D.

Finally, we consider that the DM observes the pair (S_t^M, Θ_t) in all periods. In this case, the updating rule is given by

$$b_t = \left[1 + \frac{\bar{b}_{t-1} \mathbb{P}^{\text{W}}(S_t^M | \Theta_t = \theta)}{b_{t-1} \mathbb{P}^{\text{B}}(S_t^M | \Theta_t = \theta)} \right]^{-1} \quad (97)$$

We drop the decision rule in (97) because it does not affect the belief updating there.

In fact, when the DM can observe the pair (S_t^M, Θ_t) in all periods, the sum of probability of observing all potential outcomes equals one. For those cases, the frequentist consistency of Bayesian updating (Diaconis and Freedman 1986) thus implies that the DM's belief converges almost surely to one when the machine's type is B (and respectively to zero when the machine's type is W).

E.1. Partial Relaxation of the Verification Bias

Note that the modified belief updating rule in Equation (96) is a fully Bayesian belief updating rule in the sense that the DM uses Bayes' rule to update her belief by accounting for all available information. Indeed, the frequentist consistency of Bayesian updating (Diaconis and Freedman 1986) is not violated by the censorship (not observing θ_t in case of no action) because the total probabilities of all potential observations, which are (S_t^M, S_t^H) and (S_t^M, θ_t) pairs add up to one. Hence, the belief updating rules in (6) and (96) represent two extremes such that the DM fully ignores and respectively fully accounts for the information when the correctness of the machine's prediction is not verified. In this section, we analyze the cases in between these two extremes by using the following belief update rule for $0 < \varepsilon < 1$.

For any value of $\varepsilon < 1$, the DM overlooks the information generated by unverified cases. Indeed, Equation (15) is reduced to (6) (to (96)) for $\varepsilon = 0$ (and resp., $\varepsilon = 1$). Hence, Proposition 1 and Theorems 3-4 can be thought of special cases of the following theorem.

Proof of Theorem 8. We prove this result in a way akin to the proof of Proposition 1. Thus, we first start with the benchmark cases even though the statement of the theorem is for the main set-up. Finally, we combine the results for benchmark cases to derive the results in the statement of the theorem.

No-Interaction Benchmark. The log-likelihood ratio process L_t becomes a random walk such that the mean of the i.i.d. random jumps is given by

$$\begin{aligned} \mathbb{E}^\Gamma[R_t] = & \varepsilon(p\bar{\alpha}^H\bar{\alpha}^\Gamma + \bar{p}\beta^H\beta^\Gamma) \log\left(\frac{p\bar{\alpha}^H\bar{\alpha}^B + \bar{p}\beta^H\beta^B}{p\bar{\alpha}^H\bar{\alpha}^W + \bar{p}\beta^H\beta^W}\right) + \varepsilon(p\bar{\alpha}^H\alpha^\Gamma + \bar{p}\beta^H\bar{\beta}^\Gamma) \log\left(\frac{p\bar{\alpha}^H\alpha^B + \bar{p}\beta^H\bar{\beta}^B}{p\bar{\alpha}^H\alpha^W + \bar{p}\beta^H\bar{\beta}^W}\right) \\ & + p\left[\alpha^H\left(\bar{\alpha}^\Gamma \log\left(\frac{\bar{\alpha}^B}{\bar{\alpha}^W}\right) + \alpha^\Gamma \log\left(\frac{\alpha^B}{\alpha^W}\right)\right)\right] + \bar{p}\left[\bar{\beta}^H\left(\bar{\beta}^\Gamma \log\left(\frac{\bar{\beta}^B}{\bar{\beta}^W}\right) + \beta^\Gamma \log\left(\frac{\beta^B}{\beta^W}\right)\right)\right] \end{aligned}$$

We next show that $\mathbb{E}^B[R_t] > 0$ and $\mathbb{E}^W[R_t] < 0$. The terms inside the square brackets above are, in fact, the same as in (19). In the proof of Theorem 1, we prove that the sign of those terms are aligned with the true type of the machine. Thus, we focus on the following term.

