

A Control Theoretic Model of Adaptive Learning in Dynamic Environments

Harrison Ritz, Matthew R. Nassar, Michael J. Frank, and Amitai Shenhav

Abstract

■ To behave adaptively in environments that are noisy and nonstationary, humans and other animals must monitor feedback from their environment and adjust their predictions and actions accordingly. An understudied approach for modeling these adaptive processes comes from the engineering field of control theory, which provides general principles for regulating dynamical systems, often without requiring a generative model. The proportional–integral–derivative (PID) controller is one of the most popular models of industrial process control. The proportional term is analogous to the “delta rule” in psychology, adjusting estimates in proportion to each error in prediction. The integral and derivative terms augment this update to simultaneously improve accuracy and stability. Here, we tested whether the PID algorithm can describe how people sequentially adjust their predictions in response to new information. Across three

experiments, we found that the PID controller was an effective model of participants’ decisions in noisy, changing environments. In Experiment 1, we reanalyzed a change-point detection experiment and showed that participants’ behavior incorporated elements of PID updating. In Experiments 2–3, we developed a task with gradual transitions that we optimized to detect PID-like adjustments. In both experiments, the PID model offered better descriptions of behavioral adjustments than both the classical delta-rule model and its more sophisticated variant, the Kalman filter. We further examined how participants weighted different PID terms in response to salient environmental events, finding that these control terms were modulated by reward, surprise, and outcome entropy. These experiments provide preliminary evidence that adaptive learning in dynamic environments resembles PID control. ■

INTRODUCTION

To behave adaptively, we must adjust our behavior in response to the dynamics of our environment (Pezzulo & Cisek, 2016; Ashby, 1956). Achieving this goal requires us to collect feedback about the outcomes of our recent actions and incorporate this feedback into decisions about how to adjust future actions. Within research on learning and decision-making, a popular approach for achieving this feedback-based control is the “delta-rule model”¹ ($\Delta x = \alpha \delta$; Widrow & Hoff, 1960; cf. Maxwell, 1868). This model adjusts expectations (x) proportionally to the discrepancy between observed and predicted outcomes (i.e., prediction error, δ), depending on the learning rate (α).

Although there is substantial cross-species evidence for delta-rule controlled behavior (e.g., Garrison, Erdeniz, & Done, 2013; Mirenowicz & Schultz, 1994; Rescorla & Wagner, 1972), this algorithm has major limitations. The delta rule is sensitive to any noise that will cause persistent errors, leading to either oscillatory behavior (at a high learning rate) or a sluggish response (at a low learning rate; Aström & Murray, 2008; Rumelhart, Hinton, &

Williams, 1986). However, one of the greatest limitations of this algorithm is that it performs poorly in environments that are nonstationary (i.e., that change discontinuously over time; Aström & Murray, 2008; Pearce & Hall, 1980).

More elaborate feedback control mechanisms have been developed within a branch of engineering called Control Theory that studies the regulation of dynamical systems. Many control theoretic algorithms augment the basic delta rule with additional control terms that greatly improve accuracy, stability, and responsivity. The most popular variant of these control theoretic models is the popular proportional–integral–derivative (PID) controller (Figure 1). This model is simple, accurate, and robust, with response properties that have been well characterized over the last century (Aström & Murray, 2008; Franklin, Powell, & Emami-Naeini, 1994). The PID controller takes the error from a reference signal as input, and it outputs a control signal consisting of a linear combination of control signals proportional to the error (P-Term), the integral of the error (I-Term), and the derivative of the error (D-Term; Figure 1). These three terms minimize deviations from the reference based on errors in the present, past, and expected future, respectively.

Proportional control (cf. delta-rule control) directly minimizes deviation from the reference and is often the primary driver of the control process. Integral control

This paper is part of a Special Focus deriving from a symposium at the 2017 annual meeting of Cognitive Neuroscience Society, entitled, “Multiple neurocomputational, motivational, and mnemonic mechanisms for decision-making.”
Brown University

provides low-frequency compensation for residual steady-state errors, allowing the controller to reduce noise and track gradual changes in the environment. Derivative control provides high-frequency compensation that increases stability, such as by dampening control adjustments when the controller is approaching the reference or increasing adjustments when the reference or environment suddenly changes (see Aström & Murray, 2008). Intuitively, integral control provides low-frequency compensation by combining several time points, whereas derivative control provides high-frequency compensation by tracking the instantaneous change. Here, we test whether this popular model of industrial control can account for adjustments in human behavior within a dynamic environment.

PID control has algorithmic properties that make it useful for most control systems. For instance, relative to algorithms that require an explicit representation of task dynamics, PID can provide an effective, and computationally cheaper, model-free alternative to adjusting cognitive or behavioral processes over time, particularly for natural environments that require particularly complex world models. Moreover, convergent evidence suggests that the PID algorithm may help account for the variety of feedback-related findings observed in humans and other primates. Behavioral and neural correlates of feedback-controlled choice provide preliminary evidence that participants transform decision-relevant variables in a manner predicted by the PID algorithm. Consistent with proportional control, there is substantial evidence that participants adjust their behaviors based on recent errors or conflict (Ullsperger, Danielmeier, & Jochem, 2014; Lau & Glimcher, 2005; Gratton, Coles, & Donchin, 1992; Rescorla & Wagner, 1972; Rabbitt, 1966), with corresponding signals observed most prominently in the striatum and ACC (Smith et al., 2015; Garrison et al., 2013; Matsumoto, Matsumoto, Abe, & Tanaka, 2007; Seo & Lee, 2007; Kennerley, Walton, Behrens, Buckley, & Rushworth, 2006; Ito, Stuphorn, Brown, & Schall, 2003; Mirenowicz & Schultz, 1994; Niki & Watanabe,

1979). Previous work has found that people are also sensitive to the extended history of errors or conflict (Aben, Verguts, & Van den Bussche, 2017; Wittmann et al., 2016; Alexander & Brown, 2015; Bugg & Crump, 2012; Botvinick, Braver, Barch, Carter, & Cohen, 2001; Logan & Zbrodoff, 1979; Laming, 1968), with proposals that this specifically involves integrating over recent errors (Wittmann et al., 2016; Alexander & Brown, 2015). Accordingly, experiments have found neural signals in pFC and ACC that reflect this feedback history (Wittmann et al., 2016; Bernacchia, Seo, Lee, & Wang, 2011; Blais & Bunge, 2010; Carter et al., 2000), and recent models of ACC have emphasized the role that integrative, recurrent activity in this region plays in executive control (Shahnazian & Holroyd, 2018; Hunt & Hayden, 2017; Wang, 2008). Finally, consistent with derivative control, prior work has found that participants track the environmental rate of change when making decisions, with associated neural correlates in the anterior pFCs and ACC (Wittmann et al., 2016; Jiang, Beck, Heller, & Egner, 2015; McGuire, Nassar, Gold, & Kable, 2014; Kovach et al., 2012; Bernacchia et al., 2011; Behrens, Woolrich, Walton, & Rushworth, 2007). Although some of these results have been attributed to participants' representations of environmental dynamics (e.g., Jiang et al., 2015; McGuire et al., 2014; Behrens et al., 2007), PID control may offer a more parsimonious account of these behaviors.

Despite the success of the PID model as a simple and effective algorithm for implementing control in other fields, as well as suggestive evidence for relevant neural signatures in circuits involved in adaptive control, PID has yet to be formally tested as a model of human adaptive learning. In the current set of experiments, we directly tested whether a PID model can describe human performance in adaptive learning tasks. In Experiment 1, we reanalyzed a recent study that examined predictive inference in an environment with discrete change points. Behavior on this task confirmed key predictions of the PID model but was limited in its ability to adjudicate between candidate models. Informed by our findings in

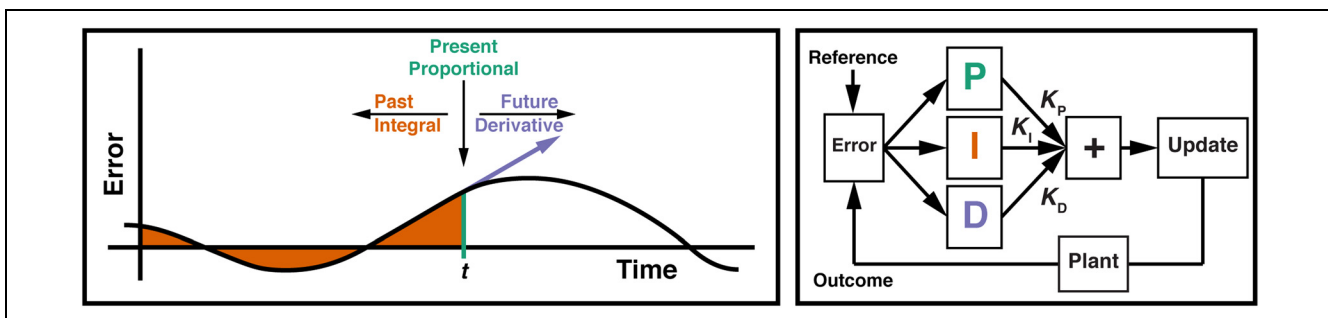


Figure 1. (Left) The PID controller uses the current error (P-Term), the integral of the error (I-Term), and the derivative of the error (D-Term) to provide compensation based on the present, past, and future, respectively. (Right) The PID controller specifies a control signal based on the weighted sum of the PID terms, with each term weighted by their respective gain. Similar to a thermostat-controlled furnace, the plant implements the control signal, moving the measured processes closer to the reference (e.g., the desired temperature). The figure on the left was adapted from Aström and Murray (2008) with permission from Princeton University Press.

Experiment 1, for Experiments 2–3, we developed a novel task that was optimized for PID control, using gradual rather than sudden change points. We found that the PID model was a strong predictor of participants' choices in both experiments. Experiment 3 replicated the predictive power of our model and further examined whether participants dynamically adjust their control terms based on rewards, surprise, and outcome entropy. Across these tasks, participants' performance confirmed key predictions of the PID model, demonstrating that this simple model provides a promising account of adaptive learning.

EXPERIMENT 1

The PID model is designed to adapt the behavior of a system in response to changes in the environment. We therefore began by testing whether this model could explain behavioral adjustments in an existing change-point detection task, one that was designed to assess how humans can adapt their learning rate to uncertain and volatile outcomes (McGuire et al., 2014). In this experiment, participants predicted where a target stimulus would appear (horizontal location on the screen) and then received feedback about where the true location of the outcome had been (see Figure 2A). Outcome locations were normally distributed around a mean, and the mean of this distribution changed suddenly throughout the experiment (at change points), according to a predetermined hazard rate. This task allows us to measure participants' choices and feedback in a continuous space with high precision, making it desirable for studying PID control. We can therefore use the PID model to predict trial-to-trial adjustments in participant behavior (predicted locations) based on their history of error feedback. In other respects, this task is not ideally suited for testing our model: The dramatic changes in target distributions may “reset” adaptive control processes (Tervo et al., 2014; Karlsson, Tervo, & Karpova, 2012; Bouret & Sara, 2005), and so this experiment serves as a preliminary test of our hypothesized control dynamics. We will address these concerns in Experiments 2–3.

