

AI-Assisted Decision Making with Human Learning

Gali Noti
Cornell University

Kate Donahue
MIT

Jon Kleinberg
Cornell University

Sigal Oren
Ben-Gurion University

Abstract

AI systems are increasingly used to support human decision-making. In many cases, despite the algorithm’s superior performance, the final decision remains in human hands. For example, an AI may assist doctors in determining which diagnostic tests to run, but the doctor ultimately makes the diagnosis. Focusing on these scenarios, this paper studies AI-assisted decision-making where the human learns through repeated interactions with the algorithm. In our framework, the algorithm – designed to maximize decision accuracy according to its own model – determines which features the human can consider. The human then makes a prediction based on their own less accurate model. Additionally, we consider the possibility of a constraint on the number of features that can be taken into account.

We observe that the discrepancy between the algorithm’s model and the human’s model creates a fundamental trade-off. Should the algorithm prioritize recommending more informative features, encouraging the human to recognize their importance, even if it results in less accurate predictions in the short term until learning occurs? Or is it preferable to forgo educating the human and instead select features that align more closely with their existing understanding, minimizing the immediate cost of learning? This trade-off is shaped by both the algorithm’s patience (the time-discount rate of its objective function over multiple periods) and the human’s willingness and ability to learn.

Our results show that optimal feature selection has a surprisingly clean combinatorial characterization, reducible to a stationary sequence of feature subsets that is tractable to compute. As the algorithm becomes more patient or the human’s learning improves, the algorithm increasingly selects more informative features, enhancing both prediction accuracy and the human’s understanding of the world. Notably, early investment in learning leads to the selection of more informative features compared to a later investment. We complement our analysis by showing that the impact of errors in the algorithm’s knowledge is limited as it does not make the prediction directly.

1 Introduction

AI systems are increasingly used to assist humans in making decisions. In many situations, although the algorithm has superior performance, it can only assist the human by providing recommendations, leaving the final decision to the human. For example, algorithms assist medical doctors in assessing patients’ risk factors and in targeting health inspections and treatments ([23, 34, 36, 49]), and assist judges in making pretrial release decisions, as well as in sentencing and parole determinations ([17, 40, 41]). However, the final decision ultimately remains at the discretion of the doctors or judges [15, 18]. This paper studies such AI-assisted decision-making scenarios where the human decision-maker learns from experience in repeated interactions with the algorithm. In our setting, an algorithm assists a human by telling them which information they should use for making a prediction, with the goal of improving the *human’s* prediction accuracy.

Consider, for instance, a doctor assessing a patient’s risk of a bacterial infection. To improve their diagnosis, the doctor can order tests that reveal the values of unknown variables. However, the number of

tests the doctor can run is limited (e.g., due to costs or time constraints), and they rely on an algorithm to determine which tests to perform. The algorithm, trained on vast amounts of data, is more accurate than the doctor in estimating the statistical relationships between the unknown variables and the disease risk. The doctor has their own (potentially inaccurate) beliefs about how test results relate to risk assessment and will interpret the results according to their beliefs. The doctor orders the tests recommended by the algorithm – either because they recognize its superior data processing capabilities or because an insurance provider conditions funding on following the algorithm’s selections. The algorithm’s objective is to select the tests that lead the doctor to make the most accurate prediction.

While we phrase the problem in terms of a doctor and an algorithm, similar problems appear in other domains as well. For example, in bail decisions, a judge may use an algorithm to determine which aspects of a defendant’s record to scrutinize. In hiring, an algorithm may guide which aspects of an applicant’s file to review or suggest questions for the interview. Similarly, an algorithm may assist investors in due diligence by highlighting key aspects of a firm to review before making an investment. Beyond these domains, in which features correspond to tests of different types, a similar challenge arises in product design, e.g., when designing dashboards and deciding on a fixed subset of features to display to assist humans in various prediction tasks. The choice prioritizes features that the human can interpret correctly over those that may be more useful but are harder to interpret. For a concrete example, many weather apps display humidity which people in general know how to interpret instead of displaying the dew point, even though the latter better captures the discomfort caused by humidity [38].

We begin with a qualitative discussion, identifying two fundamental tradeoffs in AI-assisted decision-making. Then, we provide a brief summary of our model and demonstrate how both tradeoffs manifest within our framework with a concrete example. Finally, we summarize our results.

Fundamental Tradeoffs in AI-Assisted Decision Making. A key question an algorithm faces when assisting a human decision-maker is determining what information will be *useful to the human*. If the algorithm were making the prediction on its own, the answer would be clear: it would use all available information to maximize accuracy. When access to information is constrained (e.g., due to costs associated with acquiring additional observations), the algorithm would prioritize the most informative data, selecting the observations expected to contribute most to prediction accuracy. However, when the algorithm does not make the prediction on its own but instead selects information to assist a *human* in making a prediction, it must also account for the human’s ability to use that information correctly. The algorithm seeks to balance increasing *informativeness* while providing information where the human’s and the algorithm’s interpretations of the data *diverge* less. Our starting point is identifying this tradeoff, which we call the “informativeness vs. divergence” tradeoff.

Tradeoff 1. The Informativeness vs. Divergence Tradeoff: *When selecting information for human use, the algorithm faces a tradeoff between selecting the most informative data and selecting data that minimizes divergence between the human’s model and the algorithm’s model of the ground truth.*

This tradeoff implies that the algorithm may not always select the most informative features for the prediction at hand, but may instead choose less-informative features that fit the human’s level of understanding. For example, a medical algorithm might recommend that a doctor performs a less-accurate throat culture, which the doctor has used frequently and knows how to interpret, rather than a newer, more precise blood test for specific proteins that the doctor does not yet know how to interpret its results. The need to balance informativeness and divergence is fundamental to any human-algorithm interaction where the algorithm provides information to optimize human performance. Recent research reflects a growing awareness of this challenge. In chess, for instance, skill-compatible AI is designed so that a powerful algorithm playing alongside a less-skilled human does not only choose the best move

but one that the human can understand and build upon [29]. [51] considers the question of when an algorithm should delegate a decision to a human and what information it should provide. They find that more information does not always lead to better decisions.

A second fundamental question in human-algorithm interactions arises when moving beyond a single interaction. In this repeated-interaction setting, the human naturally learns from experience by repeatedly making predictions. When considering human learning, the algorithm faces a tradeoff between optimizing for short-term versus long-term outcomes: should it optimize for the human’s performance in the present, or should it guide them to learn more toward better performance in the future? Here, the dilemma of whether to provide more informative or less divergent data is amplified. If the algorithm chooses to optimize based on the human’s current understanding (e.g., selecting the less-accurate throat culture which the human already knows how to use), doing so repeatedly comes at the risk of preventing learning opportunities and sustaining incorrect beliefs, which may lead to worse performance in the long run. Conversely, if the algorithm chooses to provide information that encourages learning (e.g., instructing the doctor to use the new blood test), it can improve long-term outcomes, but at the cost of an initial learning phase during which the human may make errors while adjusting their beliefs. We call this the “fixed vs. growth mindset” tradeoff.

Tradeoff 2. *The Fixed vs. Growth Mindset Tradeoff: When a human decision-maker learns through repeated interactions with an algorithm, the algorithm faces a tradeoff between maximizing immediate performance based on the human’s current beliefs and fostering long-term performance by getting the human to learn and improve their beliefs over time.*

The fixed vs. growth mindset tradeoff naturally arises in teaching environments. The teacher can rely on the student’s current knowledge to achieve the best performance in the short term. Alternatively, the teacher can instill a new and better skill, which may take longer to master but ultimately leads to improved performance and a stronger skill set in the long run. This tradeoff is related to the classic question of “giving a fish vs. giving a fishing rod.”¹ Providing a fish ensures immediate success but does not contribute to long-term skill development. In contrast, teaching someone to fish requires an initial investment of time and effort but ultimately equips them with the ability to succeed even more in the future – whether by catching more fish or developing greater self-sufficiency.

Note that neither of these two tradeoffs has a clear “correct” answer. The fixed vs. growth mindset tradeoff (Tradeoff 2), as we will see, depends on the human’s ability to learn and the time preferences the algorithm was set to optimize for. In cases of emergency, where immediate outcomes are critical, such as helping a patient in urgent need, it may be best to optimize based on the human’s current abilities, even if they are capable of learning. By contrast, situations where time constraints are less strict present an opportunity to focus on skill development and long-term growth.

The informativeness-divergence tradeoff (Tradeoff 1), which may lead the algorithm to provide suboptimal information according to the human’s level of understanding, raises the question of whether the algorithm is merely simplifying the problem into a form the human can comprehend or actively manipulating them. While there is no formal distinction between these two cases – since in both, the algorithm adapts to human limitations at the cost of being less informative – we tend to perceive them quite differently depending on the context. For example, it seems reasonable to introduce a doctor to a new blood test (low stake scenario), but undesirable when the algorithm’s selections lead a doctor to order costly or invasive tests that are less informative than other available options (high stake scenario).

¹This is also reminiscent of the “sneakers vs. coaching” metaphor [32] for thinking about types of AI-Assistance. While “coaching” aligns closely with the growth mindset, “sneakers” differs from the fixed mindset: in the fixed mindset approach, the algorithm does not provide any additional assistance to the human, but rather optimizes solely with their existing knowledge and abilities.

Model Summary. Consider a human who needs to predict the outcome of a variable y . The true outcome is given by a function of n features: $y = f(x)$, where the features $x = (x_1, \dots, x_n)$ are independent random variables that are standardized to have zero mean and unit variance. We assume the widely used linear functional form:² $y = f(x) = c + \sum_{i=1}^n a_i x_i$.

- The human’s belief about the coefficient of feature i at time $t \geq 0$ is $h_{i,t}$, where $h_0 = (h_{1,0}, h_{2,0}, \dots, h_{n,0})$ is their initial belief. Their belief about the constant term is \bar{c} .
- The algorithm’s estimate of the true coefficients is $a' = (a'_1, \dots, a'_n)$ and its estimate of c is c' . The algorithm’s estimates of the humans’ initial belief is denoted by $h'_0 = (h'_{1,0}, h'_{2,0}, \dots, h'_{n,0})$ and \bar{c}' .
- At every time step $t \geq 0$,
 - The algorithm selects a subset of features $A_t \subseteq [n]$, with $|A_t| \leq k$, where $k \leq n$ is a budget parameter.
 - The human and the algorithm observe the realization of the features in A_t .
 - The human makes a prediction $\bar{c} + \sum_{i \in A_t} h_{i,t} x_i$ and exhibits a loss according to the Mean Squared Error (MSE) of their prediction. We denote this loss by $MSE(A_t, h_t)$.
 - The human updates $h_{i,t+1}$ for all observed features $i \in A_t$ according to some arbitrary learning rule that converges to the true a_i as the number of times that feature i is selected goes to infinity.

The algorithm selects a sequence $\mathcal{S} = (A_t)_{t=0}^\infty$ aiming to minimize the discounted loss of the human’s prediction over time: $\sum_{t=0}^\infty \delta^t MSE(A_t, h_t)$, where $\delta \in (0, 1)$ is a discounting parameter that was chosen by the entity deploying the algorithm. We also refer to δ as the “patience” parameter: the higher δ is, the more patient the algorithm is in considering future outcomes. For the majority of the paper, we analyze the case where the algorithm’s model of the ground truth is correct (i.e., $a' = a$ and $c' = c$) and the algorithm knows the human’s initial beliefs (i.e., $h_0 = h'_0$, and $\bar{c}' = \bar{c}$) and the human’s convergence rate. Given these assumptions, the algorithm has all the information required for optimizing this discounted loss. We provide more details about our model and elaborate on our modeling choices in Section 3.

As we will see, many of our results are driven by two quantities: the magnitude of a coefficient a_i representing the *informativeness* of feature i and the distance between a_i and $h_{i,0}$ representing the *divergence* for that feature at time 0.

Example. To build intuition about our model and how the fundamental tradeoffs arise within this framework, let us revisit our medical diagnosis example in a single-decision setting. Suppose that the doctor has three possible tests they can run, with true coefficients $a_1 = 0.3$, $a_2 = 0.2$, and $a_3 = 0.1$. That is, test 1 is the most informative, followed by test 2, and test 3 is the least informative. However, the doctor overestimates the importance of tests 1 and 3, with $h_1 = 0.8$ and $h_3 = 0.15$, while accurately interpreting test 2, with $h_2 = a_2 = 0.2$. Suppose that the algorithm knows both the ground truth and the human’s model and can select any subset of tests (i.e., there is no budget constraint, $k = n = 3$).

Table 1 summarizes the MSE of the human’s prediction, for each of the 2^3 possible subsets of tests that the algorithm can select. As can be seen, the best human performance (i.e., the least MSE) is achieved by selecting tests 2 and 3, and therefore this is the algorithm’s optimal feature selection in the single-decision setup. This optimal subset includes test 2, which the human perfectly understands ($h_2 = a_2$). Test 3 is also selected, despite some divergence between the human’s beliefs and the ground truth ($h_3 \neq a_3$), because the divergence is small enough relative to its informativeness to still make it beneficial. Notably,

²Note that this is similar to the linear regression that organizations (e.g., hospitals or city governments) often use to weigh factors that human experts (e.g., doctors or judges) are considering when making decisions (e.g., [2, 8, 17]).

Feature Subset	\emptyset	{1}	{2}	{3}	{1, 2}	{1, 3}	{2, 3}	{1, 2, 3}
MSE	0.14	0.3	0.1	0.1325	0.26	0.2925	0.0925	0.2525

Table 1: Mean Squared Error (MSE) for all possible subsets of features for the example in Section 1. The example has three features $\{1, 2, 3\}$, the algorithm’s model of the true coefficients is $a' = a = (0.3, 0.2, 0.1)$, the algorithm’s model of the human’s coefficients is $h' = h = (0.8, 0.2, 0.15)$, and there is no limit on the number of features that can be selected (i.e., $k = n = 3$).

the optimal subset does not include the most informative test (test 1), as the high divergence in the human’s interpretation of this test ($h_1 \gg a_1$) outweighs its informativeness. In Section 4, we analyze the exact condition under which the algorithm selects features for an optimal subset in a single prediction instance. This illustrates Tradeoff 1: when minimizing the MSE loss, the algorithm balances between high informativeness and low divergence of the selected tests.

This example also demonstrates the importance of modeling the human’s decision-making process in addition to modeling the ground truth. A naïve algorithm, which does not model human decisions but instead bases its feature selection solely on its own estimates, would recommend considering all features to minimize the error from its own perspective. In our example, this would result in almost the worst possible error, as shown in Table 1.

