

# ONLINE APPENDIX FOR SELECTIVE ATTENTION AND LEARNING

---

**Joshua Schwartzstein**  
Dartmouth College

---

## Appendix A: Further Definitions and Results

This appendix presents formal results that are referenced in the main text, as well as definitions that are used in the appendices. Proofs of all results are presented in Appendix C.

### A.1. Supplementary Material to Section 2: Model

*A.1.1. Prior.* Here I give an alternative description of the agent's prior, which will be useful in presenting the proofs. The prior can compactly be expressed as

$$\mu(\theta) = \sum_{i \in \{X, -X\}} \sum_{j \in \{Z, -Z\}} \pi_{i,j} \mu^{i,j}(\theta).$$

Fix a model  $M \in \mathcal{M}$  and define  $c^M(x, \hat{z})$  as the set of covariates  $(x', \hat{z}') \in X \times \hat{Z}$  such that, under that model, any  $y_t$  given covariates  $(x_t, \hat{z}_t) = (x, \hat{z})$  is exchangeable with any  $y_{t'}$  given covariates  $(x_{t'}, \hat{z}_{t'}) = (x', \hat{z}')$ ; i.e., under model  $M$ ,  $\theta(x', \hat{z}') = \theta(x, \hat{z})$  with probability one if and only if  $(x', \hat{z}') \in c^M(x, \hat{z})$ . For example, under  $M_{X, -Z}$  (where only  $x$  is important to predicting  $y$ ),  $c^{X, -Z}(x, \hat{z}) = \{(x', \hat{z}') \in X \times \hat{Z} : x' = x\}$  equals the set of covariates that agree on  $x$ . With a slight abuse of notation, label the common success probability across members of  $c^M$  under model  $M$  by  $\theta(c^M)$ . Intuitively,  $c^M(x, \hat{z})$  equals the set of covariates that, under model  $M$ , can be lumped together with  $(x, \hat{z})$  without affecting the accuracy of the agent's predictions.

Let  $C^M$  denote the collection of  $c^M$ , so  $C^M$  is a partition of  $X \times \hat{Z}$ , and define  $\Theta(M) = [0, 1]^{\#C^M}$  as the effective parameter space under model  $M$  with generic element  $\theta(M)$ .  $\mu^M$  is defined by the joint distribution it assigns to the  $\#C^M$  parameters  $\theta(c^M)$ . These parameters are taken as independent with respect to  $\mu^M$  and distributed according to density,  $\psi(\cdot)$ . To take an example, if  $\psi(\theta) = \mathbf{1}_{\theta \in [0, 1]}$ , then  $\theta(c^M) \sim U[0, 1]$  for each  $M \in \mathcal{M}$ ,  $c^M \in C^M$ .

---

*The editor in charge of this paper was George-Marios Angeletos.*

E-mail: josh.schwartzstein@dartmouth.edu

*A.1.2. Forecasts.* This subsection describes the forecasts of an individual with selective attention in some detail (standard Bayesian forecasts are a special case) and presents some definitions which will be useful later.

Given the individual's prior, his period- $t$  forecast given recalled history  $\hat{h}^t$  is given by

$$\hat{E}[y|x, z, \hat{h}^t] = \sum_{M^t \in \mathcal{M}} \hat{\pi}_{M^t}^t E_{\xi}[\theta(c^{M^t}(x, \hat{z})) | \hat{h}^t, M^t], \quad (\text{A.1})$$

where

$$E_{\xi}[\theta(c^M) | \hat{h}^t, M] = \int \tilde{\theta} \psi(\tilde{\theta} | \hat{h}^t, c^M) d\tilde{\theta}$$

$$\psi(\tilde{\theta} | \hat{h}^t, c^M) = \frac{\tilde{\theta}^{\kappa(c^M | \hat{h}^t)} (1 - \tilde{\theta})^{N(c^M | \hat{h}^t) - \kappa(c^M | \hat{h}^t)} \psi(\tilde{\theta})}{\int \tau^{\kappa(c^M | \hat{h}^t)} (1 - \tau)^{N(c^M | \hat{h}^t) - \kappa(c^M | \hat{h}^t)} \psi(\tau) d\tau}.$$

$N(c^M | \hat{h}^t)$  denotes the number of times the covariates have taken on some value  $(x, \hat{z}) \in c^M$  along history  $\hat{h}^t$  and  $\kappa(c^M | \hat{h}^t)$  denotes the number of times that both the covariates have taken on such a value and  $y = 1$ . I will sometimes abuse notation and write  $N(x, \hat{z} | \hat{h}^t)$  and  $\kappa(x, \hat{z} | \hat{h}^t)$  instead of  $N(\{(x, \hat{z})\} | \hat{h}^t)$  and  $\kappa(\{(x, \hat{z})\} | \hat{h}^t)$ , respectively. Likewise, when convenient I will write  $N(x | \hat{h}^t)$  instead of  $N(\{(x', \hat{z}') : x' = x\} | \hat{h}^t)$ , etc.

To illustrate, (A.1) takes a particularly simple form when  $\psi(\theta) \sim U[0, 1]$ :  $\hat{E}[y|x, z, \hat{h}^t]$  equals

$$\hat{\pi}_{X,Z}^t \frac{\kappa(x, \hat{z} | \hat{h}^t) + 1}{N(x, \hat{z} | \hat{h}^t) + 2} + \hat{\pi}_{X,-Z}^t \frac{\kappa(x | \hat{h}^t) + 1}{N(x | \hat{h}^t) + 2} + \hat{\pi}_{-X,Z}^t \frac{\kappa(\hat{z} | \hat{h}^t) + 1}{N(\hat{z} | \hat{h}^t) + 2} + \hat{\pi}_{-X,-Z}^t \frac{\bar{\kappa}(\hat{h}^t) + 1}{t + 1}, \quad (\text{A.2})$$

where  $\bar{\kappa}(\hat{h}^t) = \sum_{x', \hat{z}'} \kappa(x', \hat{z}' | \hat{h}^t)$ .