Methods

Participants and Procedure

Experiment 1 consisted of a reanalysis of a change-point task used by McGuire and colleagues (2014). Briefly, 32 participants (17 women; mean age = 22.4 years, $SD = 3.0$ years) selected a location on a line using a joystick and then were shown the correct location for that trial. On a random subset of trials, participants were rewarded according to their accuracy, dissociating the trial value (which depended on whether the trial was rewarded) from errors. Target locations were drawn from a Gaussian distribution with a variable

mean and variance. The mean was stable for three trials and then, on a weighted coin flip (hazard rate = 0.125), was uniformly redrawn from the line coordinates; the variance alternated between high and low levels across blocks. Participants performed 160 training trials followed by four blocks of 120 trials during fMRI. Two participants were excluded for having an inconsistent number of trials per block, leaving 30 participants for the final analysis. See McGuire et al. (2014) for additional details.

Lagged Regression Analysis

A critical prediction of the PID model is that participants' updates (i.e., the change in their location guesses) should depend on their history of errors. Whereas a delta-rule model predicts that only the error on the current trial will directly influence updates, a PID controller integrates errors over a longer history, enabling the controller to correct for a consistent bias in errors. Integral control will manifest as an exponentially decaying influence over previous errors, whereas derivative control will place a positive weight on the current error and a negative weight on the $t-1$ error. These two terms make different predictions for the $t-1$ error: Integral control will place a high weight on this error, whereas derivative control will place a lower weight on $t-1$ than it does on earlier trials.

To measure the independent contribution of each trial's feedback in the recent past, we used a simple lagged regression analysis to test how prediction updates (change in predicted location from the current trial to the next trial) depended on the errors from the current and 10 previous trials ($u_t \sim 1 + e_t + e_{t-1} + \dots + e_{t-10}$; Wilkinson notation). We assessed the influence of previous trials' feedback by testing whether the sum of previous trials' betas was significantly different from zero, using a nonparametric sign-randomization test at the group level (comparing the observed results with a null distribution that we generated by randomly assigning positive or negative signs to each participant's summed betas). Throughout the article, all randomization tests used 10^5 simulations, and all statistical tests were two-tailed with $\alpha = .05$.

PID Model

The PID algorithm controls a system to maintain a desired reference signal (Figure 1). It takes as input the signed error relative to this reference ($e_t = \text{reference} - \text{output}$) and produces a control signal (u_t) that specifies the adjustment for the next time point (here, the next trial). The control signal is defined by a linear combination of three terms: the P-Term (reflecting the error), the I-Term (reflecting the leaky integration of the error), and the D-Term (reflecting the derivative of the error). Each of these terms was weighted by its own gain parameter

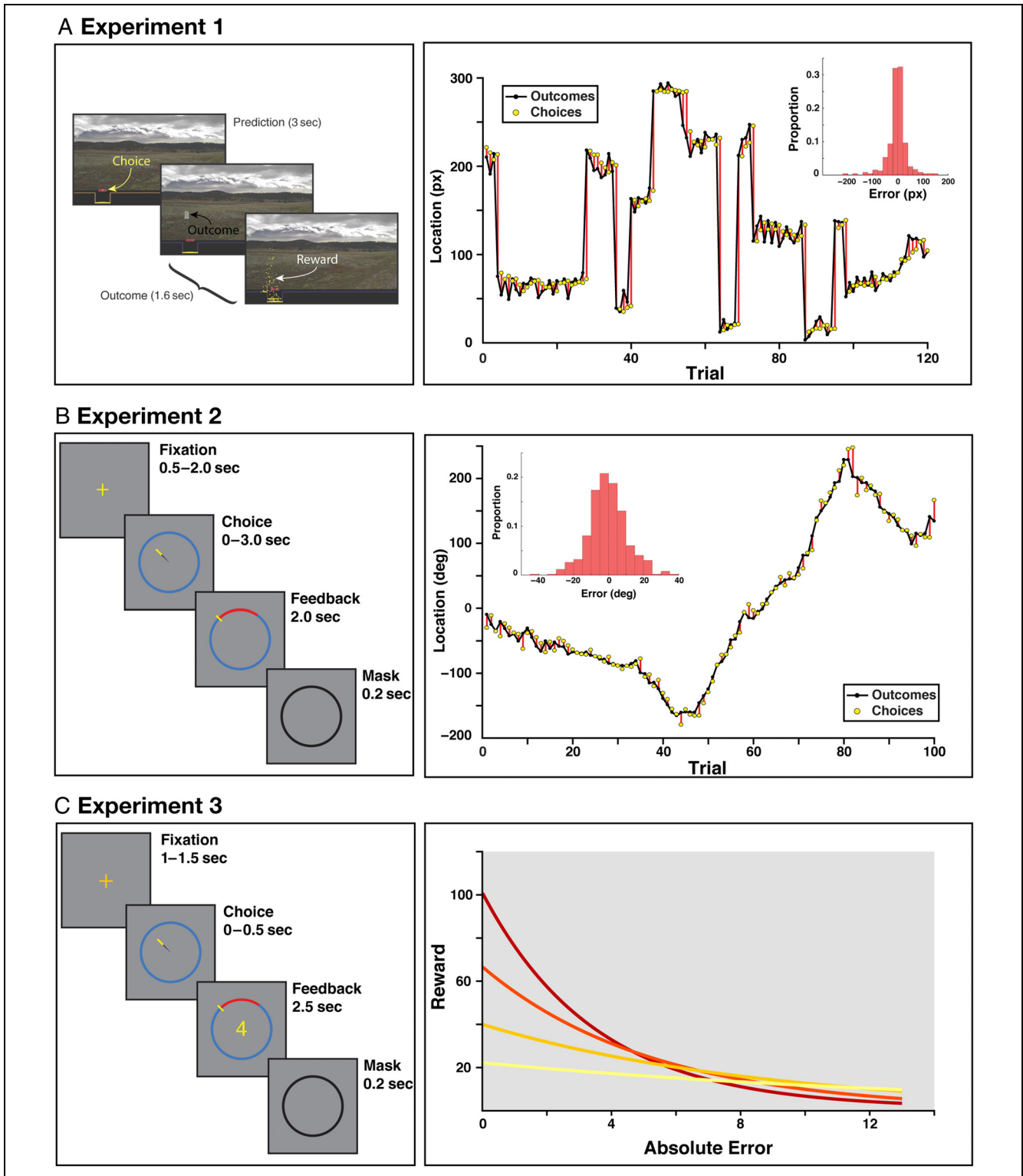


Figure 2. Experimental tasks. (A, left) On each trial of Experiment 1, participants selected a horizontal location with a joystick and were then shown the correct location. On a random subset of trials, participants received performance-contingent rewards (shown as gold coins). Figure was adapted from McGuire et al. (2014). (Right) A representative block of trials from an example participant. The mean correct location was stable for a variable number of trials and then was uniformly resampled. (Inset histogram) The distributions of errors for this participant across their session. (B, left) Participants in Experiment 2 selected a location on the circle with their mouse and then were shown the correct location. (Right) A representative block of trials, demonstrating that the mean correct location changed gradually over time. As seen in the histogram for an example participant (inset), the gradual changes in location for this task resulted in error distributions that were less peaked than in Experiment 1 (compare A inset). (C, left) Experiment 3 was identical to Experiment 2, but participants were rewarded based on their accuracy, according to one of four reward–error functions. They were informed of the current reward mode during fixation, and during feedback, they received the reward corresponding to their accuracy on that trial (conditional on the current reward mode). (Right) Error–reward slopes for the four reward modes.

(K_P , K_I , and K_D). For trial t , the control signal (u_t) was generated by transforming the error (e_t) as follows:

$$u_t = K_P e_t + K_I \sum_{n=1}^t \lambda^{t-n} e_n + K_D (e_t - e_{t-1})$$

where λ represents a memory persistence parameter, with larger values leading to longer retention. On the first trial of a block, the I-Term is e_t and the D-Term is 0, producing a control action similar to proportional control. In the following tasks, u_t was defined as the difference in the choice location between trial $t + 1$ and trial t (hereafter, the “update”), and e_t was the difference between the correct location and the chosen location. Although PID is not traditionally a state estimation algorithm, it can serve this function by regulating performance to maintain a desired accuracy (i.e., no error), such as what occurs in autoencoder learning systems (Denève, Alemi, & Bourdoukan, 2017).

PID Model Fit

We used each participant’s time course of errors within each block to generate hypothesized P, I, and D values based on the raw errors, the integral of the errors, and the first derivative of the error, respectively. Our regression model consisted of an intercept and the three PID terms ($u \sim 1 + P + I + D$), and we fit this model with iteratively reweighted robust regression (using MATLAB’s `fitlm` function; bisquare weighting factor) to minimize overfitting to the outliers that can occur when participants make scalar responses. Fit statistics were generated based on the non-reweighted residuals from the robust model, to avoid undue bias in favor of complex models.

Because the λ parameter (memory persistence) interacted with the identifiability of our PID terms when estimated jointly (e.g., when $\lambda = 0$, P and I are identical), we chose to fit this single term at the group rather than individual level. We fit λ with a grid search (range = 0.5–1, increments of 0.001), using median R^2 across participants as our measure of fit (normalizing individual differences in update variability). Regression models were estimated at the individual level, and regression weights were tested for deviance from zero at the group level with a sign-randomization test (see above).

We compared the P (i.e., delta rule), PI, PD, and PID models, as these are the most common instantiations of the PID algorithm. To quantify model performance, we calculated each participant’s Akaike information criterion (AIC; Akaike, 1983), an index of goodness-of-fit that is penalized for model complexity.² We compared the AIC at the group level using Bayesian model selection (Rigoux, Stephan, Friston, & Daunizeau, 2014), quantifying the Bayesian omnibus risk (BOR; probability that all models are equally good across the population) and each model’s protected exceedance probability (PXP; the probability that this model is more frequently the best

fit than any of the competing models, controlling for the chance rate). BOR tests whether there is an omnibus difference between models, whereas PXP describes which models fit the best.

PID Controller Simulations

To better understand the expected range of behavior under our candidates’ models, we simulated delta-rule and PID controller’s performance for the outcome histories that each participant encountered during the experiment. We used a restricted maximum likelihood estimation procedure (MATLAB’s `fmincon`) to determine the values of the PID gains, λ , and choice bias (i.e., intercept) that perform best given the outcomes of each participant’s task. We then compared these best-performing delta-rule and PID gains with the gains estimated from participants’ behaviors.

We also tested whether simulated behavior from our models produce the same pattern of behavior that we measured with our lagged regression. We simulated an ideal observer that used each participant’s estimated PID parameters and outcome history to generate a sequence of updates and then fit our lagged regression to this simulated behavior, separately for the P, PI, and PID candidate controllers. This analysis allows us to qualitatively determine the extent to which the PID model can act as a generative model of participants’ decision-making behaviors (Nassar & Frank, 2016; Gelman, Meng, & Stern, 1996). If participants are using PID control, then simulated updates from a PID controller should similarly weight the feedback received over previous trials.