Now, consider the scenario in which a learning human repeatedly interacts with the algorithm to make predictions. Suppose the human is a very fast learner, such that after using a feature once, they learn its true coefficient for all subsequent predictions. If the algorithm selects all features in the first prediction, it incurs a loss of 0.2525 at that time. However, with repeated selections of all features in subsequent steps, the error drops to zero. By contrast, if the algorithm repeatedly selects only features 2 and 3 (which were optimal in the single-shot scenario), it initially incurs a smaller loss of 0.0925, but in subsequent predictions, it incurs a loss of 0.09. That is, the error improves due to learning feature 3, but only to a suboptimal result of 0.09 in each prediction. The reason for this is that the human was never given an opportunity to learn feature 1. Thus, if the algorithm equally weights each repetition, for three steps or more it is better off enduring the initial learning period and allowing the human to make mistakes and improve their model over time. The choice between these two sequences depends on how the algorithm weights short-term losses versus long-term losses. When learning is more gradual, the learning phase lasts longer and has higher costs, which, along with the weight assigned to future outcomes, influence the algorithm’s choice as well. This second example captures Tradeoff 2: the algorithm trades off the value of teaching the human (the “growth” mindset) vs. helping the human perform as best as they can with their current beliefs (the “fixed” mindset). The contrast between the algorithm’s choices in the one-shot and the learning scenarios demonstrates the importance of taking human learning into account when considering human decision-making in repeated interactions.

Results Summary. We analyze the interaction between algorithmic assistance and a learning human decision-maker. Recall that the algorithm selects a sequence of feature subsets with the objective of minimizing the discounted loss of the human’s prediction. We begin by characterizing optimal sequences of feature selections. Initially, one might suspect that feature selection is a hard problem due to the large search space: exponential in the single-interaction setting and unbounded in the repeated-interaction setting. Our analysis reveals a surprisingly clean combinatorial structure for this problem. In Theorem 5.5, we show that there exists an optimal sequence that is a stationary sequence – a sequence in which the algorithm consistently selects the same subset of features at each step (see Section 5.3). This insight allows us to restrict our focus to stationary sequences, reducing the problem to a finite space, though still

exponential in the number of features n . Then, we show in Proposition 5.7 that for a given value of δ , it is possible to compute an optimal stationary sequence in $\Theta(n \log n)$ time (see Section 5.4). Moreover, we find that across the full range of $\delta \in (0, 1)$, the total number of stationary sequences that can be optimal is at most $\Theta(n^2)$ (Proposition 5.10, Section 5.5). Notably, our analysis imposes no restrictions on human learning, except that it satisfies a natural convergence property (see Section 3.1). For the full details of our analysis, see Section 5.

Following these results, we focus our attention on optimal stationary sequences and study the conditions under which the algorithm selects more informative feature subsets, and how this choice is influenced by the time preferences in the algorithm’s objective function and the efficiency of the human’s learning. First, holding human learning at a fixed rate, we show that as the algorithm’s patience parameter δ increases, it increasingly selects more informative feature subsets (Theorem 5.11). This improves both prediction accuracy and the human’s understanding of the world in the long term. Additionally, we show that there always exists a sufficiently large δ value above which the algorithm’s optimal selection is the most informative feature set (Theorem 5.11). Second, we fix δ and vary the efficiency of the learning rule. We show that as the human learns more efficiently, the informativeness of the feature set selected by the algorithm increases (Proposition 5.15). Moreover, our analysis highlights that it is more beneficial for the human to invest in learning during earlier time steps rather than later ones, as this allows the algorithm to select more informative features and enables the human to extract greater benefits from the interaction with the algorithm. For the full analysis, see Section 5.5.

Finally, in Section 6, we study the impact of errors in the algorithm’s knowledge of the ground-truth coefficients and its models of the human’s coefficients and learning rate. Roughly speaking, we translate these modeling errors into the maximum possible error in a quantity that we later denote as the *value* of a feature, and is used to select an optimal feature set. We show that this maximum possible error quantifies a level of tolerance to algorithmic modeling errors: when the gaps between feature values are large there is a wide error tolerance range and when they are small the impact of suboptimal choices that the algorithm makes is small. This behavior results from the structure of algorithmic assistance, in which the human makes the actual predictions, and the algorithm’s role is limited to selecting the feature sets for the human to use.

2 Related Literature

Our work is situated in the literature on designing algorithms for assisting human decision-makers (e.g., [7, 24, 25, 27, 43, 50]). In particular, we consider a setting in which the algorithm selects for the human what to learn for making the best prediction as part of a repeated interaction.

The majority of the literature assume that the algorithm has direct access to information and can give the human a decision recommendation [1, 26] or display the human the relevant information for making the decisions (e.g., [20, 21]). A notable exception in regard to the algorithm’s informational structure is [33] that theoretically studies a reverse setting, complementary to ours, where humans have discretion to choose, based on situational information, which features to use at each time step, and algorithmic tools are obliged to use the same features.

Empirical papers in this area of AI-assisted human decision-making study the ability of human decision makers to correctly rely on the algorithm [1, 12, 25, 26, 50, 52]. Typically, they do not consider human learning. Two exceptions are [42] that showed experimentally the advantage of an algorithm to not always provide a recommendation and provided evidence that human decision makers learn through repeated interaction with the algorithm, and [13] that present an experimental study that applies reinforcement

learning in repeated decisions with algorithmic advice.

Our work contributes to the growing literature on designing algorithms that interact with changing human agents. In a position paper, Dean et al. [19] call for further development of “formal interaction models” between algorithms and humans who change over time. Topics in this literature include work on human-algorithm collaboration in multi-armed bandits, where the goal is to jointly identify some best arm (e.g. [9, 16]). Additionally, some works on recommendation systems take into account the fact that human preferences may change over time (e.g. [10, 11]). Tian et al. [48] studies a dynamical human-robot interaction setting, where the human’s mental model of the robot changes over time. Performative prediction [44] and learning in Stackelberg games (e.g., [28]) set up a Stackelberg game to study how to make predictions when people respond to them in a way that shifts the data distribution used for the prediction.

Also related are works on human-AI teams showing that to increase the overall performance of the team, the AI should take the human model into account and make sure its choices are understandable for the human (chess, [29], the video game *Overcooked* [14]). Moreover, it is not always the case that a more accurate algorithm [5, 6] or one that provides more information [51] are better. Our findings highlight that algorithms need to balance between using features that the human understands and features that the human needs to learn. This tension is related to the broader field of explainable or interpretable machine learning (XAI), (e.g., [22, 30, 45] and the survey [3]) that aims to explain to a human what the model has learned and develop techniques for explaining the model’s predictions [4, 46], or to learn problem representations that are more easily interpretable by the human [31, 35, 37, 39].

3 Model and Preliminaries

In this section, we provide additional details about the model introduced in Section 1. Recall that we consider a human tasked with predicting the outcome of a variable y . The true outcome is given by a linear function of a set of n features $x = \{x_1, \dots, x_n\}$, such that $y = c + \sum_{i=1}^n a_i x_i$, where the coefficients a_i are non-zero. The features x_i are independent random variables drawn from distributions F_i with known means and finite standard deviations. Without loss of generality, throughout our analysis we assume the features are standardized (such that they have zero mean and unit variance; see Appendix A). The human and the algorithm interact repeatedly as described in Section 1. We divide the next discussion into two perspectives: the human’s perspective and the algorithm’s perspective.

3.1 The Human

At each time step t , the human observes the realization of the features A_t that the algorithm selected and makes a prediction $\bar{c} + \sum_{i \in A} h_{i,t} x_i$ (predicting the mean of zero for any unobserved features). This minimizes the Mean Squared Error (MSE) from the human’s perspective. As the MSE is a function of A_t and the human’s coefficient vector h_t , we denote it by $MSE(A_t, h_t)$ and get:

Claim 3.1. $MSE(A_t, h_t) = (c - \bar{c})^2 + (\sum_{i \notin A_t} a_i^2 + \sum_{i \in A_t} (a_i - h_{i,t})^2)$

Proof.

$$\begin{aligned} \text{MSE}(A_t, h_t) &= \mathbb{E}\left[\left(c + \sum_{i=1}^n a_i x_i - \bar{c} - \sum_{i \in A_t} h_{i,t} x_i\right)^2\right] \\ &= (c - \bar{c})^2 + \mathbb{E}\left[\left(\sum_{i=1}^n a_i x_i\right)^2\right] - \mathbb{E}\left[\left(\sum_{i \in A_t} h_{i,t} x_i\right)^2\right] \end{aligned}$$

Observe that

$$\mathbb{E}\left[\left(\sum_{i \in A_t} h_{i,t} x_i\right)^2\right] = \sum_{i \in A_t} \sum_{j \in A_t} \mathbb{E}[h_{i,t} h_{j,t} x_i x_j]$$

Since features are independent we have that for $j \neq i$, $\mathbb{E}[x_i x_j] = 0$ and since the variance is normalized to 1, we have that $\mathbb{E}[x_i^2] = 1$. Hence, $\mathbb{E}\left[\left(\sum_{i \in A_t} h_{i,t} x_i\right)^2\right] = \sum_{i \in A_t} h_{i,t}^2$. Similarly, $\mathbb{E}\left[\left(\sum_{i=1}^n a_i x_i\right)^2\right] = \sum_{i=1}^n a_i^2$. Putting this together, we get that:

$$\text{MSE}(A_t, h_t) = (c - \bar{c})^2 + \left(\sum_{i \notin A_t} a_i^2 + \sum_{i \in A_t} (a_i - h_{i,t})^2\right) \quad (1)$$

as required. \square

Recall that after the human observes a realization of a feature, they learn and update their coefficient of the feature according to some learning rule. Basic properties that are natural in a learning setting include:

- Initial beliefs: At the beginning of the interaction, the human starts with some initial beliefs.
- Improvement with experience: With additional observations, the human's beliefs become closer to the true values.
- Asymptotic learning: with infinite observations, the human's beliefs converge to the truth.

Formally, in our context, $h_{i,0}$ denotes the human's initial belief about feature i and $m_i(t)$ denotes the number of times that i was selected until time t . The human's beliefs at time t , $h_{i,t}$, is a function of $m_i(t)$. The sequence $|h_{i,t} - a_i|$ is decreasing in $m_i(t)$, and $\lim_{m_i(t) \rightarrow \infty} |h_{i,t} - a_i| = 0$. Note that since the variables are independent it is reasonable to assume that the learning process is also independent for each variable and hence we make this assumption.

The following definition of ϕ -convergence of learning functions captures the above properties.

Definition 3.2. Let $\phi : \mathbb{N} \rightarrow [0, 1]$ be a monotone decreasing function with $\phi(0) = 1$, $\lim_{m \rightarrow \infty} \phi(m) = 0$. A learning dynamic is ϕ -convergent, if for every t : $(a_i - h_{i,t})^2 = \phi(m_i(t)) \cdot (a_i - h_{i,0})^2$.

In our analysis, we consider human learning in the general sense of ϕ -convergence. This abstracts away the exact mechanism that leads to learning, i.e., what feedback the human receives and how they use this feedback to update their model.

3.2 The Algorithm

The objective of the algorithm is determined by its designer who cares about minimizing the loss of the human's prediction. The designer sets a budget $0 \leq k \leq n$ on the number of features that the human

can use for prediction. The limitation to k features may arise, for example, from a cost associated with revealing features' values.

Given that we are dealing with an infinite stream of losses, it is natural to account for the timing of each loss and apply an appropriate discount. Here, we adopt the widely used exponential discounting approach, where future losses are consistently discounted according to a parameter $\delta \in (0, 1)$. Putting this together we have that the objective of the algorithm is to select a sequence of feature subsets $\mathcal{S} = (A_t)_{t=0}^\infty$ such that $|A_t| \leq k$ for every $t \geq 0$, that minimizes the discounted loss of the human:

$$\begin{aligned} L(\mathcal{S} = (A_t)_{t=0}^\infty, \phi, h_0) &= \sum_{t=0}^\infty \delta^t \text{MSE}(A_t, h_t) = \sum_{t=0}^\infty \delta^t \left((c - \bar{c})^2 + \left(\sum_{i \notin A_t} a_i^2 + \sum_{i \in A_t} (a_i - h_{i,t})^2 \right) \right) \\ &= \sum_{t=0}^\infty \delta^t \left((c - \bar{c})^2 + \left(\sum_{i \notin A_t} a_i^2 + \sum_{i \in A_t} \phi(m_i(t))(a_i - h_{i,0})^2 \right) \right) \end{aligned} \quad (2)$$

As δ approaches 1, the designer places greater emphasis on future losses, whereas smaller values of δ indicate a stronger preference for minimizing immediate losses. Thus, δ can be interpreted as a ‘‘patience’’ parameter. Note that such an infinite-horizon discounted loss can also represent situations with an uncertain interaction length, where the interaction ends in each round with probability $1 - \delta$. In this case, an interaction occurring t steps in the future has a probability of δ^t of taking place and is therefore discounted by that factor.

We are mainly interested in settings where the algorithm is, on average, more accurate than the human. Hence, we primarily analyze the case where the algorithm has accurate coefficients (i.e., $a' = a$ and $c' = c$), and accurate estimates of the human’s initial coefficients ($h'_0 = h_0$ and $\bar{c}' = \bar{c}$), and the ϕ governing their learning dynamics ($\phi' = \phi$, where ϕ' is the algorithm’s estimate of ϕ). This means that the algorithm has all the required information to choose a sequence minimizing $\sum_{t=0}^\infty \delta^t \text{MSE}(A_t, h_t)$. To simplify notations, we omit the prime notation from the algorithm’s estimates. In Section 6, we consider algorithms that have inaccurate model of the human or of the ground truth. As we will see, the structure of the problem – where the algorithm does not make the prediction directly but instead selects the features on which the human will base their prediction – limits the impact of errors in modeling both the human and the ground truth.

Note that the algorithm is not obligated to exhaust the budget of k features. When the algorithm chooses not to select any features at all (i.e., $A = \emptyset$), the human still needs to make a prediction. In this case, since no additional information is available, the human makes the same prediction of $\hat{y}_h = \bar{c}$ for every instance of the problem, which represents their belief about the average of the predicted value.

4 Warm-up: Feature Selection with Fixed Human Beliefs

We start by characterizing the algorithm’s optimal selection of features in the static case where the human holds fixed beliefs h about feature coefficients. In this case, according to Claim 3.1 in Section 3, the algorithm will choose the subset of features minimizing:

$$\text{MSE}(A, h) = (c - \bar{c})^2 + \left(\sum_{i \notin A} a_i^2 + \sum_{i \in A} (a_i - h_i)^2 \right)$$

With only one feature, by Claim 3.1 we have: $\text{MSE}(\{1\}, h) = (c - \bar{c})^2 + (a_1 - h_1)^2$, whereas if we do not use this feature, $\text{MSE}(\emptyset, h) = (c - \bar{c})^2 + a_1^2$. Thus, selecting the feature reduces error whenever $a_1^2 > (a_1 - h_1)^2$,

which holds if $h_1^2 < 2a_1h_1$, or equivalently, $h_1 \in (0, 2a_1)$. Note that if the human and algorithm agree on the interpretation of the feature ($a_1 = h_1$), then this is always satisfied, and so the algorithm would always select the feature. If a_1 and h_1 have different signs (the human and algorithm completely disagree), then this is never satisfied. If they do have the same sign, the feature is useful only if the human’s coefficient h_1 is not too large compared to the algorithm’s coefficient a_1 . When h_1 is too large, it means that the human “overshoots,” i.e., overuses the feature and ends up too far on the other side of the truth. The factor of 2 limits overshooting to keep the feature beneficial to the human.