For future reference,

$$\begin{aligned} \hat{\pi}_{i,j}^t &= \Pr_{\xi}(M_{i,j} | \hat{h}^t) \\ &= \frac{\Pr_{\xi}(\hat{h}^t | M_{i,j}) \pi_{i,j}}{\sum_{i',j'} \Pr_{\xi}(\hat{h}^t | M_{i',j'}) \pi_{i',j'}} \\ &= \frac{\alpha_{i,j} \mathcal{B}_{i,j}^t}{\sum_{i',j'} \alpha_{i',j'} \mathcal{B}_{i',j'}^t}, \end{aligned}$$

where

$$\begin{aligned} \mathcal{B}_{i,j}^t &= \frac{\Pr_{\xi}(\hat{h}^t | M_{i,j})}{\Pr_{\xi}(\hat{h}^t | M_{X,Z})} \\ &= \frac{\int \Pr_{\xi}(\hat{h}^t | \theta) \mu^{i,j}(d\theta)}{\int \Pr_{\xi}(\hat{h}^t | \theta) \mu^{X,Z}(d\theta)} \end{aligned}$$

is the *Bayes factor* comparing model  $M_{i,j}$  to model  $M_{X,Z}$  (Kass and Raftery 1995 provide a review of Bayes factors) and

$$\alpha_{i,j} = \frac{\pi_{i,j}}{\pi_{X,Z}} \quad (\text{A.3})$$

is the prior odds for  $M_{i,j}$  against  $M_{X,Z}$ .

## A.2. *Supplementary Material to Subsection 3.3: Systematically Biased Stereotypes and Beliefs*

This section fleshes out the results from Subsection 3.3. First, selective attention may lead people to form persistent and systematically incorrect beliefs about what causes variation in the data. Second, such biased beliefs are quite robust in the sense that, even if somebody points out the possibility that some other variable is what is truly important, it will take the agent some time to evaluate this claim because he did not previously attend to the relationship between that variable and the outcome.

*A.2.1. The Formation of Systematically Biased Beliefs.* To formalize, it is useful to specialize to the case where  $x$  is a binary random variable and  $X = \{0, 1\}$ . Define

$$\begin{aligned}\mathcal{R}_x(z') &= E_{\theta_0}[y|x = 1, z'] - E_{\theta_0}[y|x = 0, z'] \\ \mathcal{R}_x &= E_z[\mathcal{R}_x(z)|x = 1] \\ \varphi &= \text{Cov}_z(E_{\theta_0}[y|x = 0, z], g(x = 1|z)),\end{aligned}$$

where

- $\mathcal{R}_x(z')$  is the *standard Bayesian's limiting reaction to  $x$  conditional on  $z = z'$* : It equals the gap between the true conditional expectation of  $y$  given  $(x, z) = (1, z')$  and that given  $(x, z) = (0, z')$ .
- $\mathcal{R}_x$  is the *standard Bayesian's average limiting reaction to  $x$* : It equals the expected gap between the true conditional expectation of  $y$  given  $(x, z) = (1, z')$  and that given  $(x, z) = (0, z')$ , where the expectation is taken over  $z'$  conditional on  $x = 1$ .<sup>1</sup>
- $\varphi$  is the covariance between the likelihood that  $y = 1$  given  $(x, z) = (0, z')$  and the likelihood that  $x = 1$  given  $z'$ .  $\varphi > 0$  means that  $z$  which are associated with  $x = 1$  are also associated with  $y = 1$ ;  $\varphi < 0$  means that  $z$  which are associated with  $x = 1$  are also associated with  $y = 0$ . The magnitude  $|\varphi|$  measures the degree to which variation in  $z$  induces a relationship between the expected value of  $y$  and the likelihood that  $x = 1$ .

Additionally, let  $\hat{E}[y|x, z] \equiv \lim_{t \rightarrow \infty} \hat{E}[y|x, z, \hat{h}^t]$  denote the selectively attentive agent's limiting forecast given  $(x, z)$ , which almost surely exists by Propositions 1 and 3.

---

1.  $\mathcal{R}_x$  is formally equivalent to what is referred to as the population average treatment effect for the treated in the statistical literature on treatment effects, where  $x = 1$  corresponds to a treatment and  $x = 0$  to a control.

PROPOSITION A1. *Suppose  $b_k \equiv b$ , the agent settles on not encoding  $z$ , and  $X = \{0, 1\}$ . Then*

$$\hat{\mathcal{R}}_x(z') \equiv \hat{E}[y|x = 1, z'] - \hat{E}[y|x = 0, z'] = \mathcal{R}_x + \frac{\varphi}{\text{Var}(x)} \quad (\text{A.4})$$

*almost surely for all  $z'$ .*

Proposition A1 says that when the agent settles on not encoding  $z$ , his limiting reaction to  $x$  conditional on  $z = z'$ ,  $\hat{\mathcal{R}}_x(z')$ , differs from the standard Bayesian's,  $\mathcal{R}_x(z')$ , in two key ways corresponding to the two terms on the right hand side of (A.4). When  $\varphi = 0$ , the agent's limiting reaction reduces to the first term,  $\mathcal{R}_x$ : By persistently failing to encode  $z$ , the agent's limiting conditional reaction equals the standard Bayesian's limiting *average* reaction. Thinking of  $z$  as a situation, this is one of the distortions exploited in Mullainathan, Schwartzstein, and Shleifer (2008): By grouping distinct situations together in forming beliefs, an agent transfers the informational content of data across situations. For example, the agent may react to a piece of information which is uninformative in a given situation,  $z$ , because it is informative in another situation,  $z'$ .

When  $\varphi \neq 0$ , the agent's limiting conditional reaction differs from the standard Bayesian's limiting average reaction in an amount and direction determined by  $\varphi$ , which can be thought of as the magnitude and direction of omitted variable bias. A non-zero  $\varphi$  creates the possibility that, by settling on not encoding  $z$ , an agent will conclude a relationship between  $y$  and  $x$  that (weakly) reverses the true relationship conditional on *any*  $z'$  (e.g., that group  $A$  members are always more likely to be friendly than group  $B$  members when, in reality, they are equally likely conditional on the situation).

It is easy to see from Proposition A1 that a selectively attentive agent who fails to learn to pay attention to  $z$  can either over- or underreact to  $x$  at  $z'$ , depending on features of the joint distribution over  $(y, x, z)$ . The simplest case to consider is the one where  $x$  is unimportant to predicting  $y$  and there can only be overreaction, the topic of the next corollary.

COROLLARY A1. *Suppose the conditions of Proposition A1 hold and, additionally,  $x$  is unimportant to predicting  $y$ . Then, so long as  $\varphi \neq 0$ ,*

1.  $|\hat{\mathcal{R}}_x(z')| = \frac{|\varphi|}{\text{Var}(x)} \neq 0$  almost surely for all  $z'$ : The agent overreacts to  $x$  and the extent of overreaction is increasing in  $\frac{|\varphi|}{\text{Var}(x)}$ .
2.  $\hat{\pi}_x \xrightarrow{a.s.} 1$ : The agent becomes certain that  $x$  is important to predicting  $y$  even though it is not.