$$\varepsilon \left\{ (p\bar{\alpha}^H\bar{\alpha}^\Gamma + \bar{p}\beta^H\beta^\Gamma) \log \left(\frac{p\bar{\alpha}^H\bar{\alpha}^B + \bar{p}\beta^H\beta^B}{p\bar{\alpha}^H\bar{\alpha}^W + \bar{p}\beta^H\beta^W} \right) + (p\bar{\alpha}^H\alpha^\Gamma + \bar{p}\beta^H\bar{\beta}^\Gamma) \log \left(\frac{p\bar{\alpha}^H\alpha^B + \bar{p}\beta^H\bar{\beta}^B}{p\bar{\alpha}^H\alpha^W + \bar{p}\beta^H\bar{\beta}^W} \right) \right\}$$

In the proof of Proposition 1, we showed that the sign of the term inside the curly brackets above is aligned with the machine's type. Since $\varepsilon > 0$, it does not change it. Hence, the DM's belief converges to $1_{\{\Gamma=B\}}$ in the no-interaction benchmark when the belief updating rule is as in (15).

No-Overriding Benchmark. In this case, the mean of the random jumps that govern the log-likelihood ratio process is given by

$$\begin{aligned} \mathbb{E}^\Gamma[R_t] = & \varepsilon(p\bar{\alpha}^\Gamma\alpha^H + \bar{p}\beta^\Gamma\bar{\beta}^H) \log \left(\frac{p\bar{\alpha}^B\alpha^H + \bar{p}\beta^B\bar{\beta}^H}{p\bar{\alpha}^W\alpha^H + \bar{p}\beta^W\bar{\beta}^H} \right) + \varepsilon(p\bar{\alpha}^\Gamma\bar{\alpha}^H + \bar{p}\beta^\Gamma\beta^H) \log \left(\frac{p\bar{\alpha}^B\bar{\alpha}^H + \bar{p}\beta^B\beta^H}{p\bar{\alpha}^W\bar{\alpha}^H + \bar{p}\beta^W\beta^H} \right) \\ & + p\alpha^\Gamma \log \left(\frac{\alpha^B}{\alpha^W} \right) + \bar{p}\bar{\beta}^\Gamma \log \left(\frac{\bar{\beta}^B}{\bar{\beta}^W} \right). \end{aligned}$$

We first consider $\Gamma = B$. If $p \geq p^B$, then we have the following inequality.

$$p\alpha^B \log \left(\frac{\alpha^B}{\alpha^W} \right) + \bar{p}\bar{\beta}^B \log \left(\frac{\bar{\beta}^B}{\bar{\beta}^W} \right) \geq 0$$

Using this, we obtain the following positive lower bound on $\mathbb{E}^B[R_t]$.

$$\begin{aligned} \mathbb{E}^B[R_t] \geq & \varepsilon(p\bar{\alpha}^B\alpha^H + \bar{p}\beta^B\bar{\beta}^H) \log \left(\frac{p\bar{\alpha}^B\alpha^H + \bar{p}\beta^B\bar{\beta}^H}{p\bar{\alpha}^W\alpha^H + \bar{p}\beta^W\bar{\beta}^H} \right) + \varepsilon(p\bar{\alpha}^B\bar{\alpha}^H + \bar{p}\beta^B\beta^H) \log \left(\frac{p\bar{\alpha}^B\bar{\alpha}^H + \bar{p}\beta^B\beta^H}{p\bar{\alpha}^W\bar{\alpha}^H + \bar{p}\beta^W\beta^H} \right) \\ & + \varepsilon \left(p\alpha^B \log \left(\frac{\alpha^B}{\alpha^W} \right) + \bar{p}\bar{\beta}^B \log \left(\frac{\bar{\beta}^B}{\bar{\beta}^W} \right) \right) > 0. \end{aligned}$$

This lower bound is positive because we obtain the KL-divergence as discussed in the proof of Proposition 1 when we use ε as the common factor.