In the single feature case, we saw that the choice of whether to select the feature or not depends on the value of $MSE(\{1\}, h) - MSE(\emptyset, h) = a_1^2 - (a_1 - h_1)^2$. We find that this is a useful quantity for choosing which features to select in the general case, and refer to it as the *value* of feature 1. As it turns out, this quantity greatly simplifies the problem of computing an optimal subset in the general case. More generally,

Definition 4.1. *The value of a set A of features $V(A, h) = MSE(\emptyset, h) - MSE(A, h)$ is the improvement in loss from using the features in A compared to not using any feature.*

Since for a given problem instance $MSE(\emptyset, h)$ is fixed, to minimize the objective function, the algorithm should select a set A^* of up to k features with the maximum value of $V(A^*, h)$.

Lemma 4.2. *The value of a set of features $A \subseteq [n]$ satisfies $V(A, h) = \sum_{i \in A} V(\{i\}, h)$.*

Proof.

$$\begin{aligned}
V(A, h) &= MSE(\emptyset, h) - MSE(A, h) \\
&= (c - \bar{c})^2 + \sum_{i=1}^n a_i^2 - ((c - \bar{c})^2 + (\sum_{i \notin A} a_i^2 + \sum_{i \in A} (a_i - h_i)^2)) \\
&= \sum_{i \in A} a_i^2 - \sum_{i \in A} (a_i - h_i)^2 = \sum_{i \in A} 2a_i h_i - h_i^2 \\
&= \sum_{i \in A} V(\{i\}, h)
\end{aligned} \tag{3}$$

□

It is instructive to take a more careful look at Equation (3) that the algorithm aims to maximize. The contribution of each feature $i \in A$ to the value of a set A is composed of two parts: (1) a_i^2 , which captures the informativeness³ of the feature; (2) minus $(a_i - h_{i,t})^2$, which captures the divergence between the human’s and the algorithm’s beliefs about the feature. The tension between these two terms is at the heart of this paper.

Lemma 4.2 implies the following corollary for multiple features.

Corollary 4.3. *The algorithm only selects features i with a positive value (i.e., $h_i^2 < 2a_i h_i$).*

Lemma 4.2 also gives rise to a simple and efficient algorithm to compute an optimal subset of features that the algorithm should select. We first compute the value for each feature $i \in [n]$, as $V(\{i\}, h) = 2a_i h_i - h_i^2$. Then, we sort features by their values and add features to the feature selection set by descending order until reaching the budget of k features or there are no more features with positive values. Thus, we establish that:

Proposition 4.4. *An optimal feature selection $A^* \subseteq [n]$ can be computed in $n \log(n)$ time.*

³The coefficient of a feature captures the feature’s importance and variance; recall that features are standardized, and thus coefficients are scaled by the standard deviation, see Lemma A.1 in Appendix A.

5 Feature Selection with Human Learning

In this section, we turn to discussing algorithmic assistance for a human decision maker who learns and updates their beliefs through repeated interactions with the algorithm. In Section 1, we saw the fundamental tension between informativeness and divergence, which leads to the phenomenon where the algorithm does not necessarily select the most informative features for prediction but instead it may select those features that the human can effectively use. As discussed, when the human’s beliefs are not fixed but are instead updated through repeated interaction with the algorithm, selecting a less informative set of features – while may be helpful in the short term before learning occurs – may limit learning opportunities, and as a result sustain human’s incorrect beliefs and potentially lead to worse performance in the long run. We will see that in our framework, this tradeoff is governed by two parameters: δ that captures how much weight the algorithm’s designer puts on short term losses versus long terms losses; and ϕ that captures how efficient the human’s learning is.

5.1 The Value of a Sequence

For the static model (Section 4) it was very useful to define the value of using a subset of features (Definition 4.1). We can extend this definition for a sequence of feature subsets in our learning setting:

Definition 5.1. *The value of a sequence of feature subsets $\mathcal{S} = (A_t)_{t=0}^\infty$ for a given value of δ and a learning dynamic that is ϕ -convergent is the improvement in discounted loss from selecting features according to \mathcal{S} compared to not using any feature at all. That is,*

$$V_{\delta,\phi}(\mathcal{S}) = L((\emptyset)_{t=0}^\infty, \phi, h_0) - L(\mathcal{S}, \phi, h_0)$$

As in the static case, since the loss of not using any feature, $L((\emptyset)_{t=0}^\infty, \phi, h_0)$, is constant, an optimal feature selection is a sequence $\mathcal{S}^* = (A_t)_{t=0}^\infty$ with maximum value of $V_{\delta,\phi}(\mathcal{S}^*)$ that obeys the budget constraint (i.e., $A_t \subseteq [n]$, $|A_t| \leq k$ for all t). We derive an explicit expression for $V_{\delta,\phi}(\mathcal{S})$:

Claim 5.2. $V_{\delta,\phi}(\mathcal{S} = (A_t)_{t=0}^\infty) = \sum_{t=0}^\infty \delta^t \sum_{i \in A_t} (a_i^2 - \phi(m_i(t))(a_i - h_{i,0})^2)$

Proof. Using the value notion for the non-learning setting (Definition 4.1) and Lemma 4.2,

$$\begin{aligned} V_{\delta,\phi}(\mathcal{S}) &= \sum_{t=0}^\infty \delta^t V(A_t, h_t) = \sum_{t=0}^\infty \delta^t \sum_{i \in A_t} V(\{i\}, h_t) \\ &= \sum_{t=0}^\infty \delta^t \sum_{i \in A_t} (a_i^2 - \phi(m_i(t))(a_i - h_{i,0})^2) \end{aligned} \tag{4}$$

Recall that $m_i(t)$ is the number of times that feature i was selected prior to time step t . □

Before we delve into our general results, in the next subsection we illustrate the learning setting in a simple setup with a specific learning dynamic that is ϕ -convergent. The example demonstrates the basic tradeoff between the fixed and growth mindsets (Tradeoff 2) in the algorithm’s choice of feature selections. We then continue with the general analysis of the implications of considering a learning decision maker on the algorithm’s optimal sequence of feature selections.

5.2 Example: Exponential Learning Dynamics

We consider a ϕ -convergent learning rule given by the following simple dynamic:

$$h_{i,t+1} = f(h_{i,t}, a, A_t) = \begin{cases} w \cdot h_{i,t} + (1 - w) \cdot a_i, & \text{if } i \in A_t, \\ h_{i,t}, & \text{otherwise.} \end{cases} \quad (5)$$

where $w \in [0, 1]$ represents a learning rate parameter. Large values of w indicate a human who is a slow learner, while smaller values of w correspond to a faster learner. Equivalently, we can write a cumulative version of this step-by-step learning rule, where the human’s belief about feature i at time t is given by:

$$h_{i,t} = w^{m_i} \cdot h_{i,0} + (1 - w^{m_i}) \cdot a_i \quad (6)$$

and $m_i = m_i(t)$ denotes the number of times feature i has been selected up to time t . At $t = 0$, the human starts with their initial beliefs $h_{i,0}$ about each feature i . As $m_i \rightarrow \infty$, the human’s coefficient of i exponentially converges to the true value a_i .

This learning dynamic is ϕ -convergent with an exponential convergence rate of $\phi(m_i) = w^{2m_i}$. To see why, substitute Equation (6) and get that:

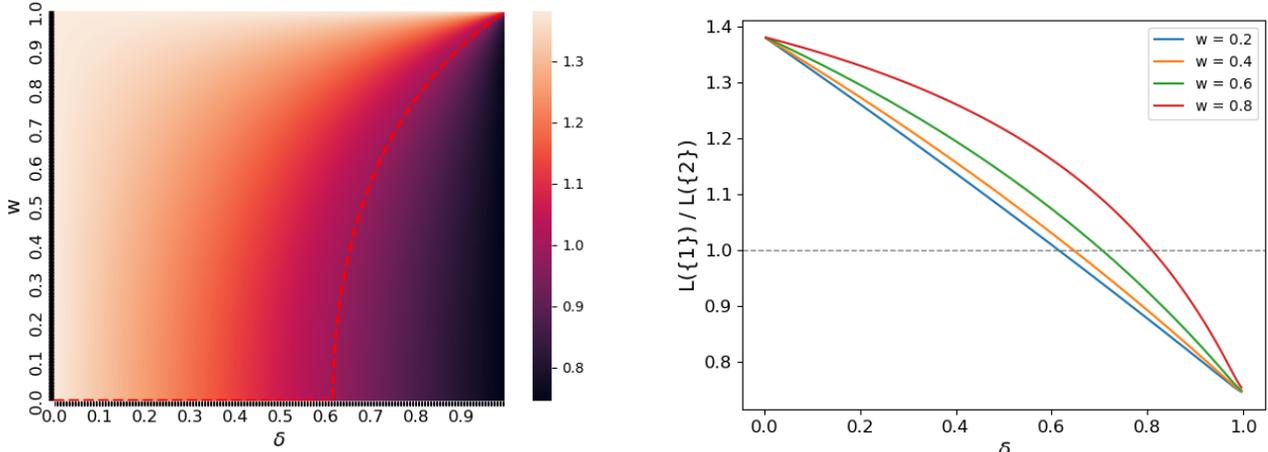
$$(a_i - h_{i,t})^2 = (a_i - (w^{m_i} \cdot h_{i,0} + (1 - w^{m_i})a_i))^2 = (w^{m_i}(a_i - h_{i,0}))^2 = w^{2m_i}(a_i - h_{i,0})^2$$

With this learning rule, consider the following toy example. Suppose there are two features and the algorithm can select only one of them at each step. Feature 1 is more informative than feature 2, with true coefficients $a_1 = 1$, $a_2 = 0.4$. Initially, the human has a larger error on the more informative feature and a smaller one on the less informative feature, with $h_{1,0} = -0.5$, $h_{2,0} = 0.75$.

We compare two simple sequences: one in which the algorithm repeatedly selects feature 1 (the “informative” sequence) and another in which it repeatedly selects feature 2 (the “non-divergent” sequence). We ask, when does the informative sequence have lower discounted loss than the non-divergent sequence? The answer depends on the discount factor δ and the learning parameter w . When δ is very small, the loss is dominated by the first time step, making the non-divergent sequence the better choice, as it provides immediate value. In this case, learning is essentially ignored since its gains occur at a later stage and are heavily discounted. On the other end of the spectrum, when δ is close to 1, the algorithm assigns significant weight to later steps where learning plays an essential role. In this case, using the informative sequence is preferable, as it provides greater value once learning has converged. Thus, we expect a transition point – an intermediate value of δ – above which the informative sequence is preferred and below which the non-divergent sequence is preferred. For higher w values, which indicate slower learning, larger δ values are required to make the informative sequence preferable.

Figure 1a shows a heatmap of the ratio of losses between selecting the informative sequence (feature 1) and selecting the non-divergent sequence (feature 2), over the δ, w plane. The red dashed line indicates a loss ratio of 1, which is the empirical transition curve; for δ values to the right of this curve, selecting the informative sequence yields lower loss, while for values to the left of the curve, the non-divergent sequence performs better.

Specifically, we can see that for every value of $w < 1$ there is a transition point $\delta^*(w)$ above which the algorithm prefers selecting the informative feature. As w increases (slower learning), $\delta^*(w)$ is increasing. The $w = 0$ line corresponds to a scenario where the human learns the coefficient of the feature they observed in one step. In this example, $\delta^*(w = 0) \approx 0.6$. The empirical curve in Figure 1a coincides with our theoretical results characterizing such curves (see Figure 2 in Appendix B). Figure 1b shows



(a) A heat map of the δ, w plane

(b) Fixing specific values of w

Figure 1: The ratio of discounted losses of the more informative sequence over the less informative sequence for the example in Section 5.2.

projections for several w values to further illustrate the transition. For slower learning (higher w) lines intersect the ratio of 1 (dashed gray line) at higher δ values. At the extreme values of δ (either close to 0 or 1), we see that the different w values give the same loss ratio; this is because, as explained above, for δ close to 0, the loss is dominated by the first step, which is independent of w , and for δ close to 1 the long-term effects of learning dominate, diminishing the influence of finite learning times.

5.3 Optimality of Stationary Sequences

In this section, we analyze optimal sequences of feature selections in the setting of human learning. We assume that the true coefficients a are distinct.⁴ We begin by defining stationary sequences:

Definition 5.3. A sequence of feature subsets $\mathcal{S} = (A_t)_{t=0}^{\infty}$ is stationary if $A_t = A \subseteq [n]$ for all t . We use the notation $(A)_{t=0}^{\infty}$ to denote a stationary sequence that always selects a subset A .

We are particularly interested in stationary sequences due to their simplicity. We show that for general $n > 1$ and $k \leq n$, if the learning dynamic is ϕ -convergent, then there exists an optimal sequence that is stationary. This fact dramatically decreases the size of the optimization space and the computational complexity of the problem; Later in Proposition 5.7 we show that we can find an optimal stationary sequence in $n \log(n)$ time.

The proof builds on the following lemma showing that if the learning dynamic is ϕ -convergent, then there exists an optimal sequence that has an infinite suffix that starting at time step t always selects the same subset (i.e., the sequence has an infinite stationary suffix).

Lemma 5.4. If the learning dynamic is ϕ -convergent, then for any sequence \mathcal{S} that does not end in an infinite stationary suffix, there exists a sequence \mathcal{S}' of higher value with a stationary suffix.

Proof. Consider a sequence \mathcal{S} that does not end in an infinite stationary suffix. Denote by $R(\mathcal{S})$ the set of features selected infinitely many times. Observe that since we consider ϕ -convergent learning dynamics,

⁴This assumption simplifies the analysis and is essentially without loss of generality, as we can always introduce a small amount of noise to ensure this property holds without significantly changing the instance.

there exists time step T after which all features in $R(\mathcal{S})$ have been selected sufficiently many times for the coefficient of each such feature i to have converged to a_i . This implies that their contribution to the accuracy of the prediction is positive. Hence, if $R(\mathcal{S}) \leq k$, the sequence that ends in the stationary suffix that always selects $R(\mathcal{S})$ has a higher value than the original sequence. Similarly, if $R(\mathcal{S}) > k$ and there are some steps in which \mathcal{S} selects less than k features, we can construct a sequence of a higher value by making sure we select k features from $R(\mathcal{S})$ in each step of the suffix.

We now handle the case that $|R(\mathcal{S})| > k$ and there exists a suffix in which \mathcal{S} selects k features at each step. We will show that there exists another sequence \mathcal{S}' that has a higher value and $|R(\mathcal{S}')| = |R(\mathcal{S})| - 1$. By applying this argument $|R(\mathcal{S})| - k$ times we get that there is a sequence \mathcal{S}^* with a higher value than \mathcal{S} such that $R(\mathcal{S}^*) = k$. This implies that a sequence with a stationary suffix that always selects $R(\mathcal{S}^*)$ has a value at least as high as the value of \mathcal{S}^* .