Corollary A1 considers the situation where  $x$  is completely unimportant to prediction and the selectively attentive agent settles on not encoding  $z$ . The first part says that, as a result of the possibility that the selectively attentive agent will

settle on not encoding  $z$ , he may come to overreact to  $x$ .<sup>2</sup> The degree to which the agent overreacts depends on the extent to which there is a tendency for  $z$ 's that are associated with  $x = 1$  to have relatively high (or low) corresponding success probabilities. Weakening this tendency will mitigate overreaction.

The second part of Corollary A1 says that, as a result of the possibility that the selectively attentive agent will settle on not encoding  $z$ , he may eventually become overconfident in the importance of  $x$ . This is true whenever  $z$  is associated with both  $x$  and  $y$  and the agent effectively suffers from omitted variable bias. Again, in this case, the agent mistakenly comes to view  $x$  as more than a proxy for selectively unattended to predictors.

These results relate to experimental findings that individuals attribute more of a causal role to information that is the focus of attention and to salient information more generally (Fiske and Taylor 2008, Chapter 3; also see Nisbett and Ross 1980, Chapter 6). To take an example, Taylor and Fiske (1975, Experiment 2) had participants watch a videotape of two people interacting in conversation. In the most relevant experimental condition, a third of the participants were instructed to pay particular attention to one of the conversationalists, a third were instructed to pay particular attention to the other, and the final third were told only to observe the conversation (i.e., they were not instructed to attend to anything in particular). Later, participants rated the extent to which each conversationalist determined the kind of information exchanged, set the tone of the conversation, and caused the partner to behave as he did. An aggregate score served as the dependent measure. The interaction between instructions and conversationalist was highly significant: Participants were more likely to see the conversationalist they attended to as causal in the interaction.<sup>3</sup>

*A.2.2. The Robustness of Systematically Biased Beliefs.* We have seen that selective attention allows for systematically biased beliefs to emerge and persist in the absence of communication across agents. What if an agent can credibly communicate the importance of a variable that the other has selectively failed to notice? A key feature of selective attention is that, following such “debiasing”, it will still take the agent time to learn to incorporate information about that variable in making predictions, since he did not keep track of it before, and to mitigate his misreaction to associated variables.

Suppose  $b_k \equiv b$ , the agent starts off not encoding  $z$  ( $\pi_z < b$ ), and  $X = \{0, 1\}$ . At some large  $t$ , the agent begins attending to  $z$  because there is an unmodeled shock to his belief that  $z$  is important to predicting  $y$  ( $\pi_z$  shifts to  $\pi'_z > \pi_z$ ) or to the degree to which he is cognitively busy ( $b$  shifts to  $b' < b$ ).<sup>4</sup> The main thing to note is that, even if this shock leads the agent to settle on encoding  $z$  and to make unbiased forecasts

2. Whenever  $X = \{0, 1\}$  and  $x$  is unimportant to predicting  $y$ , Proposition A1 establishes that Assumption 2 holds if and only if  $\varphi \neq 0$ , so it is technically redundant to include this condition in the statement of Proposition A1; it is included for clarity.

3. Participants also retained more information about the conversationalist they attended to.

4. We can think of shocks to  $\pi_z$  as resulting from a media report or something learned in a class and shocks to  $b$  as resulting from some (not modeled) reason why the agent would begin caring more about

in the limit, he continues to make systematically biased forecasts for a long time. The reason is that it takes some time for the agent to learn how to use both  $x$  and  $z$  to make predictions since he has not attended to  $z$  in the past (there is “learning by encoding”). To see this, consider how the agent reacts to  $x$  given  $z$  at the time of the shock  $t$ . Since  $t$  is assumed to be large,  $E_\xi[\theta(x, z)|M_{x, -z}, \hat{h}^t] \approx E_{\theta_0}[y|x]$  and  $\hat{\pi}_x^t \approx 1$  by the results of Section 3, so the agent’s reaction to  $x$  given  $z$  in that period equals

$$\begin{aligned} E_\xi[\theta(1, z)|\hat{h}^t] - E_\xi[\theta(0, z)|\hat{h}^t] &\approx \pi_Z'[\tau - \tau] + (1 - \pi_Z')(E_{\theta_0}[y|x = 1] - E_{\theta_0}[y|x = 0]) \\ &= (1 - \pi_Z')(E_{\theta_0}[y|x = 1] - E_{\theta_0}[y|x = 0]), \end{aligned} \tag{A.5}$$

where  $\tau = E_\psi[\theta]$  equals the prior success probability under density  $\psi$ . From (A.5), the agent’s reaction to  $x$  in period  $t$  is approximately proportional to his reaction the period before when he did not attend to  $z$ . This is intuitive: By not having attended to  $z$  in the past he has not learned that  $z$  is important to predicting  $y$  or how  $z$  is important to predicting  $y$ . As a result, even when he attends to  $z$ , his forecast places substantial weight on the empirical frequency of  $y = 1$  given only  $(x)$ .

### A.3. Supplementary Material to Section 4: Continuous Attention

This section formalizes two observations from Section 4. First, when the continuous attention assumptions hold, the agent eventually learns to attend to  $z$  and to make accurate forecasts since he always attends to  $z$  with positive probability and almost surely encodes  $z$  an infinite number of times.

PROPOSITION A2. *Suppose the continuous attention assumptions hold. Then*

1.  $\eta(\hat{\pi}_Z^t) \rightarrow 1$  almost surely.
2. For each  $x, z$ ,  $\hat{E}[y|x, z, \hat{h}^t]$  converges to  $E_{\theta_0}[y|x, z]$  in probability.

Second, even though the agent eventually learns to attend to  $z$  and to make accurate forecasts with arbitrarily large probability in the limit, he may continue not to attend to  $z$  and to make biased forecasts for a long time. In particular, a partial analog to Proposition 2 holds in this setting: the agent’s ability to recognize empirical relationships within a reasonable time horizon depends on his prior.

PROPOSITION A3. *Suppose the continuous attention assumptions hold. Then, for all  $t \geq 2$ , the probability that the agent never encodes  $z$  before period  $t$  tends towards 1 as  $\pi_Z \rightarrow 0$ .*

---

predicting  $y$  (e.g., he begins caring more about what leads to weight gain if he recently had a heart attack) or it becomes easier for the agent to attend to  $z$  (perhaps an attribute of a product is suddenly unshrouded).

## Appendix B: Sophistication

This appendix describes how the predictions of the discrete attention model would change if the agent was assumed to be sophisticated rather than naive. Again, proofs of all formal results are presented in Appendix C.