If $p < p^B$, then it follows that

$$\begin{aligned} p\alpha^B \log \left(\frac{\alpha^B}{\alpha^W} \right) + \bar{p}\bar{\beta}^B \log \left(\frac{\bar{\beta}^B}{\bar{\beta}^W} \right) & < 0, \\ (p\bar{\alpha}^B\alpha^H + \bar{p}\beta^B\bar{\beta}^H) \log \left(\frac{p\bar{\alpha}^B\alpha^H + \bar{p}\beta^B\bar{\beta}^H}{p\bar{\alpha}^W\alpha^H + \bar{p}\beta^W\bar{\beta}^H} \right) + (p\bar{\alpha}^B\bar{\alpha}^H + \bar{p}\beta^B\beta^H) \log \left(\frac{p\bar{\alpha}^B\bar{\alpha}^H + \bar{p}\beta^B\beta^H}{p\bar{\alpha}^W\bar{\alpha}^H + \bar{p}\beta^W\beta^H} \right) & > 0, \\ (p\bar{\alpha}^B\alpha^H + \bar{p}\beta^B\bar{\beta}^H) \log \left(\frac{p\bar{\alpha}^B\alpha^H + \bar{p}\beta^B\bar{\beta}^H}{p\bar{\alpha}^W\alpha^H + \bar{p}\beta^W\bar{\beta}^H} \right) + (p\bar{\alpha}^B\bar{\alpha}^H + \bar{p}\beta^B\beta^H) \log \left(\frac{p\bar{\alpha}^B\bar{\alpha}^H + \bar{p}\beta^B\beta^H}{p\bar{\alpha}^W\bar{\alpha}^H + \bar{p}\beta^W\beta^H} \right) & \\ & > p\alpha^B \log \left(\frac{\alpha^W}{\alpha^B} \right) + \bar{p}\bar{\beta}^B \log \left(\frac{\bar{\beta}^W}{\bar{\beta}^B} \right) \end{aligned}$$

Here, the first inequality is implied by the definition of p^B . The second and third inequalities follow from the first inequality and the fact that the KL-divergence for the better machine is positive.

Thus, we need to analyze how ε balances these two terms with opposing signs. In particular, a sufficiently high ε ensures that the positive term dominates the negative and hence $\mathbb{E}^B[R_t] > 0$ which implies $b_t \xrightarrow{\text{a.s.}} 1$. Similarly, a sufficiently low ε implies the opposite. These thresholds over ε depend on p . Consider the following break-even point of ε such that $\mathbb{E}^B[R_t] = 0$.

$$\varepsilon^B \triangleq \frac{p\alpha^B \log\left(\frac{\alpha^W}{\alpha^B}\right) + \bar{p}\bar{\beta}^B \log\left(\frac{\bar{\beta}^W}{\bar{\beta}^B}\right)}{(p\bar{\alpha}^B\alpha^H + \bar{p}\beta^B\bar{\beta}^H) \log\left(\frac{p\bar{\alpha}^B\alpha^H + \bar{p}\beta^B\bar{\beta}^H}{p\bar{\alpha}^W\alpha^H + \bar{p}\beta^W\bar{\beta}^H}\right) + (p\bar{\alpha}^B\bar{\alpha}^H + \bar{p}\beta^B\beta^H) \log\left(\frac{p\bar{\alpha}^B\bar{\alpha}^H + \bar{p}\beta^B\beta^H}{p\bar{\alpha}^W\bar{\alpha}^H + \bar{p}\beta^W\beta^H}\right)} \quad (98)$$

If $p < p^B$ and $\varepsilon = \varepsilon^B$, then $\mathbb{E}^B[R_t] = 0$ and hence the log-likelihood ratio process L_t and the DM's belief b_t oscillates (as in the case of Theorem 2 for $p = p^B$). On the other hand, $p < p^B$ and $\varepsilon < \varepsilon^B$ ($\varepsilon > \varepsilon^B$) imply $b_t \xrightarrow{\text{a.s.}} 0$ (and respectively $b_t \xrightarrow{\text{a.s.}} 1$). This break-even point ε^B is in $(0, 1)$ for $p < p^B$.