Results

Model-agnostic Analysis

To identify the degree to which behavioral adjustments were influenced by recent feedback, we regressed participants’ current and previous errors on their update. We found that although the current error was the strongest predictor of updates, errors from previous trials also influenced adjustment (Figure 3A). The sum of leading trials’ betas was reliably less than zero (summed beta: mean = -0.040 , $SD = 0.064$; $p = .00017$). This suggests that, although immediate feedback was the most influential factor for participants’ updates, they also incorporated an extended history of error feedback into their adjustments. Whereas the current trial had a positive influence on updates, these previous trials instead had a negative influence on the current update, potentially compensating for the extreme errors that participants made at a change point (see Figure 2A).

To verify that our model can generate performance that captures the behavior observed in this task, we simulated behavior on this task using parameters estimated for our PID model and reduced versions thereof (P and

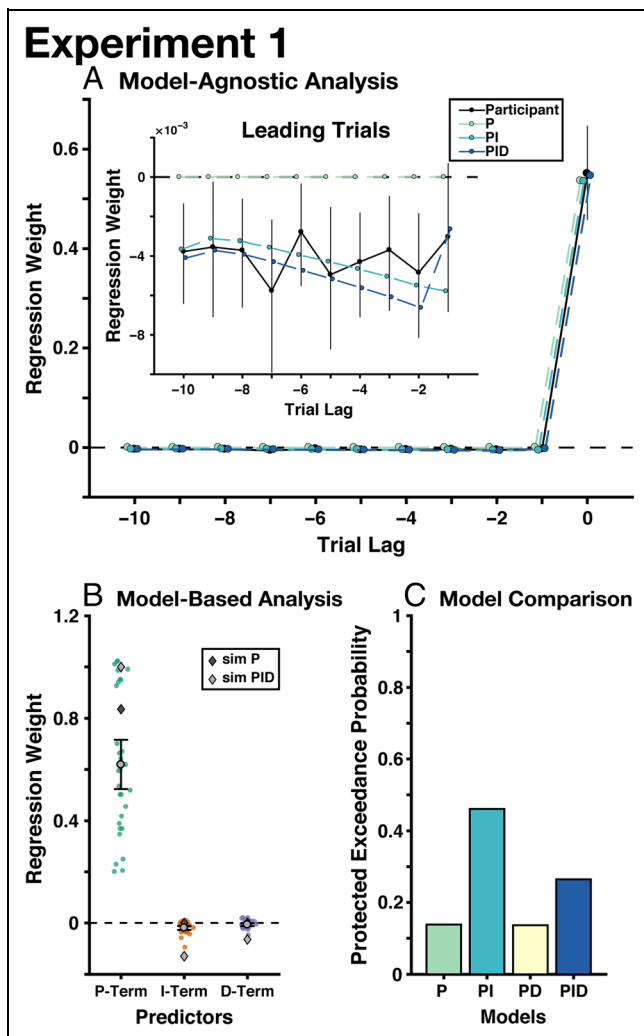


Figure 3. Experiment 1 results. (A) For our model-agnostic analyses, we regressed the errors that participants made on the current and 10 leading trials on their current update (black: participants' regression weights). Next, we used the PID parameters estimated in our regression analysis to generate behavior from P, PI, and PID controllers and fitted our lagged regression to this simulated behavior (colored lines). (Inset) The regression weights from only the leading trials (i.e., before the current trial), controlling for the effect of the current trial. (B) For our model-based analyses, we regressed the trial-wise P-, I-, and D-Terms on participants' updates and found that all three terms were significantly different from zero. Colored circles indicate individual participants' regression weights. Dark gray diamonds indicate the mean gains from the best-performing delta-rule models, based on each participant's outcome history; light gray diamonds indicate the mean gains from the best-performing PID controllers. (C) We used Bayesian model selection to adjudicate between our candidate models, finding that the PI model best explained the data, albeit with moderate support (see text). Error bars throughout indicate mean and between-participant bootstrapped 95% confidence intervals.

PI). We then performed the same lagged regression on these simulated data that we used on real data (Figure 3A). As expected, we found that the simulated PI and PID models captured the influence of leading errors, unlike the P-only model (which predicts that there should be no influence of leading errors).

PID Model Fit

We first performed a search to identify the PID gains that optimized task performance (minimizing mean squared prediction error) for the outcome sequence observed by each participant. We found that the optimal PID gains were all reliably different from zero (mean [SD] PID gain: $K_P = 1.0$ [0.18], $K_I = -0.16$ [0.12], $K_D = -0.070$ [0.077], $\lambda = 0.85$ [0.071]). Consistent with the lagged regression analysis, the optimized integral and derivative gains were negative.

Fitting our PID model to participants' updates, we found that the best-fit models accounted for a substantial amount of this variance (median $R^2 = .92$), with parameters for all terms being significantly different from zero (mean [SD] standardized betas: $\beta_P = 0.62$ [0.27], $p \leq 10^{-5}$; $\beta_I = -0.019$ [0.022], $p \leq 10^{-5}$; $\beta_D = -0.0068$ [0.016], $p = .022$; see Figure 3B). The group level λ (memory persistence) was also quite high (0.9430), suggesting that participants retained a great deal of information regarding past feedback. Participants' estimated gains qualitatively resembled the gains produced by the simulated PID controller, sharing the same sign and rank order (compare gray diamonds and circles in Figure 3B).

We used Bayesian model selection to compare the fit of each model (PXP) and tested whether there was an omnibus difference between models (BOR). We found that the PI model had the highest PXP ($PXP_P = 0.14$, $PXP_{PI} = 0.46$, $PXP_{PD} = 0.14$, $PXP_{PID} = 0.26$; Figure 3C) but that there is altogether insufficient evidence to support one model over another ($BOR = 0.55$, providing roughly equal evidence that the models are the same or different). These data therefore do not allow us to rule out the possibility that a simple delta-rule (P-only) model parsimoniously accounts for participant behaviors. The PID models did not predict behavior better than the Bayesian change-point model in the original publication (original median $R^2 = .97$; McGuire et al., 2014), which incorporated information about the generative structure of the statistical environment.

Discussion

We found preliminary evidence that participants performing a change-point detection task are influenced by their history of error feedback, consistent with the predictions of a PID controller. Participants' updates could also be predicted from the integral and derivative of their errors. Despite these promising indications of PID control, we were unable to confidently differentiate between candidate models. Furthermore, this model did not explain behavior better than the change-point detection model from this original experiment.

Although this experiment offers mixed evidence in favor of the PID algorithm, this may be because this task was designed for change-point models, with sudden,

dramatic shifts in the outcome distribution. These change points introduce extreme errors that participants might treat categorically differently from normal prediction errors, evoking a “reset” in their decision process or causing the representation of a different context (McGuire et al., 2014; O’Reilly et al., 2013; Nassar, Wilson, Heasly, & Gold, 2010; Courville, Daw, & Touretzky, 2006; Bouret & Sara, 2005). This experiment also involved a training paradigm designed to make participants aware of the generative structure of the task, an advantage not typically afforded in the real world or exploited by PID control systems. With these concerns in mind, we developed a novel adaptive learning task in which outcomes changed smoothly over time, encouraging participants to treat outcomes as arising from a single, changing context. Participants were not instructed on this generative structure explicitly, reducing the potential for the use of structured inference strategies that best characterized learning in Experiment 1.

EXPERIMENT 2

Although Experiment 1 provided promising evidence that participants use their history of feedback in a way that resembled PID control, it did not provide definitive evidence as to whether this is the best explanation for the data. However, participants may have strategically reset their predictions at extreme change points, making it more difficult to measure history-dependent predictions. To address this, we developed a task with gradual transitions in which participants tracked an outcomes distribution whose mean linearly changed from one location to another and whose variance changed randomly throughout the block. To make these location transitions seem more continuous, outcomes appeared along a circle rather than a straight line, thus also avoiding edge effects that can occur at either end of a screen. This design allowed us to precisely measure participants’ predictions, errors, and adjustments within an environment whose dynamics are more fluid and predictable than Experiment 1. This task was explicitly designed to emulate an environment for which a PID controller is well suited and specifically to maximize our power to detect differences between PID control and proportional (delta-rule) control.

Methods

Participants and Procedure

Twenty-nine Brown University undergraduate students (25 women; mean age = 18.6 years, $SD = 0.83$ years) performed a supervised learning task in which they predicted an outcome location on a circular display (see Figure 2B).

Participants completed five blocks of 100 trials in which they used a mouse cursor to guess a location on the circumference of the circle. They were then shown

the correct location, with an arc indicating the magnitude and direction of their error. Participants completed 50 training trials before the main experiment. Participants had up to 3 sec to make their guess, or else their final cursor angle would be chosen as the guess for that trial, and feedback was presented for 2 sec. Our final analysis excluded any trials where participants did not move their cursor to the edge of the circle as well as a subset of trials after aberrant feedback due to a technical issue (1.8% of the total trials).

The target location for each trial was drawn from a Gaussian distribution over arc degrees, with a mean and a standard deviation that systematically changed over time. On a weighted coin flip (hazard rate = 0.80), the distribution’s mean shifted based on a random draw from $U(-180, 180)$ degrees. After the new mean was drawn, the mean transitioned from the old mean to the new mean over $U(8, 20)$ trials, with the means during transition trials linearly interpolated between the old and new means. The standard deviation varied independently of the mean and was redrawn from $U(1, 8)$ degrees on a weighed coin flip (hazard rate = 0.40). These task parameters were selected through simulation to maximally differentiate the performance of PID and delta-rule models. Unless otherwise indicated, methods of analysis and model selection for this study are identical to Experiment 1.

Results

Model-agnostic Analysis

Regressing the current and 10 leading errors onto the current update (see Methods under Experiment 1), we again found that the sum of leading errors was significantly different from zero (summed leading betas: mean = 0.27, $SD = 0.31$, $p \leq 10^{-5}$; Figure 4A). This replicates the observation in Experiment 1 that participants incorporate the extended history of errors into their prediction process.