Let A^* denote a set that includes the k features in $R(\mathcal{S})$ that have the highest values of a_i . Note that since the features in A^* have the highest true coefficients, for any feature $i \in A^*$ and $j \in R(\mathcal{S}) \setminus A^*$, there exists $\varepsilon > 0$ such that

$$a_i^2 - \varepsilon(a_i - h_{i,0})^2 > a_j^2 \quad (7)$$

Let $\varepsilon > 0$ be such that the inequality above holds for any pair $i \in A^*$ and $j \in R(\mathcal{S}) \setminus A^*$. Now, since the learning dynamic is convergent, there exists m_ε such that for any $m \geq m_\varepsilon$, $\phi(m) \leq \varepsilon$. Let T denote the time step by which, for all of the features $i \in A^*$ we have that $m_i \geq m_\varepsilon$. If $R(\mathcal{S}) > k$, then there exists some feature $j \in R \setminus A^*$ that was selected infinitely many times. Since $|A^*| = k$, at each time step $t > T$ that j was chosen, there exists a feature $i_t \in A^*$ that was not chosen. We create the new sequence \mathcal{S}' , by replacing each selection of feature j in step $t > T$, by feature $i_t \in A^*$. By construction $|R(\mathcal{S}')| = |R(\mathcal{S})| - 1$ and by inequality (7), the value of the sequence \mathcal{S}' is greater than the value of sequence \mathcal{S} . \square

Theorem 5.5. *If the learning dynamic is ϕ -convergent, then there exists an optimal sequence and any optimal sequence is stationary (i.e., always selects the same subset of features).*

Proof. By Lemma 5.4, we have that if an optimal sequence exists, it must end with a stationary suffix. Consider a sequence $\mathcal{S} = (A_t)_{t=0}^\infty$ that its stationary suffix begins at $T + 1$: $A_T \neq A_{T+1}$ and for all $t \geq T + 1$, $A_t = A_{T+1}$. For ease of notation, let $A = A_{T+1}$ and $B = A_T$. We will show that there is another sequence \mathcal{S}' that its stationary suffix begins at $T' \leq T$ and has a higher value. The fact that for every sequence there exists a stationary sequence attaining a higher value implies that there exists an optimal sequence. This is because the number of stationary sequences is bounded. Hence an optimal stationary sequence is also a globally optimal sequence.

We will first compare between three sequences starting from step T (all with the same prefix A_0, \dots, A_{T-1}).

1. \mathcal{S}_A - a sequence that selects A starting at T .
2. $\mathcal{S}_{B \rightarrow A} = S$ - a sequence that selects B at time T and then selects A for all subsequent time steps.
3. $\mathcal{S}_{B,B \rightarrow A}$ - a sequence that selects B at time steps T and $T + 1$ and then selects A for all subsequent time steps.

Lemma 5.6. $V_{\delta,\phi}(\mathcal{S}_{B \rightarrow A}) < V_{\delta,\phi}(\mathcal{S}_A)$ or $V_{\delta,\phi}(\mathcal{S}_{B \rightarrow A}) \leq V_{\delta,\phi}(\mathcal{S}_{B,B \rightarrow A})$.

Proof. Assume towards contradiction that $V_{\delta,\phi}(\mathcal{S}_{B \rightarrow A}) \geq V_{\delta,\phi}(\mathcal{S}_A)$ and $V_{\delta,\phi}(\mathcal{S}_{B \rightarrow A}) > V_{\delta,\phi}(\mathcal{S}_{B,B \rightarrow A})$.

Let $t_i = m_i(T)$ denote the number of times that feature i was selected prior to time step T . We compute the expected discounted values of the sequences we presented:

$$\begin{aligned}
V_{\delta,\phi}(S_A) &= \sum_{t=T}^{\infty} \delta^t \cdot \sum_{i \in A} (a_i^2 - \phi(t_i + t - T) \cdot (a_i - h_{i,0})^2) \\
V_{\delta,\phi}(S_{B \rightarrow A}) &= \delta^T \sum_{j \in B} (a_j^2 - \phi(t_j)(a_j - h_{j,0})^2) \\
&\quad + \sum_{t=T+1}^{\infty} \delta^t \cdot \left(\sum_{i \in A \setminus B} (a_i^2 - \phi(t_i + t - T - 1) \cdot (a_i - h_{i,0})^2) + \sum_{i \in A \cap B} (a_i^2 - \phi(t_i + t - T) \cdot (a_i - h_{i,0})^2) \right) \\
&= \delta^T \sum_{j \in B} (a_j^2 - \phi(t_j)(a_j - h_{j,0})^2) \\
&\quad + \sum_{t=T+1}^{\infty} \delta^t \cdot \left(\sum_{i \in A} a_i^2 - \sum_{i \in A \setminus B} (\phi(t_i + t - T - 1) \cdot (a_i - h_{i,0})^2) - \sum_{i \in A \cap B} (\phi(t_i + t - T) \cdot (a_i - h_{i,0})^2) \right) \\
V_{\delta,\phi}(S_{B,B \rightarrow A}) &= \delta^T \sum_{j \in B} (a_j^2 - \phi(t_j)(a_j - h_{j,0})^2) + \delta^{T+1} \sum_{j \in B} (a_j^2 - \phi(t_j + 1) \cdot (a_j - h_{j,0})^2) \\
&\quad + \sum_{t=T+2}^{\infty} \delta^t \cdot \left(\sum_{i \in A} a_i^2 - \sum_{i \in A \setminus B} (\phi(t_i + t - 2 - T) \cdot (a_i - h_{i,0})^2) - \sum_{i \in A \cap B} (\phi(t_i + t - T) \cdot (a_i - h_{i,0})^2) \right)
\end{aligned}$$

By our assumption $V_{\delta,\phi}(S_{B \rightarrow A}) \geq V_{\delta,\phi}(S_A)$, and hence:

$$\begin{aligned}
&\delta^T \left(\sum_{j \in B \setminus A} (a_j^2 - \phi(t_j)(a_j - h_{j,0})^2) - \sum_{i \in A \setminus B} (a_i^2 - (\phi(t_i))(a_i - h_{i,0})^2) \right) \\
&\geq \sum_{t=T+1}^{\infty} \delta^t \sum_{i \in A \setminus B} (\phi(t_i + t - 1 - T) - \phi(t_i + t - T)) \cdot (a_i - h_{i,0})^2
\end{aligned}$$

If we divide by δ^T and adjust the indices accordingly, we get that:

$$\begin{aligned}
&\sum_{j \in B \setminus A} (a_j^2 - \phi(t_j)(a_j - h_{j,0})^2) - \sum_{i \in A \setminus B} (a_i^2 - (\phi(t_i))(a_i - h_{i,0})^2) \\
&\geq \sum_{t=1}^{\infty} \delta^t \sum_{i \in A \setminus B} (\phi(t_i + t - 1) - \phi(t_i + t)) \cdot (a_i - h_{i,0})^2
\end{aligned} \tag{8}$$

Also by our assumption $V_{\delta,\phi}(S_{B,B \rightarrow A}) < V_{\delta,\phi}(S_{B \rightarrow A})$, we have that:

$$\begin{aligned}
&\delta^{T+1} \left(\sum_{j \in B \setminus A} (a_j^2 - \phi(t_j + 1)(a_j - h_{j,0})^2) - \sum_{j \in A \setminus B} (a_j^2 - \phi(t_j)(a_j - h_{j,0})^2) \right) \\
&< \sum_{t=T+2}^{\infty} \delta^t \sum_{i \in A \setminus B} (\phi(t_i + t - 2 - T) - \phi(t_i + t - 1 - T)) \cdot (a_i - h_{i,0})^2
\end{aligned}$$

If we divide by δ^{T+1} and adjust the indices accordingly, we get that:

$$\begin{aligned} & \sum_{j \in B \setminus A} (a_j^2 - \phi(t_j + 1)(a_j - h_{j,0})^2) - \sum_{j \in A \setminus B} (a_j^2 - \phi(t_j)(a_j - h_{j,0})^2) \\ & < \sum_{t=1}^{\infty} \delta^t \sum_{i \in A \setminus B} (\phi(t_i + t - 1) - \phi(t_i + t)) \cdot (a_i - h_{i,0})^2 \end{aligned} \quad (9)$$

If we add inequalities (8) and (9) and do some rearranging, we get that:

$$\sum_{j \in B \setminus A} (a_j^2 - \phi(t_j)(a_j - h_{j,0})^2) > \sum_{j \in B \setminus A} (a_j^2 - \phi(t_j + 1)(a_j - h_{j,0})^2)$$

By definition $\phi(\cdot)$ is decreasing and hence $(a_j^2 - \phi(t_j + 1)(a_j - h_{j,0})^2) \geq (a_j^2 - \phi(t_j)(a_j - h_{j,0})^2)$, in contradiction to the above inequality. \square

Continuing the proof of Theorem 5.5. Note that \mathcal{S}_A is a sequence that its infinite suffix starts at T or earlier (if, $A_{T-1} = A$). Therefore, if $V_{\delta, \phi}(\mathcal{S}_{B \rightarrow A}) < V_{\delta, \phi}(\mathcal{S}_A)$, we are done. Else, $V_{\delta, \phi}(\mathcal{S}_{B \rightarrow A}) \leq V_{\delta, \phi}(\mathcal{S}_{B, B \rightarrow A})$. Observe that Lemma 5.6 can be also applied to the sequence of A 's that starts at $T + 2$. That is,

$$V_{\delta, \phi}(\mathcal{S}_{B, B \rightarrow A}) < V_{\delta, \phi}(\mathcal{S}_{B \rightarrow A}) \text{ or } V_{\delta, \phi}(\mathcal{S}_{B, B \rightarrow A}) \leq V_{\delta, \phi}(\mathcal{S}_{B, B, B \rightarrow A}).$$

Thus, $V_{\delta, \phi}(\mathcal{S}_{B \rightarrow A}) \leq V_{\delta, \phi}(\mathcal{S}_{B, B \rightarrow A})$, implies that $V_{\delta, \phi}(\mathcal{S}_{B, B \rightarrow A}) \leq V_{\delta, \phi}(\mathcal{S}_{B, B, B \rightarrow A})$. Repeating the same argument recursively, we can always improve the original sequence by increasing the length of the B sequence. In particular, for any T_B the sequence that selects B , T_B times and then selects A for all subsequent time steps, has at least the same value as the sequence that selects B only $T_B - 1$ times. We denote each such sequence by $\mathcal{S}_{T_B} = \mathcal{X}, B, \dots, B, (A)_{T+T_B}^{\infty}$, where \mathcal{X} is the prefix. We will show that here exists T' such that for any $T_B \geq T'$ there exists a stationary sequence starting at $T + 1$ that has a higher value than \mathcal{S}_{T_B} .

Let $C = A \cap B$ denote the features that appear both in A and in B . Let A^* denote an optimal subset of size at most $k - |C|$ for the following optimization problem.

$$\sum_{t=T+T_B}^{\infty} \delta^t \sum_{i \in A^*} (a_i^2 - \phi(m_i(t)) \cdot (a_i - h_{i,0})^2)$$

Implying that the selected features are the $k - |C|$ features with highest positive values of:⁵

$$U_i(T + T_B, \infty) = \frac{1}{1 - \delta} a_i^2 - \sum_{t=0}^{\infty} \delta^t \phi(m_i(t + T + T_B)) \cdot (a_i - h_{i,0})^2$$

Similarly, we can define a value for feature i and part of the sequence between time step T and $T + T_B$. An optimal stationary subsequence for steps T through $T + T_B$, will choose the $k - |C|$ features with highest values of:

$$U_i(T, T + T_B - 1) = \sum_{t=0}^{T_B-1} \delta^t a_i^2 - \sum_{t=0}^{T_B-1} \delta^t \phi(m_i(t + T)) \cdot (a_i - h_{i,0})^2$$

⁵If some of the values are negative then less than $k - |C|$ features will be selected.

Observe that by construction we only consider features $i \notin C$. Hence, for every such feature we have that $m_i(t+T) = m_i(t+T+T_B)$. Now, since $U_i(T+T_B, \infty)$ is convergent, this means that for every ε there exists T_B such that $U_i(T, T+T_B-1) \geq U_i(T+T_B, \infty) - \varepsilon$. This implies that there exist T_∞ such that for $T_B \geq T_\infty$, the values of $U_i(T, T+T_B-1)$ for any feature $i \notin C$ will be sufficiently close to $U_i(T+T_B, \infty)$ and hence the same subset will be optimal for a stationary sequence starting at $T+T_B$ and a stationary sequence between T and T_B . □

It is interesting to note that the proof of Theorem 5.5 does not require ϕ to be decreasing. In fact, as long as ϕ converges to 0, there exists an optimal sequence that is stationary even if ϕ oscillates.

5.4 Computing Optimal Stationary Sequences

Recall that in Claim 5.2, we showed that an optimal sequence for the learning setting maximizes:

$$V_{\delta, \phi}(\mathcal{S} = (A_t)_{t=0}^\infty) = \sum_{t=0}^{\infty} \delta^t \sum_{i \in A_t} (a_i^2 - \phi(m_i(t))(a_i - h_{i,0})^2)$$

Thus, an optimal stationary sequence chooses a subset $A^* \subseteq [n]$ of at most k features that maximizes:

$$\begin{aligned} V_{\delta, \phi}((A^*)_{t=0}^\infty) &= \sum_{t=0}^{\infty} \delta^t \sum_{i \in A^*} (a_i^2 - \phi(t)(a_i - h_{i,0})^2) \\ &= \sum_{i \in A^*} \sum_{t=0}^{\infty} \delta^t (a_i^2 - \phi(t)(a_i - h_{i,0})^2) \\ &= \sum_{i \in A^*} \left(\frac{1}{1-\delta} a_i^2 - \sum_{t=0}^{\infty} \delta^t \phi(t)(a_i - h_{i,0})^2 \right) \\ &= \sum_{i \in A^*} V_{\delta, \phi}(\{i\}_{t=0}^\infty) \end{aligned} \tag{10}$$

Hence, A^* includes the top k or less features where $V_{\delta, \phi}(\{i\}_{t=0}^\infty) > 0$. To verify that $V_{\delta, \phi}(\{i\}_{t=0}^\infty) = \frac{1}{1-\delta} a_i^2 - \sum_{t=0}^{\infty} \delta^t \phi(t)(a_i - h_{i,0})^2$ is well-defined, observe that $\sum_{t=0}^{\infty} \delta^t \phi(t)$ converges. The reason for this is that since $\phi(t) \leq 1$, we have $\delta^t \phi(t) \leq \delta^t$. Then, by the comparison test, since $\sum_{t=0}^{\infty} \delta^t$ converges for $\delta \in (0, 1)$, it follows that $\sum_{t=0}^{\infty} \delta^t \phi(t)$ also converges. In this paper we make the reasonable assumption that $V_{\delta, \phi}(\{i\}_{t=0}^\infty)$ can be computed in $o(1)$.