There are two main differences if the agent is sophisticated, as will be more formally discussed below. First, a sophisticate may update his beliefs about the importance of variables he never attends to. For such an agent, a belief that  $z$  is unimportant to prediction is not necessarily self-confirming in the strong sense identified in Proposition 2. (Nevertheless, Proposition B2 below shows that a sophisticate with a “weak enough” prior belief in the importance of  $z$  will never encode  $z$ .) Second, the sophisticate entertains the possibility that unattended-to variables covary with those he attends to in a way that is responsible for the relationships he identifies. Such an agent cannot become wrongly certain of the importance of  $x$ : sophistication may alter long-run beliefs.<sup>5</sup>

### *B.1. When is a Belief that $z$ is Unimportant to Prediction Self-Confirming for a Sophisticate?*

DEFINITION B1. A belief that  $z$  is unimportant to prediction is *strongly self-confirming* if and only if the agent does not update his beliefs positively about the importance of  $z$  when he never encodes  $z$ ; that is, if and only if  $\hat{\pi}_z^t \leq \pi_z$  given every represented history of the form  $(y_{t-1}, x_{t-1}, \emptyset, \dots, y_1, x_1, \emptyset)$ .

Lemma C3 shows that, when the agent is naive, a belief that  $z$  is unimportant to prediction is strongly self-confirming. This is not true when the agent is sophisticated.

PROPOSITION B1. *When the agent is sophisticated, there are primitives (i.e., values of  $\psi(\cdot)$ ,  $\pi_x$ ,  $\pi_z$ , etc.) under which a belief that  $z$  is unimportant to prediction is not strongly self-confirming.*

The basic idea behind why the sophisticate may update his beliefs about the importance of never attended-to variables is that the variability of the outcome conditional on variables he attends to may provide a subjective signal indicating that unattended-to variables influence the success probability. Indeed, the proof of Proposition B1 provides an example. Nevertheless, even when the agent is sophisticated, a belief that  $z$  is unimportant to prediction is self-confirming in the weaker sense that the degree to which the agent can positively update his belief about the importance of  $z$  is bounded when he never encodes  $z$ . As a result, the sophisticate will never encode  $z$  when he starts off with a “weak enough” prior belief about its importance.

---

5. However, sophistication does not alter long-run forecasts in the sense that Proposition 3 continues to hold if the agent is sophisticated.

To formalize, we have the identity

$$\hat{\pi}_Z^t = \frac{\Lambda(\hat{h}^t)\pi_Z}{(\Lambda(\hat{h}^t) - 1)\pi_Z + 1},$$

where

$$\Lambda(\hat{h}^t) = \frac{1 + \frac{1-\pi_X}{\pi_X} \mathcal{B}_{-X,Z}(\hat{h}^t)}{\mathcal{B}_{X,-Z}(\hat{h}^t) + \frac{1-\pi_X}{\pi_X} \mathcal{B}_{-X,-Z}(\hat{h}^t)} \quad (\text{B.1})$$

can be thought of as a “likelihood ratio” comparing the likelihood of the mentally represented history given models where  $z$  is important to prediction as compared to the likelihood of that history given models where  $z$  is unimportant to prediction. (Unless otherwise noted, in this appendix all likelihoods and Bayes’ factors are taken with respect to the sophisticated, rather than naive, agent’s beliefs).

Given certain technical conditions, I will first establish an upper bound on  $\Lambda$  that holds uniformly across  $t$  and mentally represented histories where the agent has not encoded  $z$  through  $t$ . To elaborate on these technical conditions, define  $\theta^X(x') = \sum_{z'} \theta(x', z')g(z'|x')$  as the success probability conditional on  $x'$  given  $\theta$  and write  $\theta^X = (\theta^X(x'))_{x' \in \mathcal{X}}$ . Denote the c.d.f. over  $\theta^X$  given model  $M_{i,j}$  by  $\Psi_{i,j}^X$  and note that both  $\Psi_{X,Z}^X$  and  $\Psi_{X,-Z}^X$  have associated densities  $\psi_{X,Z}^X$  and  $\psi_{X,-Z}^X$  which are bounded above. In the remainder of this appendix, make the further technical assumptions that  $\Psi_{-X,Z}^X$  also has an associated density  $\psi_{-X,Z}^X$  which is bounded above, positive at all values in the support of  $\mu_{-X,Z}$ , and continuous in a neighborhood of  $\theta_0^X$  (the last two assumptions are not needed to prove Lemma B1, but will be convenient later).<sup>6</sup>

LEMMA B1. *There exists a constant  $K < \infty$  such that, for the sophisticated agent,*

$$\Lambda(\hat{h}^t) \leq K$$

for all  $t$  and mentally represented histories of the form  $\hat{h}^t = (y_{t-1}, x_{t-1}, \emptyset, \dots, y_1, x_1, \emptyset)$ .

It is almost immediate from Lemma B1 that a sophisticate with a weak enough prior belief about the importance of  $z$  will never encode  $z$ .

PROPOSITION B2. *Fix  $b \in (0, 1)$ . There exists some  $\bar{\pi}_Z \in (0, b)$  such that the sophisticated agent never encodes  $z$  if  $\pi_Z < \bar{\pi}_Z$ .*

6. When I say that  $\theta^X$  is in the support of  $\mu_{i,j}$  I mean that there exists a  $\theta'$  in the support of  $\mu_{i,j}$  such that  $\theta'^X = \theta^X$ . When there are only two values of  $x$ , two values of  $z$ , and  $g(z'|x') \neq g(z'|x'')$  (as in the stereotyping example), we can use the change of variables formula to analytically derive that  $\Psi_{-X,Z}^X$  has an associated density  $\psi_{-X,Z}^X$  which is bounded above and positive at all values in the support of  $\mu_{-X,Z}$ .

***B.2. What Does the Sophisticated Agent Come to Believe About the Importance of  $x$ ?***

Proposition 4 established that a naive agent who settles on not encoding  $z$  will almost surely become certain that  $x$  is important to predicting  $y$  so long as  $x$  is unconditionally (of  $z$ ) predictive, where we mean the following when we say the agent becomes certain:

**DEFINITION B2.** The agent *becomes certain that  $x$  is important to predicting  $y$*  if and only if  $\hat{\pi}_x^t \rightarrow 1$  as  $t \rightarrow \infty$ .

The conclusion of Proposition 4 does not hold when the agent is sophisticated: He will not erroneously become convinced in the importance of  $x$  if he can explain what he observes in a model where only  $z$  is important to prediction; that is, when omitted variable bias could fully explain the estimated relationship between  $y$  and  $x$ .