Fitting our lagged model to behavior generated from our models produced a similar pattern of predictions as in Experiment 1: The P-only model categorically failed to capture the influence of leading errors. The PI and PID models were similar in their ability to recreate participants’ use of previous errors; however, the PID model seemed to better capture participants’ weighting of recent leading errors (i.e., over the previous three trials). We examined whether this trend was reliable across participants by fitting linear and quadratic trends over trials to each participant’s leading betas. We found a significant quadratic trend (quadratic trend standardized beta: mean = -0.0054 , $SD = 0.009$, $p = .004$) but not a linear trend (linear trend standardized beta: mean = 0.002, $SD = 0.020$, $p = .39$). Although this finding is broadly compatible with derivative control’s nonlinear weighting of previous errors, the observed trend extended further

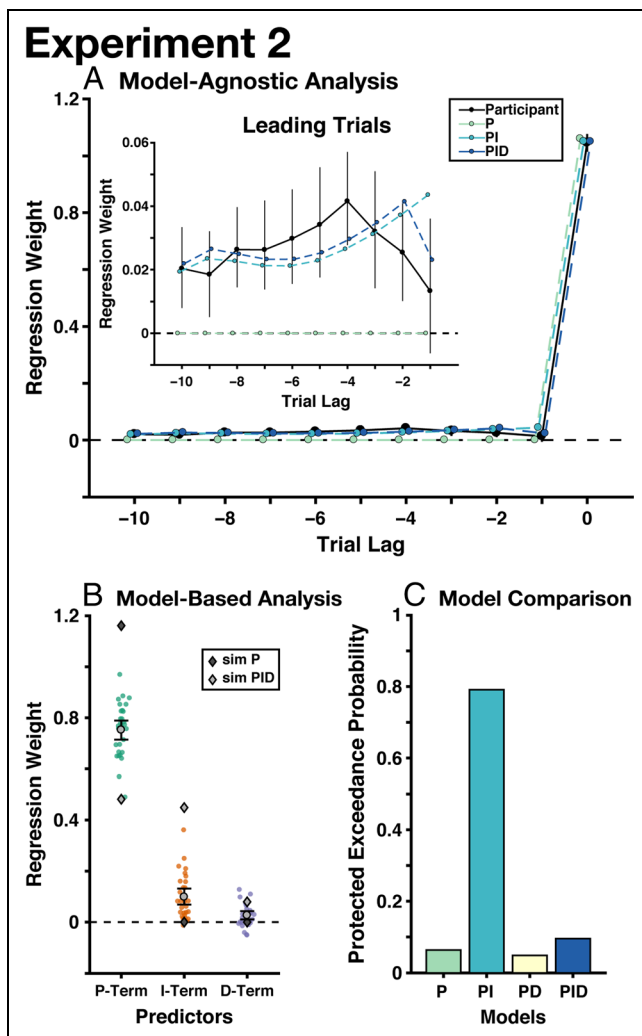


Figure 4. Experiment 2 results. (A) Participants adjusted their choices based on previous trials, unlike the predictions of a proportional controller (i.e., delta-rule model). (B) P-, I-, and D-Terms significantly predicted participants' updates. (C) The PI controller best explained participants' behavior. See Figure 3 for detailed graph legends. Error bars indicate mean and between-participant bootstrapped 95% confidence intervals.

backward in time than the previous-trial effects predicted by a simple derivative.

PID Model Fit

The best-performing PID gains were significantly different from zero (mean [SD] gain: $K_P = 0.48$ [0.13], $K_I = 0.45$ [0.085], $K_D = 0.080$ [0.066], $\lambda = 0.84$ [0.033]; all $p_s \leq 10^{-5}$; Figure 4B). Our PID model accounted for most of the variance in participants' updates (median $R^2 = .84$), and the parameters for the P-, I-, and D-Terms were all significantly different from zero (mean [SD] standardized beta: $\beta_P = 0.75$ [0.10], $p \leq 10^{-5}$; $\beta_I = 0.099$ [0.088], $p \leq 10^{-5}$; $\beta_D = 0.026$ [0.045], $p = .0034$; sign-randomization test; Figure 4B). The group level λ was 0.87. Participants' parameters were similar to the ideal

PID controller, although they had a greater reliance on proportional control and a lesser reliance on integral control than was optimal.

Bayesian model selection favored the PI model ($PXP_P = 0.064$, $PXP_{PI} = 0.79$, $PXP_{PD} = 0.049$, $PXP_{PID} = 0.095$; Figure 4C), although there was a moderate likelihood that models did not differ in their fit ($BOR = 0.20$). This was mostly due to the similarity in likelihood between the PI and PID models (excluding PID: $BOR < 0.001$, $PXP_{PI} > 0.99$). Therefore, our model selection supports the interpretation that PI control explains behavior better than the delta-rule model.

Discussion

Using a novel variant of a change-point task, we provide strong evidence that the PI control model can usefully describe participants' predictions. Our model-free analysis showed that participants incorporated previous errors in their adjustments, an observation incompatible with proportional control, but predicted by our PI and PID models. We also found that all of the PID terms were again significant predictors of participants' updates and that they were qualitatively similar to the gains of an ideal PID controller. Building on Experiment 1, this experiment provided strong evidence that PI was the best-fitting model. These data further support a role for control processes that extend beyond immediate errors.

These first two experiments have provided promising evidence that the PID framework predicts adaptive learning better than the classical delta-rule model. A striking feature of these experiments is that participants had very different estimated control gains across the two experiments, consistent with the differential gains of the best-performing PID agents. These differences suggest that participants may set their control gains in a context-specific manner, although at an unknown time scale. Popular delta-rule models have suggested that participants may in fact rapidly adjust their control gains in response to changes to the local context (Pearce & Hall, 1980). This prompted us to develop a third experiment, to replicate our results from Experiment 2 and test whether participants can adaptively adapt PID gains to their local contexts.

EXPERIMENT 3

In Experiment 3, we sought to replicate the findings from Experiment 2, while at the same time manipulating the incentives for performing accurately on the task. We additionally sought to examine three factors that might influence the weights that individuals might place on each of the PID terms over the course of an experiment.

First, we examined the influence of surprise (absolute error) on these control weights, given classic findings that such surprise signals modulate learning, indicating the degree to which the environment has been learned

(Pearce & Hall, 1980; see also McGuire et al., 2014; O'Reilly et al., 2013; Hayden, Heilbronner, Pearson, & Platt, 2011; Nassar et al., 2010). Second, given evidence that learning can be influenced by uncertainty over recent feedback (Nassar et al., 2010; Courville et al., 2006; Yu & Dayan, 2005) or related estimates of volatility (Behrens et al., 2007), we examined how PID gains were influenced by an index of the outcome entropy over the past several trials. This measure of uncertainty indexes both expected uncertainty (the variance in the generative distribution) and unexpected uncertainty (changes in the mean of the generative distribution, i.e., ramps), the latter of which is more dominant in our tasks.

We also examined the influence of reward on PID gains, given previous evidence that these can impact learning in a dissociable fashion from surprise or uncertainty alone (McGuire et al., 2014) and, more generally, that rewards may compensate for the costs of effortful control policies (Kool, Gershman, & Cushman, 2017; Manohar et al., 2015; Padmala & Pessoa, 2011; Hayden, Pearson, & Platt, 2009), including learning in particular (Shenhav, Botvinick, & Cohen, 2013; Hayden et al., 2009). For example, this could occur if integrating feedback utilizes domain-general working memory processes (Collins & Frank, 2012, 2018). Importantly, Experiments 1 and 3 were designed to de-confound reward from errors, providing us the ability to measure their influences on PID gains separately from one another and from our measure of uncertainty. In Experiment 1, performance-dependent rewards were given on a random subset of interleaved trials, whereas in Experiment 3, rewards were a nonlinear function of error that changed over time. These measures allowed us to distinguish the independent effects of surprise (absolute error) and reward on learning. For example, participants may have been motivated to perform accurately, and insofar as this motivation is further enhanced by reward, our analysis should be able to dissociate this motivation from other outcomes of error (e.g., surprise).

Finally, we compared our PID model against a popular model of adaptive learning, the Kalman filter (Kording, Tenenbaum, & Shadmehr, 2007; Kakade & Dayan, 2002; Kalman, 1960). This model performs state estimation using a delta-rule algorithm with an uncertainty-weighted learning rate. Previous experiments have found that it is a good model of behavior and it is based on the same principles that motivated the heuristic terms in our adaptive gain analysis.

Methods

Participants and Procedure

Forty-seven Brown University subject pool participants (32 women; mean age = 21.3 years, $SD = 4.07$ years) performed a rewarded supervised learning task (without monetary compensation). Apart from the reward manipulation, the structure of this task was similar to Experi-

ment 2. On each trial, the reward magnitude depended on the accuracy of the participant's guess (i.e., the absolute error between guess and outcome location; Figure 2C). These rewards decreased exponentially with increasing error magnitude. To de-correlate rewards and errors and to vary overall motivation to perform the task, we adjusted the steepness (mean) of this exponential (gamma) function over trials, resampling one of four possible means (1, 1.5, 2.5, and 4.5) at random time points, chosen with a flat hazard rate of 0.20 across all trials (Figure 2C, right). We instructed participants that these different levels of steepness defined four "reward modes." The reward mode for a given trial was indicated by the color of the fixation cross (one of four colors from equally spaced locations on a heat colormap). The input (errors) to these reward functions were divided by 3.5 to approximately match the reward that these functions returned at participants' mean performance level in Experiment 2.

Participants completed 50 training trials, followed by six blocks of 75 trials. On each trial, participants had up to 5 sec to make their guess, feedback was presented for 2.5 sec, and then the reward mode for the next trial was displayed during an ITI that was drawn from $U(1, 1.5)$. At the end of each block, participants were shown the mean reward earned during that block. Our final analysis excluded any trials where participants did not move their cursor to the edge of the circle (0.07% of the total trials). The lagged and trial-wise regression analyses were performed as described in Experiments 1 and 2.

Gain Modulation Analysis

To examine the influence of reward (Experiments 1 and 3; $n = 77$), errors (Experiments 1–3; $n = 106$), and outcome entropy (Experiments 1–3; $n = 106$) on the gains of the PID terms, we reran our PID regression analysis, including interaction terms for each type of gain modulation. In Experiment 1, the reward modulator consisted of binary reward feedback that was given on a random subset of trials, conditional on participants' error being within a prespecified threshold. This feedback was not correlated with absolute error on the task. In Experiment 3, the reward modulator was the number of points that participants received on each trial, which was a time-varying nonlinear function of absolute error (see procedure described above). In this task, participants received both error and reward feedback on every trial. Absolute error was correlated with the reward (median $r = -.68$); however, Belsley collinearity diagnostics (Belsley, Kuh, & Welsch, 1980) indicated that the collinearity between absolute error and reward was below standard tolerances, suggesting that our regression would be able to assess the independent contributions of each factor. In all three experiments, the error modulator was the absolute prediction error. Outcome entropy was defined as the natural logarithm of the outcome sample standard deviation over

the current and 10 previous trials within each block (with a truncated window for the first 10 trials in each block).

A robust regression (bisquare weighted) was run for every participant in every experiment, excluding the reward modulator for Experiment 2. The regression model included all main effects as well as the interactions between the PID terms and gain modulators [$u \sim 1 + (P + I + D) \times \text{Reward} + (P + I + D) \times \text{Absolute Error} + (P + I + D) \times \text{Outcome Entropy}$]. We mean-centered betas within their respective experiment and then re-centered the betas on their grand mean, removing between-experiment variance (Cousineau, 2005).

Kalman Filter Analysis

Our Kalman filter analysis was based on the algorithm used in Kording et al. (2007), building off the code that accompanied their publication. This Kalman filter estimated the likelihood of different states using an uncertainty-weighted delta-rule algorithm. Each state was a differential equation that defined a random walk over a specific time scale (i.e., slowly or quickly changing outcome locations). See Kording et al. (2007) for a complete description of this algorithm. Although the Kalman filter is not optimized for our task, given that the outcomes were not generated from a random walk, it has nevertheless proved to be a good model of behavior in previous experiments that used a random walk generative function (e.g., Gershman, 2015; Kording et al., 2007; Daw, O’Doherty, Dayan, Seymour, & Dolan, 2006; Kakade & Dayan, 2002).