As an example, consider the feature values for the exponential learning rule from Section 5.2, which is ϕ -convergent for $\phi(t) = w^{2t}$. For this learning rule, we obtain $\sum_{t=0}^{\infty} \delta^t w^{2t} = \frac{1}{1-\delta \cdot w^2}$ and thus the value of feature i is $V_{\delta, \phi}(\{i\}_{t=0}^\infty) = \frac{1}{1-\delta} a_i^2 - \frac{1}{1-\delta \cdot w^2} (a_i - h_{i,0})^2$.

As in the static case, this characterization of optimal subsets leads to a simple method for finding an optimal stationary sequence S^* in $O(n \log n)$ time, similar to the static case in Proposition 4.4: (1) Compute $V_{\delta, \phi}(\{i\}_{t=0}^\infty)$ for each feature i , (2) sort them in decreasing order, (3) select the top k features (or fewer) whose values are positive. This establishes the following proposition:

Proposition 5.7. *For a given δ and any ϕ -convergent learning dynamic, an optimal stationary sequence of feature selections can be computed in $n \log(n)$ time.*

5.5 When Does the Algorithm Select Informative Features?

Next, we study the following question: under what conditions does the algorithm prioritize selecting the more informative set of features allowing the human to learn (the growth mindset), rather than keep making the myopic choice of the initially less divergent features for short-term rewards (the fixed mindset)? As the example in Section 5.2 illustrates, this depends in our framework on how “patient” the algorithm is (the parameter δ) when considering future rewards, as well as the efficiency of human learning (the parameter ϕ). The definitions and set of results below formalize this intuition.

5.5.1 The algorithm’s patience

Recall that a_i^2 captures the informativeness of feature i , and $(a_i - h_{i,0})^2$ is the initial divergence of the human’s belief about the importance of the feature from the algorithm’s estimate. For a pair of features i and j such that $a_i > a_j$, we denote by $\Delta_{i,j}^I = a_i^2 - a_j^2$ the *informativeness difference* between features i and j , and by $\Delta_{i,j}^D = (a_i - h_{i,0})^2 - (a_j - h_{j,0})^2$ the *divergence difference* between i and j . The following lemma uses the informativeness difference and divergence difference between features i and j to characterize the values of δ for which $V_{\delta,\phi}(\{i\}_{t=0}^\infty) \geq V_{\delta,\phi}(\{j\}_{t=0}^\infty)$. The lemma establishes the main formal basis for analyzing the *fixed vs. growth mindset* tradeoff (Tradeoff 2). It shows that for two features i and j , the more informative feature is either always preferred for any δ , or there is a critical value of δ above which this is true. Intuitively, a feature that has high informativeness but also high divergence, may have low value in the early stages before learning takes place, but high value in later stages as learning converges. Small δ puts relatively more weight on the early stages, and so the informative feature might not be selected, but when δ is large, there is more weight on the outcomes of learning where the more informative feature has higher value.

Lemma 5.8. *Consider a pair of features i and j such that $a_i > a_j$. Then,*

- *If $\Delta_{i,j}^I \geq \Delta_{i,j}^D$, then for any $\delta > 0$, $V_{\delta,\phi}(\{i\}_{t=0}^\infty) \geq V_{\delta,\phi}(\{j\}_{t=0}^\infty)$.*
- *If $\Delta_{i,j}^I < \Delta_{i,j}^D$, there exists $\delta_{i,j}$ such that, for every $\delta < \delta_{i,j}$, $V_{\delta,\phi}(\{i\}_{t=0}^\infty) < V_{\delta,\phi}(\{j\}_{t=0}^\infty)$ and for every $\delta \geq \delta_{i,j}$, $V_{\delta,\phi}(\{i\}_{t=0}^\infty) \geq V_{\delta,\phi}(\{j\}_{t=0}^\infty)$.*

Proof. We first show that there exists at most a single value of δ , such that $V_{\delta,\phi}(\{i\}_{t=0}^\infty) = V_{\delta,\phi}(\{j\}_{t=0}^\infty)$. By Equation (10), we should show that the following equation has at most one solution:

$$\frac{1}{1 - \delta_{i,j}} a_i^2 - \sum_{t=0}^{\infty} \delta_{i,j}^t \phi(t) (a_i - h_{i,0})^2 = \frac{1}{1 - \delta_{i,j}} a_j^2 - \sum_{t=0}^{\infty} \delta_{i,j}^t \phi(t) (a_j - h_{j,0})^2$$

By rearranging and using our notations $\Delta_{i,j}^I$ and $\Delta_{i,j}^D$,

$$\frac{1}{1 - \delta_{i,j}} \cdot \Delta_{i,j}^I = \sum_{t=0}^{\infty} \delta_{i,j}^t \phi(t) \cdot \Delta_{i,j}^D \implies \frac{\Delta_{i,j}^I}{\Delta_{i,j}^D} = (1 - \delta_{i,j}) \sum_{t=0}^{\infty} \delta_{i,j}^t \phi(t)$$

Note that

$$\begin{aligned} (1 - \delta_i) \sum_{t=0}^{\infty} \delta_{i,j}^t \phi(t) &= \sum_{t=0}^{\infty} \delta_{i,j}^t \phi(t) - \sum_{t=0}^{\infty} \delta_{i,j}^{t+1} \phi(t) \\ &= 1 + \sum_{t=1}^{\infty} \delta_{i,j}^t (\phi(t) - \phi(t-1)) = 1 - \sum_{t=1}^{\infty} \delta_{i,j}^t (\phi(t-1) - \phi(t)) \end{aligned}$$

Thus, we need to show that there is at most one solution to the equation:

$$1 - \frac{\Delta_{i,j}^I}{\Delta_{i,j}^D} = \sum_{t=1}^{\infty} \delta_{i,j}^t (\phi(t-1) - \phi(t)) \quad (11)$$

Note that $\phi(t-1) - \phi(t) \geq 0$ since ϕ is decreasing, and since ϕ converges to 0 there should be at least one value of t such that $\phi(t-1) - \phi(t) > 0$. This means that the sum on the right hand-side of the equation is strictly increasing in $\delta_{i,j}$ and hence Equation (11) has at most one solution.

Note that if $\Delta_{i,j}^I \geq \Delta_{i,j}^D$ (i.e., $a_i^2 - a_j^2 \geq (a_i - h_{i,0})^2 - (a_j - h_{j,0})^2$), the left hand-side is 0 or negative and hence the equation doesn't have any solution. In this case, for any value of $\delta \in (0, 1)$, $V_{\delta, \phi}(\{i\}_{t=0}^{\infty}) \geq V_{\delta, \phi}(\{j\}_{t=0}^{\infty})$. If $\Delta_{i,j}^I < \Delta_{i,j}^D$, then the left hand-side is some number between 0 and 1. Let $F(\delta) = \sum_{t=1}^{\infty} \delta_{i,j}^t (\phi(t-1) - \phi(t))$. Note that $F(0) = 0$ and $F(1) = \sum_{t=1}^{\infty} (\phi(t-1) - \phi(t)) = 1$, where the last equality is because the right hand-side is a telescoping sum that equals $\lim_{t \rightarrow \infty} \phi(0) - \phi(t) = 1$ as ϕ converges to 0. Hence, by the intermediate value theorem Equation (11) in this case has exactly one solution. \square

To get a better understanding of the value of the switching point between a less informative feature j and a more informative feature i , it is instructive to consider the closed-form formula of $\delta_{i,j}$ for our running example of a learning dynamic which is $\phi(t) = w^{2t}$ -convergent.

Claim 5.9. *For a w^{2t} -convergent learning dynamic, the value of $\delta_{i,j}$ from Lemma 5.8 is*

$$\delta_{i,j} = \frac{\Delta_{i,j}^I - \Delta_{i,j}^D}{w^2 \cdot \Delta_{i,j}^I - \Delta_{i,j}^D}$$

Proof. Recall that $\delta_{i,j}$ is the solution of

$$1 - \frac{\Delta_{i,j}^I}{\Delta_{i,j}^D} = \sum_{t=1}^{\infty} \delta_{i,j}^t (\phi(t-1) - \phi(t))$$

Note that

$$\begin{aligned} \sum_{t=1}^{\infty} \delta_{i,j}^t (\phi(t-1) - \phi(t)) &= \sum_{t=1}^{\infty} \delta_{i,j}^t (w^{2(t-1)} - w^{2t}) = \frac{(1-w^2)}{w^2} \sum_{t=1}^{\infty} \delta_{i,j}^t (w^{2t}) \\ &= \frac{(1-w^2)}{w^2} \sum_{t=1}^{\infty} (\delta_{i,j} \cdot w^2)^t = \frac{(1-w^2)}{w^2} \cdot \frac{\delta_{i,j} w^2}{1 - \delta_{i,j} w^2} \\ &= \frac{(1-w^2) \delta_{i,j}}{1 - \delta_{i,j} w^2} \end{aligned}$$

Now, we look for $\delta_{i,j}$ such that:

$$\frac{\Delta_{i,j}^D - \Delta_{i,j}^I}{\Delta_{i,j}^D} = \frac{(1-w^2) \delta_{i,j}}{1 - \delta_{i,j} w^2}$$

Hence,

$$\begin{aligned} (\Delta_{i,j}^D - \Delta_{i,j}^I) \cdot (1 - \delta_{i,j} w^2) &= \Delta_{i,j}^D \cdot (1 - w^2) \delta_{i,j} \\ \implies \Delta_{i,j}^D - \Delta_{i,j}^I &= \delta_{i,j} (\Delta_{i,j}^D - w^2 \Delta_{i,j}^D + w^2 \Delta_{i,j}^D - w^2 \Delta_{i,j}^I) \\ \implies \delta_{i,j} &= \frac{\Delta_{i,j}^I - \Delta_{i,j}^D}{w^2 \cdot \Delta_{i,j}^I - \Delta_{i,j}^D} \end{aligned}$$

\square

Lemma 5.8 is the building block of the three results we discuss in this section. We begin by applying it to prove a bound on the number of subsets that are optimal for some value of $\delta \in (0, 1)$.

Proposition 5.10. *Fix a problem instance and ϕ -convergent learning dynamic. Let $A_\phi^*(\delta)$ denote the feature subset selected in an optimal stationary sequence for a patience parameter δ . The number of different subsets for $\delta \in (0, 1)$ is at most $\frac{n(n-1)}{2} + k \leq n^2$.*

Proof. First, we observe that for small values of δ , it is possible that there are fewer than k features with positive values, and hence $|A_\phi^*(\delta)| \leq k$. Observe that the proof of Lemma 5.8 also implies that there is at most a single δ value for which $V_{\delta, \phi}(\{i\}_{t=0}^\infty) = 0$. Hence, there are at most k switching points, where at each, one additional feature is added to $A_\phi^*(\delta)$, until we reach a value δ' such that for every $\delta \geq \delta'$, $|A_\phi^*(\delta)| = k$.

We can now focus on values of $\delta \geq \delta'$ for which $|A_\phi^*(\delta)| = k$. Recall that $A_\phi^*(\delta)$ consists of the k features with the highest values. Observe that $A_\phi^*(\tilde{\delta}) \neq A_\phi^*(\tilde{\delta} + \epsilon)$ for some $\tilde{\delta}$ and for any $\epsilon > 0$ if, for $\tilde{\delta}$, the top k values correspond to different features than for $\tilde{\delta} + \epsilon$. This implies that $\tilde{\delta}$ is a switching point as discussed in Lemma 5.8 of at least two features. Lemma 5.8 also establishes that each pair of features has at most one switching point, implying that the number of switching points and hence optimal subsets of size k is bounded by $\frac{n(n-1)}{2}$. Adding the number of potentially optimal subsets that are smaller than k , we obtain an overall bound of at most $\frac{n(n-1)}{2} + k$ optimal subsets, as required. \square

The next theorem describes the impact of δ on the algorithm's choices. Clearly, when δ is very small, almost all the value comes from the initial time step before any learning takes place, and feature selection is with a fixed mindset. We show that as δ increases and the algorithm puts more weight on future outcomes, it chooses more informative feature subsets (subsets such that $\sum_{i \in A} a_i^2$ is higher), and so feature selection tends to a growth mindset. We also show that for large enough δ , the most informative subset is selected.

Theorem 5.11. *Fix a problem instance and ϕ -convergent learning dynamic. Let $A_\phi^*(\delta)$ denote the feature subset selected in an optimal stationary sequence for a patience parameter δ .*

- *As δ increases, the informativeness of $A_\phi^*(\delta)$ increases.*
- *There exists $\delta^* \in (0, 1)$ such that for all $\delta > \delta^*$, we have that $A_\phi^*(\delta)$ is the most informative subset of features allowed by the budget k .*

Proof. We first prove the statement in the first bullet. As in the proof of Proposition 5.10, we observe that for small values of δ the set of optimal features may include less than k features. For such δ values in each switching point we add more features to the selection set and since a_i^2 is positive the informativeness of $A_\phi^*(\delta)$ can only increase. Next, we focus on values of $\delta \geq \delta'$ for which $|A_\phi^*(\delta)| = k$. As in the proof of Proposition 5.10, we observe that if $A_\phi^*(\tilde{\delta}) \neq A_\phi^*(\tilde{\delta} + \epsilon)$ for some $\tilde{\delta}$ and for any $\epsilon > 0$, then $\tilde{\delta}$ is a switching point of at least two features.

We now turn to showing that as we increase δ , $A_\phi^*(\delta)$ becomes more informative. Consider the switching point $\delta_{i,j}$, where prior to it, feature j was included in the optimal subset, and after it, feature j is no longer part of the optimal subset, but feature i is. This means that j has a higher value than i for $\delta < \delta_{i,j}$ but a lower value for $\delta > \delta_{i,j}$. By Lemma 5.8, this implies that $a_i > a_j$, as required.

As for the second bulleted statement, let A^* denote the subset of the k most informative features. Consider a feature $i \in A^*$ and $j \notin A^*$. By Lemma 5.8, since $a_i^2 > a_j^2$, it is either the case that

$V_{\delta,\phi}(\{i\}_{t=0}^\infty) \geq V_{\delta,\phi}(\{j\}_{t=0}^\infty)$ for any value of $\delta \in (0, 1)$, or there exists $\delta_{i,j} \in (0, 1)$ such that for any $\delta \geq \delta_{i,j}$, $V_{\delta,\phi}(\{i\}_{t=0}^\infty) \geq V_{\delta,\phi}(\{j\}_{t=0}^\infty)$. By setting δ^* to be the maximum of $\delta_{i,j}$'s for any $i \in A^*$ such that $i \in A^*$ and $j \notin A^*$, and using Theorem 5.5, we get that the optimal sequence is a stationary sequence that selects the most informative subset of features A^* . □

5.5.2 The learning dynamic efficiency

In this section, we examine the impact of the efficiency of human learning on the algorithm's selection. We see that, roughly speaking, an optimal subset of features for an efficient learner is more informative than that for a slower learner.