**PROPOSITION B3.** *Suppose  $x$  is unimportant to predicting  $y$  (but is predictive unconditional of  $z$ ) and the agent never encodes  $z$  because  $\pi_Z < \bar{\pi}_Z$ , where  $\bar{\pi}_Z$  is defined as in Proposition B2. Then, while the naive agent would almost surely become certain that  $x$  is important to predicting  $y$ , the sophisticated agent almost surely would not.*

## Appendix C: Proofs

See Appendix D for the development of results concerning the evolution of Bayes' factors in my model. I will frequently refer to results from that appendix in the proofs below.

### C.1. Proofs of Results from Section 2: Model

*Proof of Observation 1.* To establish the first part of the observation, recall that the standard Bayesian's period  $t$  forecast satisfies

$$E[y|x, z, h^t] = \sum_{M' \in \mathcal{M}} \pi_{M'}^t E[\theta(c^{M'}(x, z)) | h^t, M']. \quad (\text{C.1})$$

Fix an  $M \in \mathcal{M}$ . Since the marginal prior density over  $\theta(c^M(x, z))$  is non-doctrinaire under  $M$ ,  $E[\theta(c^M(x, z)) | h^t, M] \xrightarrow{a.s.} E_{\theta_0}[y | c^M(x, z)]$  by Freedman (1963).

In addition,  $E_{\theta_0}[y | c^M(x, z)] = E_{\theta_0}[y | x, z]$  for  $M = M_{X,Z}$  as well as for  $M_{-X,Z}$  when  $M_{-X,Z}$  is the true model. Consequently, it is left to show that  $\pi_{X,-Z}^t$  and  $\pi_{-X,-Z}^t$  converge almost surely to zero and that  $\pi_{-X,Z}^t$  converges almost surely to zero whenever  $M_{-X,Z}$  is *not* the true model. But these statements follow immediately from Lemma D5.

The second part of the observation follows directly from Lemma D5.  $\square$

### C.2. Proofs of Results from Section 3: Discrete Attention

*Proof of Proposition 1.* Suppose that, with positive probability under  $P_{\theta_0, \xi}(\cdot)$ , the agent does not settle on encoding or not encoding  $z$  ( $b$  must satisfy  $0 < b < 1$ ). Label this event  $NS$  and condition on  $\hat{h}^\infty \in NS$ . Since the agent encodes  $z$  infinitely often conditional on  $NS$ , by Lemma D7 we must have  $\hat{\pi}_Z^t \rightarrow 1$  with probability 1. As a result, with probability 1 there exists a  $\tilde{t}$  such that  $\hat{\pi}_Z^t > b$  for all  $t \geq \tilde{t}$  so  $e_t = 1$  for all  $t \geq \tilde{t}$ , a contradiction.  $\square$

A few Lemmas will be useful in establishing Proposition 2. First, define the likelihood ratio (or Bayes' factor) comparing the likelihood of a history under models where  $z$  is important to predicting  $y$  to the likelihood of that history under models where  $z$  is unimportant to predicting  $y$ ,  $\Lambda(h^t)$ , as in Equation (B.1).

LEMMA C1.  $\pi_Z(h^t) > b$  if and only if

$$\Lambda(h^t) > \frac{1 - \pi_Z}{\pi_Z} \frac{b}{1 - b}.$$

*Proof.* Simple algebra.  $\square$

LEMMA C2. For all  $\varepsilon > 0$  there exists  $\lambda > 0$  such that

$$P_{\theta_0} \left( \min_{t' \geq 1} \Lambda(h^{t'}) > \lambda \right) \geq 1 - \varepsilon \quad (\text{C.2})$$

*Proof.* Fix  $\varepsilon > 0$ . From Lemma D5, we know that  $\mathcal{B}_{X,-Z}(h^t) \xrightarrow{a.s.} 0$ ,  $\mathcal{B}_{-X,-Z}(h^t) \xrightarrow{a.s.} 0$ . As a result,

$$\Lambda(h^t) = \frac{\left(1 + \frac{1+\pi_X}{\pi_X} \mathcal{B}_{-X,Z}(h^t)\right)}{\left(\mathcal{B}_{X,-Z}(h^t) + \frac{1-\pi_X}{\pi_X} \mathcal{B}_{-X,-Z}(h^t)\right)} \xrightarrow{a.s.} \infty.$$

Consequently, there exists a value  $T \geq 1$  such that  $P_{\theta_0}(\min_{t' \geq T} \Lambda(h^{t'}) \geq 1) > 1 - \varepsilon$  (see, for example, Lemma 7.2.10 in Grimmett and Stirzaker 2001).

Since, in addition, there exists  $\lambda$  ( $0 < \lambda < 1$ ) such that

$$\min_h \min_{1 \leq k \leq T} \Lambda(h^k) > \lambda,$$

we have

$$P_{\theta_0}(\min_{t' \geq 1} \Lambda(h^{t'}) > \lambda) \geq 1 - \varepsilon.$$

□

The next Lemma establishes some finite sample properties of Bayes' factors and beliefs when the agent never encodes  $z$ .

LEMMA C3. *If the (naive) agent has not encoded  $z$  up to period  $t$ , then*

$$\begin{aligned} \mathcal{B}_{-X,Z}(\hat{h}^t) &= \mathcal{B}_{-X,-Z}(\hat{h}^t) \\ \mathcal{B}_{X,-Z}(\hat{h}^t) &= 1 \\ \pi_Z(\hat{h}^t) &= \pi_Z. \end{aligned}$$