Following Kording et al. (2007), states were defined as 30 diffusion time scales logarithmically spaced between two trials and the length of the experiment. We fit state noise parameters for each participant using restricted maximum likelihood estimation (MATLAB’s `fmincon`). The initial mean was set to the first outcome, and the initial covariance was set to a small variance constant (10^{-4}). As in our PID analysis, we fit the Kalman filter’s parameters so as to minimize the difference between its prediction updates and each participant’s prediction updates, based on participants’ errors on each trial (i.e., one-step look ahead).

We also compared the PI model against a variant of the Kalman filter that is less commonly used to describe adaptive behaviors but was better suited for our experiment. This position-velocity Kalman filter tracks randomly drifting changes in both the position (x) and velocity (\dot{x}) of the outcome locations:

$$\mathbf{x}_{t+1} = \mathbf{F}\mathbf{x}_t + \mathcal{N}(0, Q)$$

$$\mathbf{F}\mathbf{x}_t = \begin{pmatrix} 1 & \tau \\ 0 & v \end{pmatrix} \begin{pmatrix} x_t \\ \dot{x}_t \end{pmatrix}$$

$$Q = \begin{pmatrix} \frac{1}{4}\tau^4 & \frac{1}{2}\tau^3 \\ \frac{1}{2}\tau^3 & \tau^2 \end{pmatrix} \sigma_\alpha^2$$

We used restricted maximum likelihood estimation to fit participant-specific velocity decay (v), time delay (τ), and state noise (σ_α^2) parameters to participants’ updates using the same one-step look ahead procedure described above. The initial mean was set to the first outcome, and the initial covariance was set to the variance in outcome position and velocity, averaged across participants.

Results

Model-agnostic Analysis

Regressing the current and 10 leading errors onto the current update (see Methods under Experiment 1), we replicated the observation that participants were influenced by past errors (summed betas: mean = 0.23, $SD = 0.21$, $p \leq 10^{-5}$; see Figure 5A). Our model-generated behavior

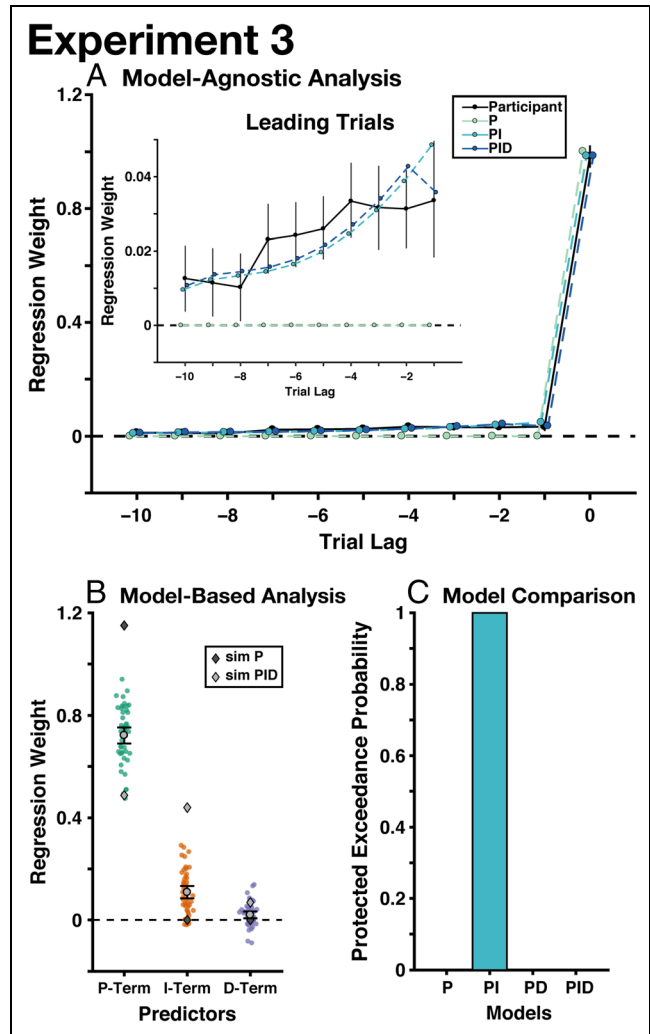


Figure 5. Experiment 3 results. (A) Participants adjusted their choices based on previous errors, unlike the predictions of a proportional controller (i.e., delta-rule model). (B) P, I, and D control significantly predicted participants’ updates. (C) The PI controller best explained participants’ behavior. See Figure 3 for detailed graph legends. Error bars indicate mean and between-participant bootstrapped 95% confidence intervals.

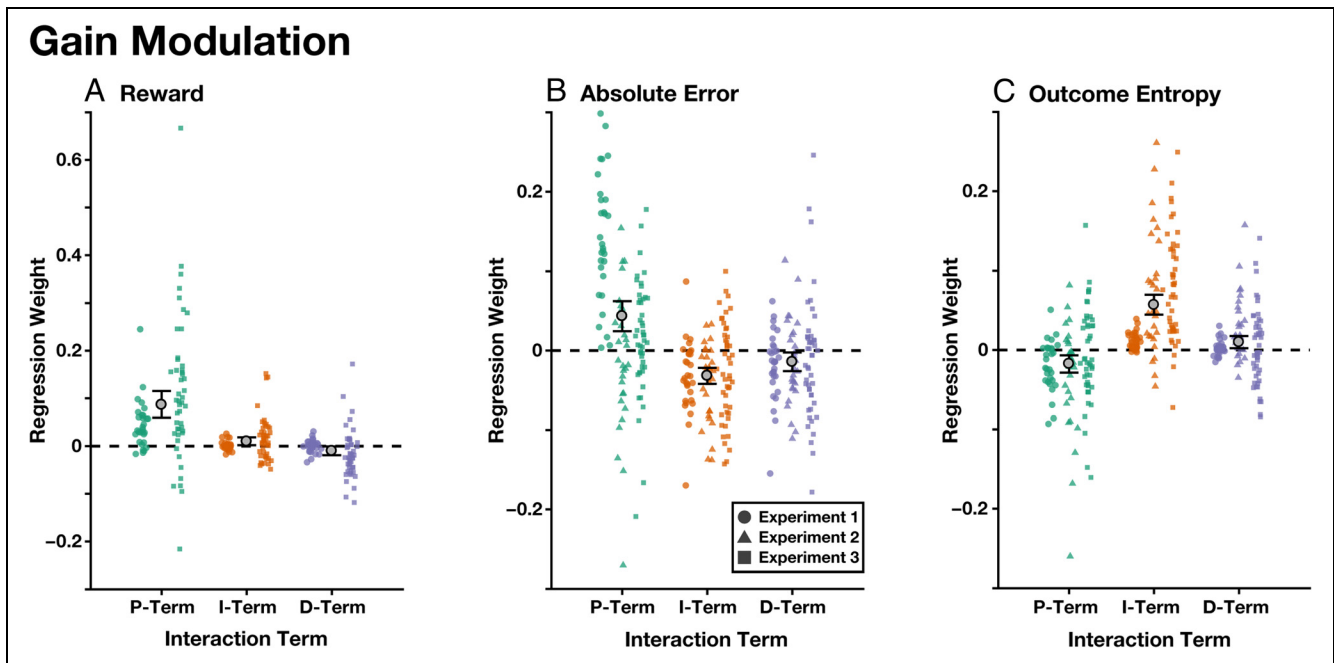


Figure 6. Gain modulation. Trial-wise reward (A), absolute error (B), and outcome entropy (C) significantly interacted with all three PID terms. All models included PID terms as main effects. Colored shapes indicate individual participant's standardized betas in each experiment (see legend). Error bars indicate mean and between-participant bootstrapped 95% confidence intervals, uncorrected for between-experiment variance.

again showed that the delta-rule model categorically fails to capture the influence of leading errors. Unlike Experiment 2, here, we found that the weighting of previous errors was best fit as a linear decay from the current trials, resembling PI control (mean [SD] trend beta: linear = 0.0087 [0.016], $p < 10^{-4}$; quadratic = -0.002 [0.016], $p = .39$; sign-randomization test). This discrepancy from Experiment 2 may be because rewards in Experiment 3 were highly dependent on accuracy. This may have biased participants more toward integral control (which favors accuracy) and away from derivative control (which favors stability; Aström & Murray, 2008).

PID Model Fit

Replicating Experiments 1 and 2, we found that our standard PID model accounted for most of the variance in participants' updates (median $R^2 = .81$). The parameters for the P-, I-, and D-Terms were all significantly different from zero (mean [SD] standardized beta: $\beta_P = 0.72$ [0.11], $p \leq 10^{-5}$; $\beta_I = 0.11$ [0.086], $p \leq 10^{-5}$; $\beta_D = 0.020$ [0.047], $p = .006$; Figure 5B). The group level λ was 0.8016. Participants' estimated gains were similar to the ideal PID controller, but they overweighted proportional control and underweighted integral control. We found that there were likely differences between the model likelihoods (BOR < 0.001) and that Bayesian model selection strongly favored the PI model ($PXP_{PI} > 0.99$) over the alternate models (all other $PXPs < 10^{-4}$; Figure 5C).

Gain Modulation

We examined the independent influence of rewards, absolute error, and outcome entropy in modulating the PID gains across our three experiments. We found that all three modulators significantly interacted with the P-, I-, and D-Terms, but in distinct ways (Figure 6): Increased reward led to increased P and I gains and a decreased D gain (Figure 6A; mean [SD] interaction beta: $\beta_{P:\text{reward}} = 0.086$ [0.12], $p \leq 10^{-5}$; $\beta_{I:\text{reward}} = 0.0098$ [0.036], $p = .016$; $\beta_{D:\text{reward}} = -0.010$ [0.040], $p = .032$; sign-randomization test). Increased absolute error led to an increased P gain and decreased I and D gains (Figure 6B; mean [SD] interaction beta: $\beta_{P:\text{error}} = 0.043$ [0.078], $p \leq 10^{-5}$; $\beta_{I:\text{error}} = -0.032$ [0.053], $p \leq 10^{-5}$; $\beta_{D:\text{error}} = -0.014$ [0.062], $p = .019$); increased outcome entropy led to a decreased P gain and increased I and D gains (Figure 6C; mean [SD] interaction beta: $\beta_{P:\text{entropy}} = -0.018$ [0.056], $p = .0016$; $\beta_{I:\text{entropy}} = 0.057$ [0.059], $p \leq 10^{-5}$; $\beta_{D:\text{entropy}} = 0.0098$ [0.039], $p = .011$).

These interactions were robust to several quality checks. First, all effects remained significant when we corrected for multiple comparisons using the Holm-Bonferroni procedure (Holm, 1979). Given the presence of outliers, we also tested our effects using a robust Wilcoxon signed-rank test (Wilcoxon, 1945), finding that all interactions remained significant ($ps \leq .014$). Finally, we also found that all interactions remained significant when we did not remove between-experiment variance ($ps \leq .035$; Figure 6 depicts participants' raw interaction betas).