Definition 5.12. *A learning dynamic that is ϕ -convergent is more efficient than a learning dynamic that is ϕ' -convergent if for every t , $\phi(t) \leq \phi'(t)$ and there is at least a single t such that $\phi(t) < \phi'(t)$.*

To take a specific example, in the exponential learning model, a learning dynamic that is w_1^{2t} -convergent is more efficient than a learning dynamic that is w_2^{2t} -convergent if $w_1 < w_2$. At a higher level, we note that learning dynamics that are ϕ -convergent for ϕ with decreasing marginals (i.e., concave) are more efficient than learning dynamics that are ϕ' -convergent for ϕ' with the same marginals as ϕ but in a different order. Intuitively, this suggests that it is more beneficial for the human to invest more in learning during earlier time steps rather than later ones, as this allows the algorithm to select more informative features. Formally, we introduce the following definition to compare the efficiency of two learning dynamics:

Definition 5.13. *Fix a learning dynamic that is ϕ -convergent. $\psi(t) = \phi(t-1) - \phi(t)$ for $t \geq 1$ is the marginal function of $\phi(t)$.*

By this definition, if ϕ has decreasing marginals then ψ is monotonically decreasing.

Claim 5.14. *Consider a learning dynamic that is ϕ' -convergent, and let ψ' denote its marginal function. If ψ' is not decreasing, then a learning dynamic that is ϕ -convergent, where the marginals of ϕ , ψ , are the same as ψ' except that they are sorted in decreasing order, is more efficient.*

Proof. We need to show that for every t , $\phi(t) \leq \phi'(t)$. By definition $\phi(t) = 1 - \sum_{x=1}^t \psi(x)$. Thus, we need to show that

$$1 - \sum_{x=1}^t \psi(x) \leq 1 - \sum_{x=1}^t \psi'(x)$$

The inequality clearly holds since in the construction of ϕ we sorted the values of ψ in decreasing order and hence for any t , $\sum_{x=1}^t \psi(x) \geq \sum_{x=1}^t \psi'(x)$ as required. □

The following proposition shows that fixing a problem instance and a δ value, an optimal sequence for a more efficient learning dynamic selects the same or more informative feature subset.

Proposition 5.15. *Let $A_\delta^*(\phi)$ denote the subset of features chosen in an optimal stationary sequence. If ϕ is more efficient than ϕ' , then $A_\delta^*(\phi)$ is more informative than $A_\delta^*(\phi')$.*

Proof. Recall that an optimal subset includes k features with highest values $V_{\delta,\phi}(\{i\}_{t=0}^\infty) = \sum_{t=0}^\infty \delta^t (a_i^2 - \phi(t)(a_i - h_{i,0})^2)$ (given that those are positive). We show that if $a_i > a_j$, then:

- $V_{\delta,\phi}(\{\{i\}\}_{t=0}^\infty) \leq V_{\delta,\phi}(\{\{j\}\}_{t=0}^\infty) \implies V_{\delta,\phi'}(\{\{i\}\}_{t=0}^\infty) \leq V_{\delta,\phi'}(\{\{j\}\}_{t=0}^\infty)$.
- $V_{\delta,\phi'}(\{\{i\}\}_{t=0}^\infty) \geq V_{\delta,\phi'}(\{\{j\}\}_{t=0}^\infty) \implies V_{\delta,\phi}(\{\{i\}\}_{t=0}^\infty) \leq V_{\delta,\phi}(\{\{j\}\}_{t=0}^\infty)$.

This means that for every pair of features, either both agree on which feature has the higher value, or only learning dynamics that are ϕ -convergent assign feature i a higher value than feature j . As a result, the set of optimal features selected by a ϕ -convergent learning dynamic is at least as informative as the set of features selected by a ϕ' -convergent learning dynamic.

Recall that the informativeness difference between features i and j is $\Delta_{i,j}^I = a_i^2 - a_j^2$ and the divergence difference is $\Delta_{i,j}^D = (a_i - h_{i,0})^2 - (a_j - h_{j,0})^2$. Observe that:

$$\begin{aligned} V_{\delta,\phi}(\{\{i\}\}_{t=0}^\infty) - V_{\delta,\phi}(\{\{j\}\}_{t=0}^\infty) &= \sum_{t=0}^{\infty} \delta^t (a_i^2 - \phi(t)(a_i - h_{i,0})^2) - \sum_{t=0}^{\infty} \delta^t (a_j^2 - \phi(t)(a_j - h_{j,0})^2) \\ &= \frac{1}{1-\delta} \Delta_{i,j}^I - \sum_{t=0}^{\infty} \delta^t \cdot \phi(t) \cdot \Delta_{i,j}^D \end{aligned}$$

Since $a_i > a_j$ we have that $\Delta_{i,j}^I > 0$. If $\Delta_{i,j}^D \leq 0$, then for any learning dynamic that is $\tilde{\phi}$ -convergent, we have that $V_{\delta,\tilde{\phi}}(\{\{i\}\}_{t=0}^\infty) - V_{\delta,\tilde{\phi}}(\{\{j\}\}_{t=0}^\infty) > 0$. The more interesting case is when $\Delta_{i,j}^D > 0$, and hence feature i is more divergent than feature j . For this case we observe that $\sum_{t=0}^{\infty} \delta^t \cdot \phi(t) \leq \sum_{t=0}^{\infty} \delta^t \cdot \phi'(t)$, and hence $V_{\delta,\phi}(\{\{i\}\}_{t=0}^\infty) - V_{\delta,\phi}(\{\{j\}\}_{t=0}^\infty) \geq V_{\delta,\phi'}(\{\{i\}\}_{t=0}^\infty) - V_{\delta,\phi'}(\{\{j\}\}_{t=0}^\infty)$ which implies that both bulleted statements hold, as required. \square

6 Misspecification

We analyze the effect of errors in the algorithm's estimates of the ground-truth coefficients a , the human's coefficients h , and the convergence rate ϕ of the learning dynamic. Modeling errors can lead to incorrect feature selection, reducing overall value. To quantify this, we express these errors as the maximum possible error margin in value per feature. As we will see, it is important to distinguish between "overshoot" errors and "undershoot" errors, as these margins are asymmetric for few parameters. The following definition captures this formally, where $V(\{i\})$ is the true value of feature i depending on whether we are in the static setting or the learning setting, and $V'(\{i\})$ is the value as computed by the imperfect algorithm.

Definition 6.1. *Let ϵ_i denote an upper bound on the magnitude of error in some coefficient of feature i (e.g., ground-truth coefficient or human's coefficient). The error margin of feature i 's value, $V(\{i\})$, is defined by two non-negative functions $\bar{\xi}_i(\epsilon_i)$ and $\xi_i(\epsilon_i)$ such that,*

$$V(\{i\}) - \bar{\xi}_i(\epsilon_i) \leq V'(\{i\}) \leq V(\{i\}) + \xi_i(\epsilon_i)$$

The error margins in estimating specific values can be used to compute the error margin in the algorithm's selection. The idea is to incorporate these margins into the algorithm's decision of whether to include feature i or feature j . Specifically, let A^* denote an optimal set, if for feature $i \in A^*$ and any feature $j \notin A^*$, we have that $V(\{i\}) - V(\{j\}) \geq \bar{\xi}_i(\epsilon_i) + \xi_j(\epsilon_j)$, then the algorithm will also prefer i over j , and the error will not affect the result. Else, the algorithm may choose j instead of i , but the error from this choice will be bounded by $\xi_i(\epsilon_i) + \xi_j(\epsilon_j)$.

Note that we do not assume anything about the structure of the error (e.g., the direction of the error or on which features the algorithm errs). Since the actual performance is affected only by the ranking of

features by their values, increasing ε can have impact only in a collection of thresholds where ranking changes. That is, the change in value is discontinuous and in particular, it is a composition of step functions.

Formally we prove the next proposition:

Proposition 6.2. *Let A^* denote an optimal feature set and A denote the feature set selected by a misspecified algorithm. Let ε_i be a bound on the magnitude of the estimation error for some coefficient of feature i , then:*

$$V(A^*) - V(A) \leq \sum_{i \in A^* \setminus A} \bar{\xi}_i(\varepsilon_i) + \sum_{j \in A \setminus A^*} \xi_j(\varepsilon_j)$$

Proof. Note that,

$$V(A^*) - V(A) = \sum_{i \in A^* \setminus A} V(\{i\}) - \sum_{j \in A \setminus A^*} V(\{j\}) = \sum_{\text{unique}(i,j), i \in A^* \setminus A, j \in A \setminus A^*} V(\{i\}) - V(\{j\})$$

Where the last step is a partition of the features in $A^* \setminus A$ and $A \setminus A^*$ to distinct pairs i and j . If $|A^*| \neq |A|$ we can introduce some dummy features that have value 0 (i.e., $a = a' = h = h' = 0$) to make the two sets equal without affecting anything else. For every pair $i \in A^* \setminus A$ and $j \in A \setminus A^*$, we get

$$V'(\{i\}) - V'(\{j\}) \geq V\{i\} - \bar{\xi}_i(\varepsilon_i) - (V(\{j\}) + \xi_j(\varepsilon_j)) = V(\{i\}) - V(\{j\}) - (\bar{\xi}_i(\varepsilon_i) + \xi_j(\varepsilon_j))$$

Thus, if the difference in values is large, $V(\{i\}) - V(\{j\}) \geq \bar{\xi}_i(\varepsilon_i) + \xi_j(\varepsilon_j)$, the misspecified algorithm should also select feature i instead of feature j . Since it did not, we know that the difference between the features is small and hence, $V(\{i\}) - V(\{j\}) \leq \bar{\xi}_i(\varepsilon_i) + \xi_j(\varepsilon_j)$, which completes the proof. \square

In the remaining of this section we will consider different misspecifications and prove bounds on $\bar{\xi}_i(\varepsilon_i)$ and $\xi_i(\varepsilon_i)$. For this section, we assume that A is a set chosen by the algorithm and A^* is an optimal set according to the true coefficients. By applying Proposition 6.2, the error margins $\bar{\xi}_i(\varepsilon_i)$ and $\xi_i(\varepsilon_i)$ turn to general error bounds on the value of the subset/sequence that the algorithm selects.

6.1 Misspecification in the Non-Learning Setting

Recall that when the human's beliefs are fixed, we have $V(\{i\}, h) = 2a_i h_i - h_i^2$. We first bound the error resulting from the algorithm misspecifying the ground-truth coefficients. As we will see, in this case the value's margin error is symmetric.

Observation 6.3. *Consider a fixed human belief setting and feature i such that $|a'_i - a_i| \leq \varepsilon_i$ and $h'_i = h_i$ for every i , then $\bar{\xi}_i(\varepsilon_i) = \xi_i(\varepsilon_i) = 2\varepsilon_i |h_i|$.*

Proof. Let $\varepsilon' = |a'_i - a_i|$. Observe that,

$$V'(\{i\}) = 2a'_i h_i - h_i^2 = 2(a_i \pm \varepsilon') h_i - h_i^2 = V(\{i\}) \pm 2\varepsilon' h_i$$

Hence, we get that: $V'(\{i\}) \geq V(\{i\}) - 2\varepsilon' |h_i|$ and $V'(\{i\}) \leq V(\{i\}) + 2\varepsilon' |h_i|$. Since both errors are increasing in ε' in their respective directions, we conclude that $\bar{\xi}_i(\varepsilon_i) = \xi_i(\varepsilon_i) = 2\varepsilon |h_i|$ as required. \square

By applying Proposition 6.2 where for every i , we have $\varepsilon_i = \varepsilon$, we get:

Corollary 6.4. *In the non-learning setting, if for every i , $|a'_i - a_i| \leq \varepsilon$ and $h'_i = h$, then $V(A^*) - V(A) \leq 2\varepsilon (\sum_{i \in A^* \setminus A} |h_i| + \sum_{j \in A \setminus A^*} |h_j|)$*

Next, we turn to bound the error resulting from the algorithm misspecifying the human's coefficients.

Observation 6.5. *Consider a fixed human belief setting and feature i such that $|h'_i - h_i| \leq \varepsilon_i$, $a'_i = a_i$, and $\varepsilon \leq |h_i - a_i|$, then $\bar{\xi}_i(\varepsilon_i) = 2\varepsilon_i|h_i - a_i| + (\varepsilon_i)^2$ and $\xi_i(\varepsilon_i) = 2\varepsilon_i|h_i - a_i| - (\varepsilon_i)^2$.*

Proof. Let $\varepsilon' = |h'_i - h_i|$. Observe that,

$$\begin{aligned} V'(\{i\}) &= 2a_i h'_i - (h'_i)^2 = 2a_i(h_i \pm \varepsilon') - (h_i \pm \varepsilon')^2 = V(\{i\}) \pm 2a\varepsilon' \mp 2\varepsilon' h_i - (\varepsilon')^2 \\ &= V(\{i\}) \pm 2\varepsilon'|h_i - a_i| - (\varepsilon')^2 \end{aligned}$$

Hence, $V'(\{i\}) \geq V(\{i\}) - (2\varepsilon'|h_i - a_i| + (\varepsilon')^2)$ as this error margin increases with ε' , we get that $\bar{\xi}_i(\varepsilon_i) = 2\varepsilon_i|h_i - a_i| + (\varepsilon_i)^2$. Now, observe that, $V'(\{i\}) \leq V(\{i\}) + \max\{2\varepsilon'|h_i - a_i| - (\varepsilon')^2, 0\}$. This error margin is increasing up to $\varepsilon' \leq |h_i - a_i|$. We assume that this holds as this simply means that the error of the algorithm is smaller than the error of the human. We conclude that $\xi_i(\varepsilon_i) = 2\varepsilon_i|h_i - a_i| - (\varepsilon_i)^2$. \square

By applying Proposition 6.2 where for every i , we have $\varepsilon_i = \varepsilon$, we get:

Corollary 6.6. *In the non-learning setting, if for every feature i , $|h'_i - h_i| \leq \varepsilon$, $\varepsilon \leq |h_i - a_i|$ and $a'_i = a_i$, then for $|A| = |A^*|$ we have that $V(A^*) - V(A) \leq 2\varepsilon (\sum_{i \in A^* \setminus A} |h_i - a_i| + \sum_{j \in A \setminus A^*} |h_j - a_j|)$. If $|A| \neq |A^*|$, the bound includes additional terms as specified by Proposition 6.2.*

It is interesting to note that when the human agrees with the algorithm on the sign of the coefficients, the error for misspecifying the ground-truth coefficient can be substantially larger than the error for misspecifying the human coefficients. The situation is reversed when the coefficients of the human and the ground truth have opposite signs.

6.2 Misspecification in the Learning Setting

Recall that in this setting:

$$V_{\delta, \phi}(\{i\}_{t=0}^\infty) = \frac{1}{1 - \delta} a_i^2 - \sum_{t=0}^{\infty} \delta^t \phi(t) (a_i - h_{i,0})^2$$

We first consider an algorithm that misspecifies the human's coefficients and observe that the error margins are direct generalization of the error we have seen for the analogous error in the non-learning setting.