*Proof.* If the agent has not encoded  $z$  up to period  $t$ , then

$$\begin{aligned} \mathcal{B}_{-X,Z}(\hat{h}^t) &= \frac{\int_0^1 \theta^{\kappa(\varnothing|\hat{h}^t)} (1-\theta)^{N(\varnothing|\hat{h}^t) - \kappa(\varnothing|\hat{h}^t)} \psi(\theta) d\theta}{\prod_{x'} \int_0^1 \theta^{\kappa(x'|\hat{h}^t)} (1-\theta)^{N(x'|\hat{h}^t) - \kappa(x'|\hat{h}^t)} \psi(\theta) d\theta} \\ &= \frac{\int_0^1 \theta^{\bar{\kappa}(\hat{h}^t)} (1-\theta)^{t-1-\bar{\kappa}(\hat{h}^t)} \psi(\theta) d\theta}{\prod_{x'} \int_0^1 \theta^{\kappa(x'|\hat{h}^t)} (1-\theta)^{N(x'|\hat{h}^t) - \kappa(x'|\hat{h}^t)} \psi(\theta) d\theta} \\ \mathcal{B}_{-X,-Z}(\hat{h}^t) &= \frac{\int_0^1 \theta^{\bar{\kappa}(\hat{h}^t)} (1-\theta)^{t-1-\bar{\kappa}(\hat{h}^t)} \psi(\theta) d\theta}{\prod_{x'} \int_0^1 \theta^{\kappa(x'|\hat{h}^t)} (1-\theta)^{N(x'|\hat{h}^t) - \kappa(x'|\hat{h}^t)} \psi(\theta) d\theta} \\ &= \mathcal{B}_{-X,Z}(\hat{h}^t) \\ \mathcal{B}_{X,-Z}(\hat{h}^t) &= \frac{\prod_{x'} \int_0^1 \theta^{\kappa(x'|\hat{h}^t)} (1-\theta)^{N(x'|\hat{h}^t) - \kappa(x'|\hat{h}^t)} \psi(\theta) d\theta}{\prod_{x'} \int_0^1 \theta^{\kappa(x'|\hat{h}^t)} (1-\theta)^{N(x'|\hat{h}^t) - \kappa(x'|\hat{h}^t)} \psi(\theta) d\theta} \\ &= 1. \end{aligned}$$

Plugging these expressions into the definition of  $\pi_Z(\hat{h}^t)$  yields

$$\begin{aligned}\pi_Z(\hat{h}^t) &= \frac{\pi_Z \left[ 1 + \frac{1-\pi_X}{\pi_X} \mathcal{B}_{-X,Z}(\hat{h}^t) \right]}{1 + \frac{1-\pi_X}{\pi_X} \mathcal{B}_{-X,Z}(\hat{h}^t)} \\ &= \pi_Z.\end{aligned}$$

□

*Proof of Proposition 2.*

Part (1). First I show that, for all  $\varepsilon > 0$ , there exists  $\pi_1 \in (0, 1)$  (or  $b_1 \in (0, 1)$ ) such that the agent settles on encoding  $z$  with probability at least  $1 - \varepsilon$  for all  $\pi_Z \geq \pi_1$  ( $b \leq b_1$ ). Fix  $\varepsilon$ . Note that, whenever  $\hat{\pi}_Z^k > b$  for all  $k < t$ ,  $\hat{h}^t = h^t$ ,  $\hat{\pi}_Z^t = \pi_Z(h^t)$ , and  $P_{\theta_0, \xi}(\hat{h}^t) = P_{\theta_0}(h^t)$ . As a result, it is sufficient to show that there exists  $\pi_1 \in (0, 1)$  ( $b_1 \in (0, 1)$ ) such that

$$P_{\theta_0}(\min_{t' \geq 1} \pi_Z(h^{t'}) > b) \geq 1 - \varepsilon$$

whenever  $\pi_Z \geq \pi_1$  ( $b \leq b_1$ ).

By Lemma C1,

$$\begin{aligned}\pi_Z(h^t) > b &\iff \\ \Lambda(h^t) &> \frac{1 - \pi_Z}{\pi_Z} \frac{b}{1 - b}.\end{aligned}$$

Consequently,  $P_{\theta_0}(\min_{t' \geq 1} \pi_Z(h^{t'}) > b) \geq 1 - \varepsilon$  if and only if

$$P_{\theta_0} \left( \min_{t' \geq 1} \Lambda(h^{t'}) > \frac{1 - \pi_Z}{\pi_Z} \frac{b}{1 - b} \right) \geq 1 - \varepsilon.$$

From Lemma C2 we know that there exists  $\lambda(\varepsilon) > 0$  such that

$$P_{\theta_0} \left( \min_{t' \geq 1} \Lambda(h^{t'}) > \lambda(\varepsilon) \right) \geq 1 - \varepsilon,$$

so the result follows from setting  $\pi$  or  $b$  to satisfy

$$\begin{aligned}\frac{1 - \pi}{\pi} \frac{b}{1 - b} &= \lambda(\varepsilon) \Rightarrow \\ \pi_1 &= \frac{b}{b + \lambda(\varepsilon)(1 - b)} \\ b_1 &= \frac{\lambda(\varepsilon)\pi}{\pi(\lambda(\varepsilon) - 1) + 1}.\end{aligned}$$

Part (2). It is left to show that the agent settles on not encoding  $z$  with probability 1 when  $\pi_Z < b$ . It is sufficient to show that, when  $\pi_Z < b$ ,  $\pi_Z(\hat{h}^t) = \pi_Z$  for all  $t > 1$ . But this follows from Lemma C3. □

*Proof of Proposition 3.* Part (1). Analogous to the proof of Observation 1.1 and hence omitted.

Part (2). If the agent settles on not encoding  $z$  then, by definition, there exists  $n$  such that  $e_t = 0$  for all  $t \geq n$ . In any period  $t \geq n$ , the agent's expectation satisfies

$$\begin{aligned}\hat{E}[y|x, z, \hat{h}^t] &= E_\xi[\theta(x, \emptyset)|\hat{h}^t] \\ &= \sum_{M' \in \mathcal{M}} \hat{\pi}_{M'}^t E_\xi[\theta(c^{M'}(x, \emptyset))|\hat{h}^t, M'].\end{aligned}$$

Fix an  $M \in \mathcal{M}$ . Since the marginal prior density over  $\theta(c^M(x, \emptyset))$  is non-doctrinaire under  $M$ ,  $E_\xi[\theta(c^M(x, \emptyset))|\hat{h}^t, M] \xrightarrow{a.s.} E_{\theta_0}[y|c^M(x, \emptyset)]$  by Freedman (1963), where  $E_{\theta_0}[y|c^M(x, \emptyset)] = E_{\theta_0}[y|x]$  for  $M \in \{M_{X,Z}, M_{X,-Z}\}$  and  $E_{\theta_0}[y|c^M(x, \emptyset)] = E_{\theta_0}[y]$  for  $M \in \{M_{-X,Z}, M_{-X,-Z}\}$ .

If  $E_{\theta_0}[y|x] = E_{\theta_0}[y]$ , then we are done. Assume  $E_{\theta_0}[y|x] \neq E_{\theta_0}[y]$  for some  $x$ . It is left to show that both  $\hat{\pi}_{-X,Z}^t$  and  $\hat{\pi}_{-X,-Z}^t$  converge almost surely to zero. Equivalently, it is left to show that both  $\mathcal{B}_{-X,Z}(\hat{h}^t)$  and  $\mathcal{B}_{-X,-Z}(\hat{h}^t)$  converge almost surely to zero.