Kalman Filter Analysis

We fit the Kalman filter to participants' behavior in both Experiments 2 and 3, finding that Bayesian model selection strongly favored the PI control model over the standard Kalman filter (pooling across experiments; $PXP_{PI} > 0.99$, $BOR < 10^{-14}$; Figure 7A). Using our lagged regression analysis approach, we also found that the standard Kalman filter's updates depended on previous errors in a qualitatively different way from participant updates. Unlike participants, the Kalman filter placed negative weights on errors made in earlier trials (Figure 7B). We also found that the standard Kalman filter also performed especially poorly when outcomes changed over time (i.e., at different outcome velocities), whereas participants and the PI model were able to accommodate such changes in outcomes (Figure 7C).

We also compared the PI control model against a Kalman filter model that tracked the position and velocity of outcomes over time. Despite the additional complexity of this model, we found that the PI model fit similarly well ($PXP_{PI} = 0.63$, $BOR = 0.70$). These models were identifiable, as we could accurately recover the correct model when either of them generated behavior, suggesting that they offer dissociable explanations of participants' behaviors. Interestingly, we found that participants' velocity estimates strongly decayed over time (mean $v = 0.32$) and that this parameter strongly correlated with participants' integral gain ($r = .79$, $p < 10^{-16}$), suggesting that these terms might serve complementary computational roles. Collectively, these results show that the PI model offers a more parsimonious account of participants' behavior than a complex, task-informed inferential model.

Discussion

In Experiment 3, we found confirmatory evidence that the PI model accurately describes participants' predictions and that participants adjust their weighting of different PID terms based on trial-wise task dynamics. We found that each PID term was uniquely sensitive to changes in reward, absolute error, and outcome entropy, extending previous observations of the role of these modulators on proportional control and providing further evidence that the PID terms represent distinct control processes. We also found that the PI model offered a better explanation of behavior than the standard Kalman filter and performed similarly to a specialized Kalman filter variant, demonstrating that the PI model is as powerful as more complex models based on explicit state space representations.

When participants received larger rewards, they modulated their gains in a way that is consistent with a preference for accuracy (P- and I-Terms) over stability (D-Term; Ang, Chong, & Li, 2005), potentially indicating an exploitive strategy for the high-reward environments (Kovach et al., 2012). Although participants' proportional gain was already larger than the best-performing gain, this may reflect the unique role of reward modulation, when controlling for the environmental changes (e.g., entropy) that make a high proportional gain less desirable. Another alternative is that the P- and/or I-Terms are effortful to implement, with rewards "paying the cost" of these control policies (Kool et al., 2017; Manohar et al., 2015). Further work will be necessary to dissociate the role of salience and motivation on reward-modulated gain adjustments.

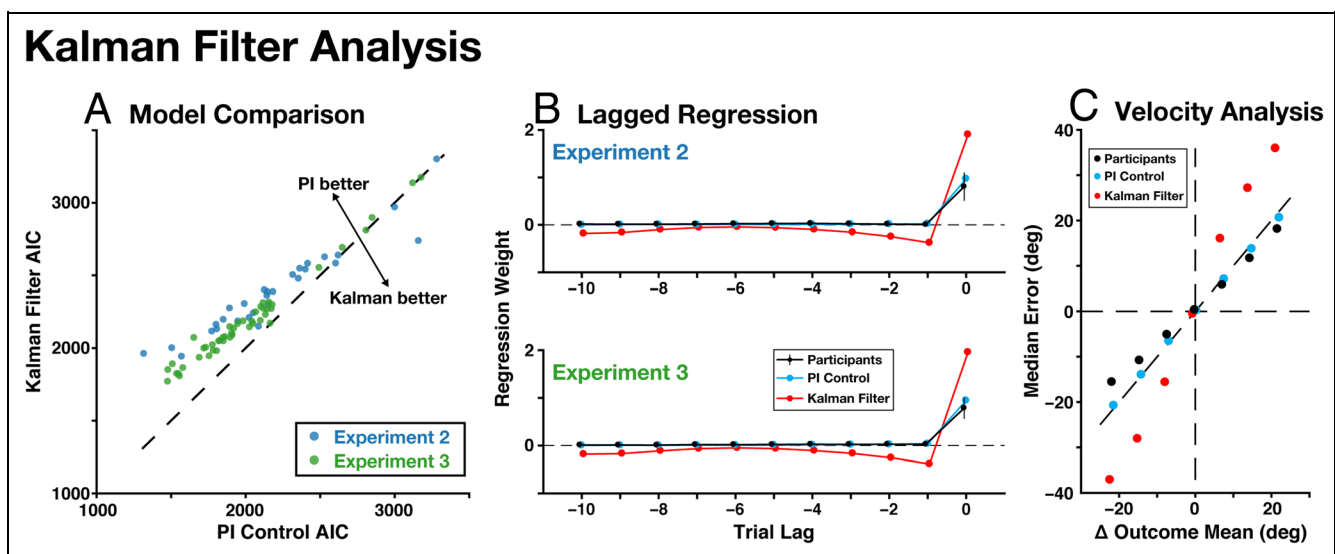


Figure 7. Kalman filter analysis. The PI control model fitted participants' behavior better than a standard Kalman filter model. (A) For most participants (colored dots), a complexity-corrected measure of fit (AIC) was better for the PI control model than the Kalman filter model. (B) Unlike the PI control model, the Kalman filter model poorly resembled how lagged errors influenced participants' updates (compare with A in Figures 4–5). (C) Unlike the PI control model, the Kalman filter did not resemble participants' accuracy when the outcome distribution changed over time (pooled across Experiments 2 and 3).

In response to absolute errors (i.e., surprise), participants increased their immediate adjustment (P-Term) and relied less on previous feedback (I- and D-Terms). This is consistent with the idea that large errors may indicate changes in the environment (Nassar et al., 2010; Pearce & Hall, 1980) and with filtering mechanisms in industrial PID control that improve robustness by limiting the long-term influence of noisy samples (Ang et al., 2005).

Although outcome entropy, and uncertainty, has traditionally been thought to increase the gain on proportional control (Nassar et al., 2010; Behrens et al., 2007; Courville et al., 2006), in our experiment, the P-Term was decreased and the I- and D-Terms were instead increased. Interestingly, when we instead implement this gain modulation in a P-only model, we do find that outcome entropy increases the gain of the P-Term (data not shown). Unlike previous experiments studying uncertainty, environmental change in Experiments 2 and 3 required tracking gradually changing outcomes, which accounted for most of the outcome entropy and for which integral and derivative control are particularly useful (Wittmann et al., 2016; Kovach et al., 2012). In Experiment 1, where these gradual changes were not present, we found that uncertainty after a change point increased control gains (McGuire et al., 2014), which may be reflected here by integrating over the trials since the change point.

We found that the PI model explained participants' behavior better than a standard Kalman filter (a powerful model of adaptive learning; Kording et al., 2007) and that the Kalman filter failed to capture participants' use of feedback history. This difference was largely due to the ability of the integral term to track ramping changes in the environment, epochs that were poorly accounted for by the Kalman filter. Interestingly, the Kalman filter's updates were negatively correlated with errors made on earlier trials (when controlling for the influence of the current error). We believe that this is due to the short diffusion time scales, which were updated the fastest (they were set to the highest state noise, as in Kording et al., 2007) and define the difference between current and recent trials. We found that the lagged influence of recent trials was more strongly negative for shorter time scales (data not shown).

We also compared the PI model against a position-velocity Kalman filter that tracked both the position and velocity of outcomes, finding that these models fit similarly well. There was a strong relationship between this Kalman filter's velocity term and the PI controller's integral term, suggesting that participants could use integral control to track ramping changes in the environment. This position-velocity Kalman filter has received little attention in the learning literature and warrants further investigation; however, it currently offers a less parsimonious explanation of behavior than PI control due to its greater computational complexity and its requirement

for explicit state representations. Although both of these Kalman filters did not offer better models than PI control, the Kalman filter embodies the same principles as our adaptive gain analysis: Control gains should be adaptive and depend on factors like environmental stability.

GENERAL DISCUSSION

Across three experiments, we found that the PI model successfully captured participants' prediction updating in a stochastic environment. By incorporating a richer model of control monitoring and adjustment, the PI controller was able to account for ways in which performance in such environments deviates from predictions of standard-error-driven (delta-rule) learning models. We also replicated and extended previous findings showing that learning parameters themselves are modulated by environmental signals (e.g., reward) and extended these findings to show that signals related to the magnitude of reward, error, and outcome entropy can differentially affect the gains on the PID model parameters.

Our findings suggest that PI control offers a good account of behavior across two fairly different task environments. Indeed, although we found that normative PID gains differed substantially between Experiment 1 (discrete transitions) and Experiments 2–3 (gradual transitions), participants' behavior continued to qualitatively match the behavior predicted by this normative controller across studies, in each case matching the sign and rank order of the best-performing control gain. This suggests that these gains adapted to the specific environment that participants were acting in. Specifically, when outcomes were prone to shift sharply and dramatically (Experiment 1), participants tended to rely less on history-dependent control processes like integral and derivative control, especially on trials in which large errors may have indicated a state shift.

Although we have focused our discussion of the PID controller on all three of its components, in industrial settings, the D-Term is often given the lowest gain or not included (Aström & Murray, 2008), as it is highly sensitive to noise. Accordingly, our own data supported little to no role for the derivative term in the current experiments, both normatively and in our model fits to participants' behavior. Although the derivative control term was significant in all of the experiments and interacted with the absolute error, it did not account for sufficient variance to outweigh complexity penalties in model comparison. This may have been compounded by the fact that the derivative term was negatively modulated by absolute error, which may have caused it to explain less of the variance on trials where there were large updates. Although the outcomes in Experiments 2 and 3 were designed to differentiate PID control from the delta-rule model, they were not designed to specifically detect derivative control. Future research should investigate cases where derivative control is especially beneficial

for good performance. Because derivative control provides high-frequency compensation to improve responsiveness, it may be the case that derivative control is generally poorly suited for tasks that depend on intertrial adjustments and favor accuracy over speed. Relative to Experiment 2, Experiment 3 emphasized accuracy through its reward structure and deemphasized responsiveness because of its longer trial length. Although there were several differences between these experiments, these factors may have contributed to the differences in derivative control between these experiments.

Some of the most promising evidence for derivative control was that, in Experiment 2, participants down-weighted recent errors (from $t-3$ and $t-1$) relative to what would be expected by error integration alone. Although basic derivative control would only compare the current and previous errors, participants' behavior resembles a common practice in control engineering to low-pass filter the derivative term to improve robustness (Ang et al., 2005). The discrepancy between the observed nonlinear influence of previous errors (predicted by the full PID model) and the model selection preference for the PI model may therefore be accounted for by alternative forms of derivative control.

We found that the PID terms depended on reward feedback (Experiments 1 and 3), absolute errors (Experiments 1–3), and outcome entropy (Experiments 1–3) on a trial-to-trial basis. Although there is substantial literature on how environmental factors should influence the standard delta-rule model, less is known on how these factors should affect PID gains. These modulation factors may offer insight as to how the control system sets different control gains, which in our experiment were fit to behavior. Although we have proposed speculative explanations for the role of each modulating factor, at a minimum, the unique pattern of interactions for each of the PID terms suggests that P, I, and D represent dissociable forms of control. Future experiment should examine the extent to which gain modulation depends on the structure of the task and environment, for instance, whether the task rewards consistency in addition to accuracy.