Next, we consider $\Gamma = W$. The same arguments with opposing signs and thresholds p^W and ε^W follow where

$$\varepsilon^W \triangleq \frac{p\alpha^W \log\left(\frac{\alpha^B}{\alpha^W}\right) + \bar{p}\bar{\beta}^W \log\left(\frac{\bar{\beta}^B}{\bar{\beta}^W}\right)}{(p\bar{\alpha}^B\alpha^H + \bar{p}\beta^B\bar{\beta}^H) \log\left(\frac{p\bar{\alpha}^B\alpha^H + \bar{p}\beta^B\bar{\beta}^H}{p\bar{\alpha}^W\alpha^H + \bar{p}\beta^W\bar{\beta}^H}\right) + (p\bar{\alpha}^B\bar{\alpha}^H + \bar{p}\beta^B\beta^H) \log\left(\frac{p\bar{\alpha}^B\bar{\alpha}^H + \bar{p}\beta^B\beta^H}{p\bar{\alpha}^W\bar{\alpha}^H + \bar{p}\beta^W\beta^H}\right)}. \quad (99)$$

The break-even point ε^W is in $(0, 1)$ for $p < p^W$.

Main Set-up. Indeed, the pairs of p and ε generate the cases that are analogous to Theorems 3-4. This is because, the DM always (regardless the value of ε) correctly learns the machine's type Γ in the no-interaction benchmark. Yet, the belief might wrongly converge to zero or one in the no-overriding benchmark as discussed in the first two steps of this proof.

Note that, the DM's decision rule in the main set-up in equation (6), and that of equation (15) are the same, specifically, $\mathbb{P}(\Theta_t = A | S_t^H = s^H, S_t^M = s^M, b_{t-1}) \geq r$. Thus, Lemma 1 characterizes the DM's decision rule again.

Overall, when the machine is better, the DM correctly learns the machine's type in the main set-up if p and ε values imply that the belief at the no-overriding benchmark converges to 1. Otherwise (more specifically $p < p^B$ and $\varepsilon \leq \varepsilon^B$), the DM's belief is recurrent and oscillates as in the proof of Theorem 3.

Similarly, when the machine is worse, the DM correctly learns the machine's type in the main set-up if p and ε values are $p \leq p^W$ or $\varepsilon \geq \varepsilon^W$ such that the belief at the no-overriding benchmark converges to zero or oscillates. Otherwise (more specifically $p > p^W$ and $\varepsilon < \varepsilon^W$), the DM's belief in the main set-up converges to a Bernoulli random variable as in the proof of Theorem 4. Q.E.D.

Recent ESMT Working Papers

	ESMT No.
Is your machine better than you? You may never know Francis de Véricourt, ESMT Berlin Huseyin Gurkan, ESMT Berlin	22-02 (R1)
Decertification in quality-management standards by incrementally and radically innovative organizations Joseph A. Clougherty, University of Illinois at Urbana-Champaign Michał Grajek, ESMT European School of Management and Technology	22-04
Do decision makers have subjective probabilities? An experimental test David Ronayne, ESMT European School of Management and Technology Roberto Veneziani, Queen Mary University of London William R. Zame, University of California at Los Angeles	22-03
Is your machine better than you? You may never know Francis de Véricourt, ESMT Berlin Huseyin Gurkan, ESMT Berlin	22-02
Mismanaging diagnostic accuracy under congestion Mirko Kremer, Frankfurt School of Finance and Management Francis de Véricourt, ESMT Berlin	22-01