For  $t \geq n$  and  $j \in \{-Z, Z\}$ ,

$$\begin{aligned}\mathcal{B}_{-X,j}(\hat{h}^t) &= \frac{\Pr_\xi(\hat{h}^t|M_{-X,j})}{\Pr_\xi(\hat{h}^t|M_{X,Z})} \\ &= \frac{\Pr_\xi(\hat{h}_n^t|M_{-X,j}, \hat{h}^n) \Pr_\xi(\hat{h}^n|M_{-X,j})}{\Pr_\xi(\hat{h}_n^t|M_{X,Z}, \hat{h}^n) \Pr_\xi(\hat{h}^n|M_{X,Z})} \\ &= \left( \frac{\int \prod_{k=n}^{t-1} \theta(x_k, \emptyset)^{y_k} (1 - \theta(x_k, \emptyset))^{1-y_k} \mu^{-X,j}(d\theta|\hat{h}^n)}{\int \prod_{k=n}^{t-1} \theta(x_k, \emptyset)^{y_k} (1 - \theta(x_k, \emptyset))^{1-y_k} \mu^{X,Z}(d\theta|\hat{h}^n)} \right) \frac{\Pr_\xi(\hat{h}^n|M_{-X,j})}{\Pr_\xi(\hat{h}^n|M_{X,Z})},\end{aligned}$$

where  $\hat{h}_n^t = (y_{t-1}, x_{t-1}, \emptyset, \dots, y_n, x_n, \emptyset)$ . Since  $\Pr_\xi(\hat{h}^n|M_{-X,j})/\Pr_\xi(\hat{h}^n|M_{X,Z})$  is fixed for all  $t \geq n$ , it is necessary and sufficient to show that

$$\left( \frac{\int \prod_{k=n}^{t-1} \theta(x_k, \emptyset)^{y_k} (1 - \theta(x_k, \emptyset))^{1-y_k} \mu^{-X,j}(d\theta|\hat{h}^n)}{\int \prod_{k=n}^{t-1} \theta(x_k, \emptyset)^{y_k} (1 - \theta(x_k, \emptyset))^{1-y_k} \mu^{X,Z}(d\theta|\hat{h}^n)} \right) \quad (\text{C.3})$$

converges to zero almost surely to establish such convergence of the Bayes' factor. But, noting that (i) (C.3) equals  $\mathcal{B}_{-X,j}^t(p_0^0)$  for some uniformly non-doctrinaire  $\mu^{-X,j}$ ,  $\mu^{X,Z}$ , and (ii)

$$(y_{t-1}, x_{t-1}, \hat{z}_{t-1}, \dots, y_n, x_n, \hat{z}_n)$$

is a random sample from  $p_0^0$ , the result follows from Lemma D6.  $\square$

*Proof of Proposition 4.* Part (1). Analogous to the proof of Observation 1.2 and hence omitted.

Part (2). The fact that  $\hat{\pi}_X^t \xrightarrow{a.s.} 1$  when the agent settles on not encoding  $z$  follows immediately from Assumption 2 and Lemma D6. That  $\hat{\pi}_Z^t \leq b$  for large  $t$  follows from the definition of settling on not encoding  $z$  and the encoding rule.  $\square$

### C.3. Proofs of Results from Appendix A.2: Supplementary Results to Subsection 3.3

*Proof of Proposition A1.* This result follows from theorems in Samuels (1993), but, for completeness, I'll provide a proof.

Let the true distribution over  $(y, x, z)$  be denoted by  $p_0(\cdot)$  (the distribution generated by  $\theta_0$  and  $g(\cdot)$ ) and let  $\mathbb{E}$  denote the expectation operator under  $p_0(\cdot)$ . With this notation,

$$\begin{aligned}\mathcal{R}_x(z') &= \mathbb{E}[y|x = 1, z'] - \mathbb{E}[y|x = 0, z'] \\ \mathcal{R}_x &= \mathbb{E}[\mathcal{R}_x(z)|x = 1] \\ \varphi &= \text{Cov}(\mathbb{E}[y|x = 0, z], g(x = 1|z)).\end{aligned}$$

From Proposition 3,  $\hat{E}[y|x = 1, z] - \hat{E}[y|x = 0, z]$  almost surely equals

$$\begin{aligned}\mathbb{E}[y|x = 1] - \mathbb{E}[y|x = 0] &= \mathbb{E}[\mathbb{E}[y|x, z]|x = 1] - \mathbb{E}[\mathbb{E}[y|x, z]|x = 0] \\ &= \frac{\mathbb{E}[\mathbb{E}[y|x = 1, z]g(x = 1|z)](1 - g(x = 1)) - \mathbb{E}[\mathbb{E}[y|x = 0, z]g(x = 0|z)]g(x = 1)}{g(x = 1)(1 - g(x = 1))} \\ &= \frac{\mathbb{E}[\mathcal{R}_x(z)g(x = 1|z)](1 - g(x = 1))}{g(x = 1)(1 - g(x = 1))} + \frac{\mathbb{E}[\mathbb{E}[y|x = 0, z](g(x = 1|z) - g(x = 1))]}{g(x = 1)(1 - g(x = 1))} \\ &= \mathcal{R}_x + \frac{\varphi}{g(x = 1)(1 - g(x = 1))} \\ &= \mathcal{R}_x + \frac{\varphi}{\text{Var}(x)}.\end{aligned}$$

□

*Proof of Corollary A1.* By the assumption that  $x$  is unimportant to predicting  $y$ ,  $\mathcal{R}_x(z') = 0$  for all  $z'$  so  $\mathcal{R}_x = 0$ . Then, by Proposition A1,

$$\begin{aligned}\hat{\mathcal{R}}_x(z') &= E_{\theta_0}[y|x = 1] - E_{\theta_0}[y|x = 0] \\ &= \frac{\varphi}{\text{Var}(x)},\end{aligned}\tag{C.4}$$

which establishes the first part of the Corollary. Additionally,  $E_{\theta_0}[y|x = 1] - E_{\theta_0}[y|x = 0] \neq 0$  whenever  $\varphi \neq 0$  (by (C.4)) and the second part of the Corollary then follows from Proposition 4. □

### C.4. Proofs of Results from Sections 4 and A.3: Continuous Attention

LEMMA C4. *If the continuous attention assumptions hold then the agent almost surely encodes  $z$  an infinite number of times.*