The PID model provides robust control without relying on an explicit model of the environment, offering a parsimonious explanation of participants' behavior. Although this model is not optimal (e.g., with respect to mean squared error), it offers an approximate solution without the computational demands of exactly modeling the nonlinear system dynamics (Motter, 2015). That said, there have been notable successes for algorithms that instead learn generative models of the environment (e.g., using Bayesian estimation) and can represent the uncertainty about upcoming choices (e.g., Franklin & Frank, 2015; Griffiths, Lieder, & Goodman, 2015; McGuire et al., 2014; Nassar et al., 2010; Daw, Niv, & Dayan, 2005; although see Duverne & Koechlin, 2017; Geana & Niv, 2014; Mathys, Daunizeau, Friston, & Stephan, 2011). To examine this possibility, we compared the PI

control model against the Kalman filter, a standard model for state estimation in the face of uncertainty. We found that the PI model better explained participants' behavior than a standard Kalman filter (Kording et al., 2007) and fitted comparably to a Kalman filter that was specialized for this experiment. In contrast to the Kalman filter, the PID controller offers a general control process that can parsimoniously account for participants' behavior with minimal knowledge about the task structure. These benefits would likely be compounded by the complex dynamics of natural environments.

Despite these promising results, we would not rule out the possibility that participants rely on a combination of both model-free (e.g., PID) and model-based control (Kool, Cushman, & Gershman, forthcoming; Korn & Bach, 2018; Momennejad et al., 2017; Daw, Gershman, Seymour, Dayan, & Dolan, 2011; Gläscher, Daw, Dayan, & O'Doherty, 2010). Previous experiments have demonstrated the utility of model-based predictions for explaining participants' behavior in other environments, and participants can report confidence in their choices. Model-based control may serve to modulate the PID controller itself (e.g., to tune gain parameters or reset control processes; McGuire et al., 2014; Nassar et al., 2010; Behrens et al., 2007; Bouret & Sara, 2005); may be selectively engaged in environments that are stable, constrained, or familiar; and/or may trade off over different stages in learning (Denève et al., 2017).

Another promising feature of the PID model is that it offers a model of behavioral control that can be plausibly implemented by a neural system. There have been several neural network implementations of PID controllers in industrial engineering (e.g., Cong & Liang, 2009), with integral and derivative control implemented as positive and negative recurrent circuits, respectively. This simple architecture demonstrates the ease with which a neural system could develop PID control dynamics. Moreover, recent studies have found neuroscientific evidence that is broadly consistent with the predictions of such an architecture. For instance, Bernacchia and colleagues (2011) found that, in rhesus macaques' cingulate cortex and pFC, large populations of neurons encoded the history of trial-epoch-selective activity, likely including error-related responses (cf. Seo & Lee, 2007). Each of these regions contained equally sized populations of neurons that tracked either the exponentially weighted sum of recent trials or the difference between recent and previous trials, putative markers of integral and derivative control, respectively. Convergent data in humans found that fMRI activity in dorsal ACC parametrically tracked a recent history of prediction errors in a changing environment (Wittmann et al., 2016), again consistent with the operations of an integral-based controller. Accordingly, these authors found that incorporating integration into their behavioral model explained choices in their task better than the traditional delta-rule model. Although these findings provide evidence for neural signatures of feedback history (see also Seo & Lee, 2007; Kennerley et al., 2006)

and are consistent with the monitoring function of PID control, future experiments are needed to formally test for the neural correlates of this model.

These experiments together provide strong evidence for the viability of control theoretic models as mechanisms of human prediction updating in dynamic environments. This class of models has been highly influential in research on motor control, including the PID controller in particular (e.g., Kawato & Wolpert, 1998). Motor control models typically describe the rapid regulation of limb movements to produce trajectories that are fast, accurate, and robust. In contrast, participants in our experiments were not motivated to make fast or accurate trajectories and instead may have used an analogous control process to adapt their predictions from trial to trial. Control theoretic algorithms (like PID control) may be a domain-general class of neural functions, involved in a diverse array of cognitive processes (Pezzulo & Cisek, 2016; Powers, 1973; Ashby, 1956), including the cognitive control functions that have been suggested to operate using both classical (Botvinick et al., 2001) and optimal (Shenhav et al., 2013) control principles. The architecture of these executive control algorithms and the nature of the references that they regulate are important areas of further research.

Acknowledgments

The authors thank Kia Sadahiro and William McNelis for their assistance in data collection.

Reprint requests should be sent to Harrison Ritz, Brown University, Providence, RI 02912, or via e-mail: harrison_ritz@brown.edu.

Notes

1. Domain-specific “delta-rule” algorithms are common in many fields, such as the Rescorla–Wagner learning rule (Rescorla & Wagner, 1972) or a delta rule algorithm used in neural networks (Widrow & Hoff, 1960). In this article, we define the delta rule as a more general class of error-based learning rules in which adjustments are proportional to errors.

2. We chose AIC over the more conservative Bayesian information criterion (BIC) because model recovery found that BIC was overly conservative: Model selection using BIC did not prefer the full PID model when this model generated behavior (i.e., when PID was the ground truth). Although AIC is not the ideal fit metric for Bayesian model selection (as it is not an approximation of model likelihood), the development team for SPM’s Bayesian model selection protocol has justified using AIC as a legitimate alternative to BIC: “Though not originally motivated from a Bayesian perspective, model comparisons based on AIC are asymptotically equivalent to those based on Bayes factors (Akaike, 1973a), that is, AIC approximates the model evidence” (Penny, Stephan, Mechelli, & Friston, 2004, p. 1162; see also Rigoux et al., 2014; Penny, 2012).

REFERENCES

Aben, B., Verguts, T., & Van den Bussche, E. (2017). Beyond trial-by-trial adaptation: A quantification of the time scale of cognitive control. *Journal of Experimental Psychology: Human Perception and Performance*, *43*, 509–517.

- Akaike, H. (1983). Information measures and model selection. *Bulletin of the International Statistical Institute*, *50*, 277–291.
- Alexander, W. H., & Brown, J. W. (2015). Hierarchical error representation: A computational model of anterior cingulate and dorsolateral prefrontal cortex. *Neural Computation*, *27*, 2354–2410.
- Ang, K. H., Chong, G., & Li, Y. (2005). PID control system analysis, design, and technology. *IEEE Transactions on Control Systems Technology*, *13*, 559–576.
- Ashby, W. R. (1956). *An introduction to cybernetics*. Chapman and Hall.
- Aström, K. J., & Murray, R. M. (2008). *Feedback systems: An introduction for scientists and engineers*. Princeton University Press.
- Behrens, T. E., Woolrich, M. W., Walton, M. E., & Rushworth, M. F. (2007). Learning the value of information in an uncertain world. *Nature Neuroscience*, *10*, 1214–1221.
- Belsley, D. A., Kuh, E., & Welsch, R. E. (1980). *Regression diagnostics: Identifying influential data and sources of collinearity*. Hoboken, NJ: John Wiley & Sons.
- Bernacchia, A., Seo, H., Lee, D., & Wang, X. J. (2011). A reservoir of time constants for memory traces in cortical neurons. *Nature Neuroscience*, *14*, 366–372.
- Blais, C., & Bunge, S. (2010). Behavioral and neural evidence for item-specific performance monitoring. *Journal of Cognitive Neuroscience*, *22*, 2758–2767.
- Botvinick, M. M., Braver, T. S., Barch, D. M., Carter, C. S., & Cohen, J. D. (2001). Conflict monitoring and cognitive control. *Psychological Review*, *108*, 624.
- Bouret, S., & Sara, S. J. (2005). Network reset: A simplified overarching theory of locus coeruleus noradrenaline function. *Trends in Neurosciences*, *28*, 574–582.
- Bugg, J. M., & Crump, M. J. C. (2012). In support of a distinction between voluntary and stimulus-driven control: A review of the literature on proportion congruent effects. *Frontiers in Psychology*, *3*, 367.
- Carter, C. S., Macdonald, A. M., Botvinick, M., Ross, L. L., Stenger, V. A., Noll, D., et al. (2000). Parsing executive processes: Strategic vs. evaluative functions of the anterior cingulate cortex. *Proceedings of the National Academy of Sciences, U.S.A.*, *97*, 1944–1948.
- Collins, A. G., & Frank, M. J. (2012). How much of reinforcement learning is working memory, not reinforcement learning? A behavioral, computational, and neurogenetic analysis. *European Journal of Neuroscience*, *35*, 1024–1035.
- Collins, A. G., & Frank, M. J. (2018). Within- and across-trial dynamics of human EEG reveal cooperative interplay between reinforcement learning and working memory. *Proceedings of the National Academy of Sciences, U.S.A.*, *115*, 2502–2507.
- Cong, S., & Liang, Y. (2009). PID-like neural network nonlinear adaptive control for uncertain multivariable motion control systems. *IEEE Transactions on Industrial Electronics*, *56*, 3872–3879.
- Courville, A. C., Daw, N. D., & Touretzky, D. S. (2006). Bayesian theories of conditioning in a changing world. *Trends in Cognitive Sciences*, *10*, 294–300.
- Cousineau, D. (2005). Confidence intervals in within-subject designs: A simpler solution to Loftus and Masson’s method. *Tutorials in Quantitative Methods for Psychology*, *1*, 42–45.
- Daw, N. D., Gershman, S. J., Seymour, B., Dayan, P., & Dolan, R. J. (2011). Model-based influences on humans’ choices and striatal prediction errors. *Neuron*, *69*, 1204–1215.
- Daw, N. D., Niv, Y., & Dayan, P. (2005). Uncertainty-based competition between prefrontal and dorsolateral striatal systems for behavioral control. *Nature Neuroscience*, *8*, 1704–1711.