Observation 6.7. *Consider the learning setting such that $\phi' = \phi$, and some feature i such that $|h'_i - h_i| \leq \varepsilon_i$, $a'_i = a_i$, and $\varepsilon \leq |h_i - a_i|$, then $\bar{\xi}_i(\varepsilon_i) = \sum_{t=0}^{\infty} \delta^t \phi(t) (2\varepsilon_i \cdot (|a_i - h_{i,0}|) + \varepsilon_i^2)$ and $\xi_i(\varepsilon_i) = \sum_{t=0}^{\infty} \delta^t \phi(t) (\varepsilon^2 + 2\varepsilon_i \cdot (|a_i - h_{i,0}|) - \varepsilon_i^2)$.*

Proof. Let $\varepsilon' = |h'_i - h_i|$. Observe that,

$$\begin{aligned}
V'_{\delta,\phi}(\{(i)\}_{t=0}^\infty) &= \frac{1}{1-\delta} a_i^2 - \sum_{t=0}^\infty \delta^t \phi(t) (a_i - h'_{i,0})^2 \\
&= \frac{1}{1-\delta} a_i^2 - \sum_{t=0}^\infty \delta^t \phi(t) (a_i - h_{i,0} \pm \varepsilon')^2 \\
&= \frac{1}{1-\delta} a_i^2 - \sum_{t=0}^\infty \delta^t \phi(t) (a_i - h_{i,0})^2 - \sum_{t=0}^\infty \delta^t \phi(t) ((\varepsilon')^2 \pm 2\varepsilon' \cdot (a_i - h_{i,0})) \\
&= V_{\delta,\phi}(\{(i)\}_{t=0}^\infty) - \sum_{t=0}^\infty \delta^t \phi(t) ((\varepsilon')^2 \pm 2\varepsilon' \cdot (a_i - h_{i,0}))
\end{aligned}$$

We skip the rest of the analysis as it is the same as in Observation 6.5. □

By applying Proposition 6.2 where for every i , we have $\varepsilon_i = \varepsilon$, we get:

Corollary 6.8. *If $\phi' = \phi$ and for every feature i , $|h'_i - h_i| \leq \varepsilon$, $\varepsilon \leq |h_i - a_i|$ and $a'_i = a_i$, then for $|A| = |A^*|$ we have that,*

$$V_{\delta,\phi}((A^*)_{t=0}^\infty) - V_{\delta,\phi}((A)_{t=0}^\infty) \leq 2\varepsilon \sum_{t=0}^\infty \delta^t \left(\sum_{i \in A^* \setminus A} |h_i - a_i| + \sum_{j \in A \setminus A^*} |h_j - a_j| \right)$$

Next, we consider an imperfect algorithm with inaccurate ground-truth coefficients. In the learning setting, the effect of such errors is more nuanced than in the static setting. First, they capture the informativeness of the features. Second, even when the algorithm initially has a correct estimate of the humans' coefficients, as the human learns based on the ground truth, the algorithm's estimate of the human's coefficients becomes inaccurate. As we will see, this leads to an interesting interplay between these two types of errors.

Our analysis can be viewed as an upper bound on the algorithm's error since, in practice, the algorithm is also likely to learn and improve its estimate of the ground-truth coefficients over time. The interaction in this case is significantly more complex, and we defer its detailed analysis to future work.

Observation 6.9. *Consider the learning setting such that $\phi' = \phi$, and some feature i such that $|a'_i - a_i| \leq \varepsilon_i$, $h'_i = h_i$ and $\varepsilon_i \leq \min\{|a_i|, |a_i - h_i|\}$, then,*

- $\xi_i(\varepsilon_i) = \varepsilon_i^2 \left(\frac{1}{1-\delta} - \sum_{t=0}^\infty \delta^t \phi(t) \right) + 2\varepsilon_i \left| \frac{1}{1-\delta} a_i - \sum_{t=0}^\infty \delta^t \phi(t) (a_i - h_{i,0}) \right|$.
- $\bar{\xi}_i(\varepsilon_i) = \varepsilon_i^2 \left(\sum_{t=0}^\infty \delta^t \phi(t) - \frac{1}{1-\delta} \right) + 2\varepsilon_i \left| \frac{1}{1-\delta} a_i - \sum_{t=0}^\infty \delta^t \phi(t) (a_i - h_{i,0}) \right|$.

Proof. Let $\varepsilon' = |a'_i - a_i|$. Observe that,

$$\begin{aligned}
V'_{\delta,\phi}(\{(i)\}_{t=0}^\infty) &= \frac{1}{1-\delta} (a'_i)^2 - \sum_{t=0}^\infty \delta^t \phi(t) (a'_i - h_{i,0})^2 \\
&= \frac{1}{1-\delta} (a_i \pm \varepsilon_i)^2 - \sum_{t=0}^\infty \delta^t \phi(t) (a_i - h_{i,0} \pm \varepsilon')^2 \\
&= \frac{1}{1-\delta} a_i^2 - \sum_{t=0}^\infty \delta^t \phi(t) (a_i - h_{i,0})^2 + \frac{1}{1-\delta} ((\varepsilon')^2 \pm 2a_i \varepsilon') - \sum_{t=0}^\infty \delta^t \phi(t) ((\varepsilon')^2 \pm 2\varepsilon' \cdot (a_i - h_{i,0})) \\
&= V_{\delta,\phi}(\{(i)\}_{t=0}^\infty) + \frac{1}{1-\delta} (\varepsilon_i^2 \pm 2a_i \varepsilon') - \sum_{t=0}^\infty \delta^t \phi(t) (\varepsilon'^2 \pm 2\varepsilon' \cdot (a_i - h_{i,0})) \\
&= V_{\delta,\phi}(\{(i)\}_{t=0}^\infty) + (\varepsilon')^2 \left(\frac{1}{1-\delta} - \sum_{t=0}^\infty \delta^t \phi(t) \right) \pm \frac{1}{1-\delta} (2a_i \varepsilon') \mp \sum_{t=0}^\infty \delta^t \phi(t) 2\varepsilon' (a_i - h_{i,0})
\end{aligned}$$

We get that:

$$\begin{aligned}
V'_{\delta,\phi}(\{(i)\}_{t=0}^\infty) &\geq V_{\delta,\phi}(\{(i)\}_{t=0}^\infty) - \left(\varepsilon'^2 \left(\sum_{t=0}^\infty \delta^t \phi(t) - \frac{1}{1-\delta} \right) + \left| \frac{1}{1-\delta} (2a_i \varepsilon') - \sum_{t=0}^\infty \delta^t \phi(t) 2\varepsilon' (a_i - h_{i,0}) \right| \right) \\
V'_{\delta,\phi}(\{(i)\}_{t=0}^\infty) &\leq V_{\delta,\phi}(\{(i)\}_{t=0}^\infty) + \varepsilon'^2 \left(\frac{1}{1-\delta} - \sum_{t=0}^\infty \delta^t \phi(t) \right) + \left| \frac{1}{1-\delta} (2a_i \varepsilon') - \sum_{t=0}^\infty \delta^t \phi(t) 2\varepsilon' (a_i - h_{i,0}) \right|
\end{aligned}$$

Note that the first term in the lower bound is negative since for any choice of ϕ which is bounded by 1, we have that $\frac{1}{1-\delta} \geq \sum_{t=0}^\infty \delta^t \phi(t)$. Hence, a-priori it is unclear that the error is increasing with ε' . By taking a derivative with respect to ε' we get that,

$$2\varepsilon' \left(\sum_{t=0}^\infty \delta^t \phi(t) - \frac{1}{1-\delta} \right) + 2 \left| \frac{1}{1-\delta} a_i - \sum_{t=0}^\infty \delta^t \phi(t) \varepsilon' (a_i - h_{i,0}) \right| \geq 0$$

Implying that:

$$\varepsilon' \leq \frac{\left| \frac{1}{1-\delta} a_i - \sum_{t=0}^\infty \delta^t \phi(t) \varepsilon' (a_i - h_{i,0}) \right|}{\frac{1}{1-\delta} - \sum_{t=0}^\infty \delta^t \phi(t)} \leq \min\{|a_i|, |a_i - h_i|\}$$

Hence, we conclude that,

$$\bar{\xi}_i(\varepsilon_i) = \varepsilon'^2 \left(\sum_{t=0}^\infty \delta^t \phi(t) - \frac{1}{1-\delta} \right) + 2\varepsilon' \left| \frac{1}{1-\delta} a_i - \sum_{t=0}^\infty \delta^t \phi(t) (a_i - h_{i,0}) \right|$$

We get that the upper bound is increasing in ε since $\sum_{t=0}^\infty \delta^t \phi(t) - \frac{1}{1-\delta} \geq 0$ and thus,

$$\xi_i(\varepsilon_i) = \varepsilon'^2 \left(\frac{1}{1-\delta} - \sum_{t=0}^\infty \delta^t \phi(t) \right) + 2\varepsilon' \left| \frac{1}{1-\delta} a_i - \sum_{t=0}^\infty \delta^t \phi(t) (a_i - h_{i,0}) \right|$$

□

By applying Proposition 6.2 where for every i , we have $\varepsilon_i = \varepsilon$, we get the following corollary:

Corollary 6.10. *If $\phi' = \phi$ and for every feature i , $|a'_i - a_i| \leq \varepsilon$, $\varepsilon \leq \min\{|a_i|, |a_i - h_i|\}$ and $h'_i = h_i$, then for $|A| = |A^*|$ we have that*

$$V(A^*) - V(A) \leq 2\varepsilon \sum_{i \in A^* \cup A \setminus A^* \cap A} \left| \frac{1}{1-\delta} a_i - \sum_{t=0}^{\infty} \delta^t \phi(t)(a_i - h_{i,0}) \right|$$

It is interesting to further examine the error estimate we obtained. Note that if a_i and $a_i - h_i$ have the same sign, the two types of errors partially cancel out, reducing the total error. However, if they have opposite signs, they compound, increasing the total error.

Lastly, we look at errors in estimating the human's learning dynamic ϕ . There are various measures that can be used to specify a bound on the distance between the true function ϕ and the estimated function ϕ' . In our case, since ϕ always appears in converging sums, it is useful to define the error ε in estimating ϕ , given the patience parameter δ , as the error in estimating the sum $\sum_{t=0}^{\infty} \delta^t \phi(t)$. As explained in Section 3, this quantity summarizes the weight of the human's initial divergence $(a_i - h_{i,0})^2$ in a feature's value estimate; when learning is fast, the weight is small, and when learning is slow, the weight is large. The error in value of a feature, given such an error in ϕ , is:

Observation 6.11. *Consider an algorithm that assumes that human learning is ϕ' -convergent, where in fact it is ϕ -convergent, such that $|\sum_{t=0}^{\infty} \delta^t \phi'(t) - \sum_{t=0}^{\infty} \delta^t \phi(t)| \leq \varepsilon$ and some feature i such that $h'_i = h_i$ and $a'_i = a_i$, then $\bar{\xi}_i(\varepsilon) = \xi_i(\varepsilon) = \varepsilon(a_i - h_{i,0})^2$.*

Proof. Let $\varepsilon' = |\sum_{t=0}^{\infty} \delta^t \phi'(t) - \sum_{t=0}^{\infty} \delta^t \phi(t)|$

$$\begin{aligned} V_{\delta, \phi'}(\{i\}_{t=0}^{\infty}) &= \frac{1}{1-\delta} (a_i)^2 - \sum_{t=0}^{\infty} \delta^t \phi'(t)(a_i - h_{i,0})^2 \\ &= \frac{1}{1-\delta} (a_i)^2 - \left(\sum_{t=0}^{\infty} \delta^t \phi(t) \pm \varepsilon \right) (a_i - h_{i,0})^2 \\ &= V_{\delta, \phi}(\{i\}_{t=0}^{\infty}) \pm \varepsilon (a_i - h_{i,0})^2 \end{aligned}$$

As it is easy to see that the marginal error is increasing in ε' , we conclude that, $\bar{\xi}_i(\varepsilon) = \xi_i(\varepsilon) = \varepsilon(a_i - h_{i,0})^2$. \square

By applying Proposition 6.2 where for every i , we have $\varepsilon_i = \varepsilon$, we get the following corollary:

Corollary 6.12. *If $|\sum_{t=0}^{\infty} \delta^t \phi'(t) - \sum_{t=0}^{\infty} \delta^t \phi(t)| \leq \varepsilon$ and for every feature i , $h'_i = h_i$ and $a'_i = a_i$, then, $V_{\delta, \phi}((A^*)_{t=0}^{\infty}) - V_{\delta, \phi}((A)_{t=0}^{\infty}) \leq \varepsilon \sum_{A^* \oplus A} (a_i - h_{i,0})^2$.*

So we see that the impact of misestimating the speed of learning on the value estimate is the product of two factors: the extent of the error in ϕ and the magnitude of the human's initial error in their belief. If the human's initial belief about a feature is already close to its true coefficient, misestimating the speed of learning does not matter much since there is little for the human to learn. Conversely, when the human's initial error is large, accurately estimating how fast they learn becomes important and errors in this estimates may have larger impact on estimating feature value.

7 Discussion

In this paper, we formally model human decision-making with algorithmic assistance when the human is learning through repeated interactions, and study optimal algorithmic strategies in this context.

We identify two fundamental tradeoffs in the space of AI-assisted decision making. The first tradeoff, “informativeness vs. divergence,” applies to any AI-assisted decision making setting in which the algorithm is designed to choose information for a human’s use. This tradeoff is amplified in repeated interactions where the human learns from experience. In this case, the algorithm that determines what the human will learn, faces a second tradeoff between “fixed vs. growth mindset” – optimize with respect to the human’s current knowledge or provide them with opportunities for learning, which may take longer but lead to better long-term performance. We propose a stylized model that unravels how time preferences and the human’s ability to learn influence the fixed vs. growth mindset tradeoff. Our results highlight the importance of modeling the human decision-making process and incorporating human learning when designing algorithms to assist decision-makers.

In our modeling, we chose assumptions that are both simple and widely used in the literature, yet sufficient to reveal fundamental principles in the interaction between algorithmic assistance and human learning. For example, we assumed that the variables in linear regression are drawn from independent distributions – a common assumption in learning over time literature (e.g., bandit problems and online learning in general [47]). As future work, it would be interesting to explore correlations between variables or decision instances over time. Exploring alternative performance metrics beyond mean squared error and learning models beyond linear regression is also of interest. For example, when performance metrics are discontinuous, such as in binary classification problems, algorithmic decisions may aim to shift the human’s prediction to cross a decision threshold, potentially raising ethical considerations regarding manipulation. An additional direction is to study a human decision maker who recognizes that the algorithm is optimizing and might not provide the features they expect. In this case, the human might act strategically according to their beliefs, leading to a game setting.

We believe that the phenomena we identify extend beyond our specific model, capturing fundamental principles of human-algorithm interactions. These insights lay the groundwork for further research, advancing our understanding of how algorithmic assistance interacts with human learning.