- Daw, N. D., O'Doherty, J. P., Dayan, P., Seymour, B., & Dolan, R. J. (2006). Cortical substrates for exploratory decisions in humans. *Nature*, *441*, 876.
- Denève, S., Alemi, A., & Bourdoukan, R. (2017). The brain as an efficient and robust adaptive learner. *Neuron*, *94*, 969–977.
- Duvernois, S., & Koechlin, E. (2017). Rewards and cognitive control in the human prefrontal cortex. *Cerebral Cortex*, *27*, 5024–5039.
- Franklin, G. F., Powell, J. D., & Emami-Naeini, A. (1994). *Feedback control of dynamic systems (Vol. 3)*. Reading, MA: Addison-Wesley.
- Franklin, N. T., & Frank, M. J. (2015). A cholinergic feedback circuit to regulate striatal population uncertainty and optimize reinforcement learning. *eLife*, *4*, e12029.
- Garrison, J., Erdeniz, B., & Done, J. (2013). Prediction error in reinforcement learning: A meta-analysis of neuroimaging studies. *Neuroscience & Biobehavioral Reviews*, *37*, 1297–1310.
- Geana, A., & Niv, Y. (2014). *Causal model comparison shows that human representation learning is not Bayesian*. Paper presented at the Cold Spring Harbor Symposia on Quantitative Biology.
- Gelman, A., Meng, X.-L., & Stern, H. (1996). Posterior predictive assessment of model fitness via realized discrepancies. *Statistica Sinica*, *6*, 733–760.
- Gershman, S. J. (2015). A unifying probabilistic view of associative learning. *PLoS Computational Biology*, *11*, e1004567.
- Gläscher, J., Daw, N., Dayan, P., & O'Doherty, J. P. (2010). States versus rewards: Dissociable neural prediction error signals underlying model-based and model-free reinforcement learning. *Neuron*, *66*, 585–595.
- Gratton, G., Coles, M. G., & Donchin, E. (1992). Optimizing the use of information: Strategic control of activation of responses. *Journal of Experimental Psychology: General*, *121*, 480.
- Griffiths, T. L., Lieder, F., & Goodman, N. D. (2015). Rational use of cognitive resources: Levels of analysis between the computational and the algorithmic. *Topics in Cognitive Science*, *7*, 217–229.
- Hayden, B. Y., Heilbronner, S. R., Pearson, J. M., & Platt, M. L. (2011). Surprise signals in anterior cingulate cortex: Neuronal encoding of unsigned reward prediction errors driving adjustment in behavior. *Journal of Neuroscience*, *31*, 4178–4187.
- Hayden, B. Y., Pearson, J. M., & Platt, M. L. (2009). Fictive reward signals in the anterior cingulate cortex. *Science*, *324*, 948–950.
- Holm, S. (1979). A simple sequentially rejective multiple test procedure. *Scandinavian Journal of Statistics*, *6*, 65–70.
- Hunt, L. T., & Hayden, B. Y. (2017). A distributed, hierarchical and recurrent framework for reward-based choice. *Nature Reviews Neuroscience*, *18*, 172–182.
- Ito, S., Stuphorn, V., Brown, J. W., & Schall, J. D. (2003). Performance monitoring by the anterior cingulate cortex during saccade countermanding. *Science*, *302*, 120–122.
- Jiang, J., Beck, J., Heller, K., & Egner, T. (2015). An insula–frontostriatal network mediates flexible cognitive control by adaptively predicting changing control demands. *Nature Communications*, *6*, 8165.
- Kakade, S., & Dayan, P. (2002). Acquisition and extinction in autoshaping. *Psychological Review*, *109*, 533.
- Kalman, R. E. (1960). A new approach to linear filtering and prediction problems. *Journal of Basic Engineering*, *82*, 35–45.
- Karlsson, M. P., Tervo, D. G., & Karpova, A. Y. (2012). Network resets in medial prefrontal cortex mark the onset of behavioral uncertainty. *Science*, *338*, 135–139.
- Kawato, M., & Wolpert, D. (1998). Internal models for motor control. *Sensory Guidance of Movement*, *218*, 291–307.
- Kennerley, S. W., Walton, M. E., Behrens, T. E., Buckley, M. J., & Rushworth, M. F. (2006). Optimal decision making and the anterior cingulate cortex. *Nature Neuroscience*, *9*, 940–947.
- Kool, W., Cushman, F. A., & Gershman, S. J. (forthcoming). Competition and cooperation between multiple reinforcement learning systems. In *Goal-directed decision making: Computations and neural circuits*. New York: Elsevier.
- Kool, W., Gershman, S. J., & Cushman, F. A. (2017). Cost-benefit arbitration between multiple reinforcement-learning systems. *Psychological Science*, *28*, 1321–1333.
- Kording, K. P., Tenenbaum, J. B., & Shadmehr, R. (2007). The dynamics of memory as a consequence of optimal adaptation to a changing body. *Nature Neuroscience*, *10*, 779.
- Korn, C. W., & Bach, D. R. (2018). Heuristic and optimal policy computations in the human brain during sequential decision-making. *Nature Communications*, *9*, 325.
- Kovach, C. K., Daw, N. D., Rudrauf, D., Tranel, D., O'Doherty, J. P., & Adolphs, R. (2012). Anterior prefrontal cortex contributes to action selection through tracking of recent reward trends. *Journal of Neuroscience*, *32*, 8434–8442.
- Laming, D. R. J. (1968). *Information theory of choice-reaction times*. Oxford, UK: Academic Press.
- Lau, B., & Glimcher, P. W. (2005). Dynamic response-by-response models of matching behavior in rhesus monkeys. *Journal of the Experimental Analysis of Behavior*, *84*, 555–579.
- Logan, G. D., & Zbrodoff, N. J. (1979). When it helps to be misled: Facilitative effects of increasing the frequency of conflicting stimuli in a Stroop-like task. *Memory & Cognition*, *7*, 166–174.
- Manohar, S. G., Chong, T. T. J., Apps, M. A. J., Batla, A., Stamelou, M., Jarman, P. R., et al. (2015). Reward pays the cost of noise reduction in motor and cognitive control. *Current Biology*, *25*, 1707–1716.
- Mathys, C., Daunizeau, J., Friston, K. J., & Stephan, K. E. (2011). A Bayesian foundation for individual learning under uncertainty. *Frontiers in Human Neuroscience*, *5*, 39.
- Matsumoto, M., Matsumoto, K., Abe, H., & Tanaka, K. (2007). Medial prefrontal cell activity signaling prediction errors of action values. *Nature Neuroscience*, *10*, 647.
- Maxwell, J. C. (1868). I. On governors. *Proceedings of the Royal Society of London*, *16*, 270–283.
- McGuire, J. T., Nassar, M. R., Gold, J. I., & Kable, J. W. (2014). Functionally dissociable influences on learning rate in a dynamic environment. *Neuron*, *84*, 870–881.
- Mirenovic, J., & Schultz, W. (1994). Importance of unpredictability for reward responses in primate dopamine neurons. *Journal of Neurophysiology*, *72*, 1024–1027.
- Momennejad, I., Russek, E. M., Cheong, J. H., Botvinick, M. M., Daw, N. D., & Gershman, S. J. (2017). The successor representation in human reinforcement learning. *Nature Human Behaviour*, *1*, 680.
- Motter, A. E. (2015). Network control. *Chaos*, *25*, 097621.
- Nassar, M. R., & Frank, M. J. (2016). Taming the beast: Extracting generalizable knowledge from computational models of cognition. *Current Opinion in Behavioral Sciences*, *11(Suppl. C)*, 49–54.
- Nassar, M. R., Wilson, R. C., Heasly, B., & Gold, J. I. (2010). An approximately Bayesian delta-rule model explains the dynamics of belief updating in a changing environment. *Journal of Neuroscience*, *30*, 12366–12378.
- Niki, H., & Watanabe, M. (1979). Prefrontal and cingulate unit activity during timing behavior in the monkey. *Brain Research*, *171*, 213–224.

- O'Reilly, J. X., Schuffelgen, U., Cuell, S. F., Behrens, T. E., Mars, R. B., & Rushworth, M. F. (2013). Dissociable effects of surprise and model update in parietal and anterior cingulate cortex. *Proceedings of the National Academy of Sciences, U.S.A.*, *110*, E3660–E3669.
- Padmala, S., & Pessoa, L. (2011). Reward reduces conflict by enhancing attentional control and biasing visual cortical processing. *Journal of Cognitive Neuroscience*, *23*, 3419–3432.
- Pearce, J. M., & Hall, G. (1980). A model for Pavlovian learning: Variations in the effectiveness of conditioned but not of unconditioned stimuli. *Psychological Review*, *87*, 532.
- Penny, W. D. (2012). Comparing dynamic causal models using AIC, BIC and free energy. *Neuroimage*, *59*, 319–330.
- Penny, W. D., Stephan, K. E., Mechelli, A., & Friston, K. J. (2004). Comparing dynamic causal models. *Neuroimage*, *22*, 1157–1172.
- Pezzulo, G., & Cisek, P. (2016). Navigating the affordance landscape: Feedback control as a process model of behavior and cognition. *Trends in Cognitive Sciences*, *20*, 414–424.
- Powers, W. T. (1973). *Behavior: The control of perception*. Chicago: Aldine.
- Rabbitt, P. (1966). Errors and error correction in choice-response tasks. *Journal of Experimental Psychology*, *71*, 264.
- Rescorla, R. A., & Wagner, A. W. (1972). A theory of Pavlovian conditioning: Variations in the effectiveness of reinforcement and nonreinforcement. In A. H. Black & W. F. Prokasy (Eds.), *Classical conditioning II: Current research and theory* (pp. 64–99). New York: Appleton-Century-Crofts.
- Rigoux, L., Stephan, K. E., Friston, K. J., & Daunizeau, J. (2014). Bayesian model selection for group studies—Revisited. *Neuroimage*, *84*, 971–985.
- Rumelhart, D. E., Hinton, G. E., & Williams, R. J. (1986). Learning representations by back-propagating errors. *Nature*, *323*, 533.
- Seo, H., & Lee, D. (2007). Temporal filtering of reward signals in the dorsal anterior cingulate cortex during a mixed-strategy game. *Journal of Neuroscience*, *27*, 8366–8377.
- Shahnazian, D., & Holroyd, C. B. (2018). Distributed representations of action sequences in anterior cingulate cortex: A recurrent neural network approach. *Psychonomic Bulletin & Review*, *25*, 302–321.
- Shenhav, A., Botvinick, M. M., & Cohen, J. D. (2013). The expected value of control: An integrative theory of anterior cingulate cortex function. *Neuron*, *79*, 217–240.
- Smith, E. H., Banks, G. P., Mikell, C. B., Cash, S. S., Patel, S. R., Eskandar, E. N., et al. (2015). Frequency-dependent representation of reinforcement-related information in the human medial and lateral prefrontal cortex. *Journal of Neuroscience*, *35*, 15827–15836.
- Tervo, D. G., Proskurin, M., Manakov, M., Kabra, M., Vollmer, A., Branson, K., et al. (2014). Behavioral variability through stochastic choice and its gating by anterior cingulate cortex. *Cell*, *159*, 21–32.
- Ullsperger, M., Danielmeier, C., & Jocham, G. (2014). Neurophysiology of performance monitoring and adaptive behavior. *Physiological Reviews*, *94*, 35–79.
- Wang, X.-J. (2008). Decision making in recurrent neuronal circuits. *Neuron*, *60*, 215–234.
- Widrow, B., & Hoff, M. E. (1960). *Adaptive switching circuits*. Technical report No. TR-1553-1. Palo Alto, CA: Stanford Electronics Labs.
- Wilcoxon, F. (1945). Individual comparisons by ranking methods. *Biometrics Bulletin*, *1*, 80–83.
- Wittmann, M. K., Kolling, N., Akaishi, R., Chau, B. K., Brown, J. W., Nelissen, N., et al. (2016). Predictive decision making driven by multiple time-linked reward representations in the anterior cingulate cortex. *Nature Communications*, *7*, 12327.
- Yu, A. J., & Dayan, P. (2005). Uncertainty, neuromodulation, and attention. *Neuron*, *46*, 681–692.