Acknowledgments

This work was supported in part by BSF grant 2018206, a MIT METEOR fellowship, a Vannevar Bush Faculty Fellowship, a Simons Collaboration grant, and a grant from the MacArthur Foundation.

References

- [1] Alex Albright. If you give a judge a risk score: evidence from kentucky bail decisions. *Harvard John M. Olin Fellow’s Discussion Paper*, 85:16, 2019.
- [2] Rohan Alur, Loren Laine, Darrick K Li, Dennis Shung, Manish Raghavan, and Devavrat Shah. Integrating expert judgment and algorithmic decision making: An indistinguishability framework. *arXiv preprint arXiv:2410.08783*, 2024.
- [3] Alejandro Barredo Arrieta, Natalia Díaz-Rodríguez, Javier Del Ser, Adrien Bennetot, Siham Tabik, Alberto Barbado, Salvador García, Sergio Gil-López, Daniel Molina, Richard Benjamins, et al. Explainable artificial intelligence (xai): Concepts, taxonomies, opportunities and challenges toward responsible ai. *Information fusion*, 58:82–115, 2020.
- [4] Sebastian Bach, Alexander Binder, Grégoire Montavon, Frederick Klauschen, Klaus-Robert Müller, and Wojciech Samek. On pixel-wise explanations for non-linear classifier decisions by layer-wise relevance propagation. *PloS one*, 10(7):e0130140, 2015.
- [5] Gagan Bansal, Besmira Nushi, Ece Kamar, Walter S Lasecki, Daniel S Weld, and Eric Horvitz. Beyond accuracy: The role of mental models in human-ai team performance. In *HCOMP 2019*, 2019.
- [6] Gagan Bansal, Besmira Nushi, Ece Kamar, Eric Horvitz, and Daniel S Weld. Is the most accurate ai the best teammate? optimizing ai for teamwork. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 35, pages 11405–11414, 2021.
- [7] Gagan Bansal, Tongshuang Wu, Joyce Zhou, Raymond Fok, Besmira Nushi, Ece Kamar, Marco Tulio Ribeiro, and Daniel Weld. Does the whole exceed its parts? The effect of AI explanations on complementary team performance. In *Proceedings of the 2021 CHI Conference on Human Factors in Computing Systems*, pages 1–16, 2021.
- [8] Oliver Blatchford, William R Murray, and Mary Blatchford. A risk score to predict need for treatment for uppergastrointestinal haemorrhage. *The Lancet*, 356(9238):1318–1321, 2000.
- [9] Sebastian Bordt and Ulrike Von Luxburg. A bandit model for human-machine decision making with private information and opacity. In *International Conference on Artificial Intelligence and Statistics*, pages 7300–7319. PMLR, 2022.
- [10] William Brown and Arpit Agarwal. Diversified recommendations for agents with adaptive preferences. *Advances in Neural Information Processing Systems*, 35:26066–26077, 2022.
- [11] William Brown and Arpit Agarwal. Online recommendations for agents with discounted adaptive preferences. In *International Conference on Algorithmic Learning Theory*, pages 244–281. PMLR, 2024.
- [12] Zana Buçinca, Maja Barbara Malaya, and Krzysztof Z Gajos. To trust or to think: cognitive forcing functions can reduce overreliance on ai in ai-assisted decision-making. *Proceedings of the ACM on Human-computer Interaction*, 5(CSCW1):1–21, 2021.
- [13] Zana Buçinca, Siddharth Swaroop, Amanda E Paluch, Susan A Murphy, and Krzysztof Z Gajos. Towards optimizing human-centric objectives in ai-assisted decision-making with offline reinforcement learning. *arXiv preprint arXiv:2403.05911*, 2024.
- [14] Micah Carroll, Rohin Shah, Mark K Ho, Tom Griffiths, Sanjit Seshia, Pieter Abbeel, and Anca Dragan. On the utility of learning about humans for human-ai coordination. *Advances in neural information processing systems*, 32, 2019.
- [15] Bryan Casey, Ashkon Farhangi, and Roland Vogl. Rethinking explainable machines. *Berkeley*

Technology Law Journal, 34(1):143–188, 2019.

- [16] Lawrence Chan, Dylan Hadfield-Menell, Siddhartha Srinivasa, and Anca Dragan. agree. In *2019 14th ACM/IEEE International Conference on Human-Robot Interaction (HRI)*, pages 354–363. IEEE, 2019.
- [17] CJA. Updating the new york city criminal justice agency release assessment, 2020. URL <https://www.nycja.org/publications/Updating-the-new-york-city-criminal-justice-agency-release-assessment>. Accessed: 2025-02-05.
- [18] European Commission. Proposal for a regulation laying down harmonised rules on artificial intelligence (artificial intelligence act) and amending certain union legislative acts. 2021.
- [19] Sarah Dean, Evan Dong, Meena Jagadeesan, and Liu Leqi. Accounting for ai and users shaping one another: The role of mathematical models. *arXiv preprint arXiv:2404.12366*, 2024.
- [20] Murat Dikmen and Catherine Burns. The effects of domain knowledge on trust in explainable ai and task performance: A case of peer-to-peer lending. *International Journal of Human-Computer Studies*, 162:102792, 2022.
- [21] Yuhan Du, Anna Markella Antoniadis, Catherine McNestry, Fionnuala M McAuliffe, and Catherine Mooney. The role of xai in advice-taking from a clinical decision support system: A comparative user study of feature contribution-based and example-based explanations. *Applied Sciences*, 12(20): 10323, 2022.
- [22] Mira Finkelstein, Lucy Liu, Yoav Kolumbus, David C Parkes, Jeffrey S Rosenschein, Sarah Keren, et al. Explainable reinforcement learning via model transforms. *Advances in Neural Information Processing Systems*, 35:34039–34051, 2022.
- [23] Carolina Garcia-Vidal, Gemma Sanjuan, Pedro Puerta-Alcalde, Estela Moreno-García, and Alex Soriano. Artificial intelligence to support clinical decision-making processes. *EBioMedicine*, 2019.
- [24] Catalina Gomez, Sue Min Cho, Shichang Ke, Chien-Ming Huang, and Mathias Unberath. Human-ai collaboration is not very collaborative yet: a taxonomy of interaction patterns in ai-assisted decision making from a systematic review. *Frontiers in Computer Science*, 6:1521066, 2025.
- [25] Ben Green and Yiling Chen. The principles and limits of algorithm-in-the-loop decision making. *Proceedings of the ACM on Human-Computer Interaction*, 2019.
- [26] Ben Green and Yiling Chen. Disparate interactions: An algorithm-in-the-loop analysis of fairness in risk assessments. In *FAT* 2019*, 2019.
- [27] Sophie Greenwood, Karen Levy, Solon Barocas, Jon Kleinberg, and Hoda Heidari. Designing algorithmic delegates. In *NeurIPS 2024 Workshop on Behavioral Machine Learning*, 2024.
- [28] Nika Haghtalab, Chara Podimata, and Kunhe Yang. Calibrated stackelberg games: Learning optimal commitments against calibrated agents. *Advances in Neural Information Processing Systems*, 36, 2024.
- [29] Karim Hamade, Reid McIlroy-Young, Siddhartha Sen, Jon Kleinberg, and Ashton Anderson. Designing skill-compatible ai: Methodologies and frameworks in chess. In *The Twelfth International Conference on Learning Representations*, 2024.
- [30] Alexandre Heuillet, Fabien Couthouis, and Natalia Díaz-Rodríguez. Explainability in deep reinforcement learning. *Knowledge-Based Systems*, 214:106685, 2021.
- [31] Sophie Hilgard, Nir Rosenfeld, Mahzarin R Banaji, Jack Cao, and David Parkes. Learning represen-

- tations by humans, for humans. In *International conference on machine learning*, pages 4227–4238. PMLR, 2021.
- [32] Jake Hofman, Daniel G. Goldstein, and David Rothschild. A sports analogy for understanding different ways to use ai. *Harvard Business Review*, December 2023. URL <https://www.microsoft.com/en-us/research/publication/a-sports-analogy-for-understanding-different-ways-to-use-ai/>.
- [33] Andrei Iakovlev and Annie Liang. The value of context: Human versus black box evaluators. *arXiv preprint arXiv:2402.11157*, 2024.
- [34] Senerath Mudalige Don Alexis Chinthaka Jayatilake and Gamage Upeksha Ganegoda. Involvement of machine learning tools in healthcare decision making. *Journal of healthcare engineering*, 2021(1): 6679512, 2021.
- [35] Aravindh Mahendran and Andrea Vedaldi. Understanding deep image representations by inverting them. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 5188–5196, 2015.
- [36] Mark A Musen, Blackford Middleton, and Robert A Greenes. Clinical decision-support systems. In *Biomedical informatics*. 2021.
- [37] Omer Nahum, Gali Noti, David C Parkes, and Nir Rosenfeld. Decongestion by representation: Learning to improve economic welfare in marketplaces. In *The Twelfth International Conference on Learning Representations*, 2024.
- [38] National Weather Service. Dew point vs humidity, 2025. URL https://www.weather.gov/arx/why_dewpoint_vs_humidity. Accessed: February 5, 2025.
- [39] Anh Nguyen, Jason Yosinski, and Jeff Clune. Multifaceted feature visualization: Uncovering the different types of features learned by each neuron in deep neural networks. *arXiv preprint arXiv:1602.03616*, 2016.
- [40] NJ-Courts. Annual report to the governor and the legislature. *New Jersey Courts*, 2020.
- [41] Inc. Northpointe. Practitioner’s guide to compas core. <https://www.equivant.com/wp-content/uploads/Practitioners-Guide-to-COMPAS-Core-040419.pdf>, 2019. Accessed: 2023-06-12.
- [42] Gali Noti and Yiling Chen. Learning when to advise human decision makers. In *Proceedings of the Thirty-Second International Joint Conference on Artificial Intelligence, IJCAI ’23*, 2023. ISBN 978-1-956792-03-4. doi: 10.24963/ijcai.2023/339. URL <https://doi.org/10.24963/ijcai.2023/339>.
- [43] Kenny Peng, Nikhil Garg, and Jon Kleinberg. A no free lunch theorem for human-ai collaboration. *arXiv preprint arXiv:2411.15230*, 2024.
- [44] Juan Perdomo, Tijana Zrnic, Celestine Mendler-Dünner, and Moritz Hardt. Performative prediction. In *International Conference on Machine Learning*, pages 7599–7609. PMLR, 2020.
- [45] Wojciech Samek, Thomas Wiegand, and Klaus-Robert Müller. Explainable artificial intelligence: Understanding, visualizing and interpreting deep learning models. *arXiv preprint arXiv:1708.08296*, 2017.
- [46] Amitojdeep Singh, Sourya Sengupta, and Vasudevan Lakshminarayanan. Explainable deep learning models in medical image analysis. *Journal of Imaging*, 6(6):52, 2020.
- [47] Aleksandrs Slivkins et al. Introduction to multi-armed bandits. *Foundations and Trends® in Machine Learning*, 12(1-2):1–286, 2019.
- [48] Ran Tian, Masayoshi Tomizuka, Anca D Dragan, and Andrea Bajcsy. Towards modeling and

influencing the dynamics of human learning. In *Proceedings of the 2023 ACM/IEEE International Conference on Human-Robot Interaction*, pages 350–358, 2023.

- [49] Nenad Tomašev, Xavier Glorot, Jack W Rae, Michal Zielinski, Harry Askham, Andre Saraiva, Anne Mottram, Clemens Meyer, Suman Ravuri, Ivan Protsyuk, et al. A clinically applicable approach to continuous prediction of future acute kidney injury. *Nature*, 2019.
- [50] Philipp Tschandl, Christoph Rinner, Zoe Apalla, Giuseppe Argenziano, Noel Codella, Allan Halpern, Monika Janda, Aimilios Lallas, Caterina Longo, Josep Malvehy, et al. Human–computer collaboration for skin cancer recognition. *Nature Medicine*, 2020.
- [51] Ruqing Xu. Persuasion, delegation, and private information in algorithm-assisted decisions. *arXiv preprint arXiv:2402.09384*, 2024.
- [52] Yunfeng Zhang, Q Vera Liao, and Rachel KE Bellamy. Effect of confidence and explanation on accuracy and trust calibration in ai-assisted decision making. In *FAT* ’20*, 2020.

A On Standardized Features

Lemma A.1. *Let $\{z_i\}_{i=1}^n$ be independent random variables, with known means $\{\mu_i\}_{i=1}^n$ and known finite standard deviations $\{\sigma_i\}_{i=1}^n$, and let $y = a_0 + \sum_{i=1}^n a_i z_i$ be a linear combination of the features that we are willing to represent, for $a_i \in \mathbb{R}$ for all $i \in [n]$. Then, w.l.o.g., we can assume $\mu_i = 0$ and $\sigma_i = 1$ for all $i \in [n]$.*

Proof. For all i , we use $x_i(z_i) = (z_i - \mu_i)/\sigma_i$. Therefore, the distribution of x_i has zero mean and unit variance. We use coefficients a'_i as follows: $a'_0 = a_0 + \sum_{i=1}^n \mu_i a_i$, and for all $i \in [n]$, $a'_i = \sigma_i a_i$. Now, we have:

$$y = a_0 + \sum_{i=1}^n a_i z_i = a_0 + \sum_{i=1}^n \mu_i a_i + \sum_{i=1}^n \sigma_i \cdot a_i (z_i - \mu_i) / \sigma_i = a'_0 + \sum_{i=1}^n a'_i x_i = y'$$

Thus, when the means and standard deviations of the z_i 's are known, using the standardized features is simply a change of variables that does not change the problem. □

B Figure for the Example in Section 5.2

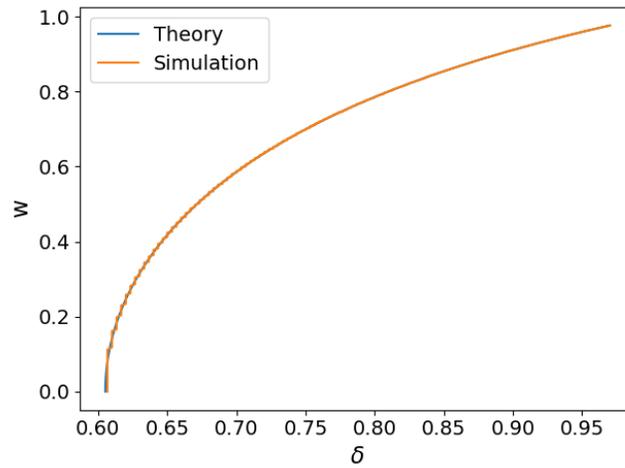


Figure 2: The transition curve for the example in Section 5.2: For δ values to the right of this curve, selecting the more informative feature (of the two) is preferred, while for δ values to the left of this curve, the less divergent feature is preferred. The plot demonstrates that the empirical curve in Figure 1a coincides with the theoretical curve derived in Claim 5.9.