*Proof.* Fix some  $t$ . I will show that the probability that the agent never encodes  $z$  after  $t$  is bounded above by 0. Independent of the history before  $t$ , the probability of not encoding  $z$  at  $t+k$  ( $k > 0$ ) given not having encoded  $z$  between  $t$  and  $t+k$  is strictly less than

$$1 - \eta \left( \frac{b\pi_Z}{a(1 - \pi_Z) + b\pi_Z} \right) < 1, \quad (\text{C.5})$$

where  $a$  and  $b$  are positive constants (do not depend on  $k$ ). As a result, the probability of never encoding  $z$  after  $t$  is less than the infinite product

$$\left( 1 - \eta \left( \frac{b\pi_Z}{a(1 - \pi_Z) + b\pi_Z} \right) \right)^\infty = 0$$

and the result follows.  $\square$

*Proof of Proposition A2.*

Part (1): From Lemma C4 we know that the agent almost surely encodes  $z$  an infinite number of times. Combining this result with Lemma D7, we have that  $\hat{\pi}_Z^t \rightarrow 1$  almost surely which implies that  $\eta(\hat{\pi}_Z^t) \rightarrow 1$  almost surely.

Part (2): Fix  $(x, z)$  and  $\varepsilon > 0$ . We want to show that

$$\lim_{t \rightarrow \infty} P_{\theta_0, \xi}(|\hat{E}[y|x, z, \hat{h}^t] - E_{\theta_0}[y|x, z]| > \varepsilon) = 0. \quad (\text{C.6})$$

Let us expand  $P_{\theta_0, \xi}(|\hat{E}[y|x, z, \hat{h}^t] - E_{\theta_0}[y|x, z]| > \varepsilon)$ , which equals

$$\begin{aligned} & P_{\theta_0, \xi}(|E_\xi[\theta(x, z)|\hat{h}^t] - E_{\theta_0}[y|x, z]| > \varepsilon)P_{\theta_0, \xi}(e_t = 1) \\ & + P_{\theta_0, \xi}(|E_\xi[\theta(x, \emptyset)|\hat{h}^t] - E_{\theta_0}[y|x, z]| > \varepsilon)(1 - P_{\theta_0, \xi}(e_t = 1)). \end{aligned}$$

The expansion indicates that, to establish (C.6), it is sufficient to show two things:

A..  $E_\xi[\theta(x, z)|\hat{h}^t] \xrightarrow{a.s.} E_{\theta_0}[y|x, z]$

B..  $P_{\theta_0, \xi}(e_t = 1) \rightarrow 1$

A. follows from now familiar arguments applying the non-doctrinaire assumption, the strong law of large numbers, the consistency properties of Bayes' factors, and the fact that the agent encodes  $z$  an infinite number of times (Lemma C4). B. follows from the fact that  $P_{\theta_0, \xi}(e_t = 1) = E_{\theta_0, \xi}[\eta(\hat{\pi}_Z^t)]$  and  $E_{\theta_0, \xi}[\eta(\hat{\pi}_Z^t)] \rightarrow 1$  because  $\eta(\hat{\pi}_Z^t)$  is bounded and tends almost surely towards 1 by Proposition A2.1.  $\square$

*Proof of Proposition A3.* By Lemma C3, the probability that the agent never encodes  $z$  before period  $t$  is given by

$$(1 - \eta(\pi_Z))^{t-1}, \quad (\text{C.7})$$

which tends towards 1 as  $\pi_Z \rightarrow 0$ .  $\square$

The next lemma demonstrates that the fraction of time that the agent spends encoding  $z$  tends towards 1 assuming continuous attention. Recall that  $\mathcal{E}(t) = \{\tau < t : \hat{z}_\tau \neq \emptyset\}$  denotes the number of times the agent encodes  $z$  prior to period  $t$ .

LEMMA C5. *If the continuous attention assumptions hold then*

$$\frac{\#\mathcal{E}(t)}{t-1} \xrightarrow{a.s.} 1.$$

*Proof.* Define

$$\gamma^t = \frac{\#\mathcal{E}(t)}{t-1}.$$

I will apply a result from the theory of stochastic approximation to show that  $\gamma^t \xrightarrow{a.s.} 1$  (Benaim 1999). We have

$$\begin{aligned} \gamma^t - \gamma^{t-1} &= \frac{e_t - \gamma^{t-1}}{t-1} \\ &= \frac{1}{t-1} (F(\gamma^{t-1}) + \varepsilon_t + u_t), \end{aligned}$$

where

$$\begin{aligned} F(\gamma^{t-1}) &= 1 - \gamma^{t-1} \\ \varepsilon_t &= e_t - \eta(\hat{\pi}_Z^t) \\ u_t &= \eta(\hat{\pi}_Z^t) - 1. \end{aligned}$$

Note that

1.  $F$  is Lipschitz continuous and is defined on a compact set  $[0, 1]$
2.  $E[\varepsilon_t | \hat{h}^t] = 0$  and  $E[|\varepsilon_t|^2] < \infty$
3.  $u_t \xrightarrow{a.s.} 0$

so Theorem A in Fudenberg and Takahashi (2008) tells us that, with probability 1, every  $\omega$ -limit of the process  $\{\gamma^t\}$  is connected and internally chain recurrent for  $\Phi$ , where  $\Phi$  is the continuous time semi-flow induced by

$$\dot{\gamma}(t) = F(\gamma(t)).$$

Since  $F'(\gamma) = -1 < 0$  and the unique steady state of the continuous time process is  $\gamma = 1$ , the only connected and internally chain recurrent set for  $\Phi$  is  $\{1\}$  by Liouville's Theorem.  $\square$

Note that the next few lemmas (Lemmas C6 - C8) rely heavily on results and notation from Appendix D.

LEMMA C6.  $d = \delta_{X, -Z}(p_0^1)$

*Proof.* Apply Lemma D1 to get  $\underline{\theta}(c^{X,-Z}(x,z)) = p_0^1(y|x) = p_{\theta_0}(y|x)$  for all  $x$ . The result then follows from the definition of  $\delta_{X,-Z}(p_0^1)$ .  $\square$

LEMMA C7. *If the continuous attention assumptions hold, then*

$$\frac{\mathcal{B}_{X,-Z}(\hat{h}^t)}{e^{-d(t-1)}} \xrightarrow{a.s.} K$$

for some  $K$  satisfying  $0 < K < \infty$ .

*Proof.* I will show that

$$\frac{1}{t-1} \log \mathcal{B}_{X,-Z}(\hat{h}^t) \rightarrow -d \quad (\text{C.8})$$

almost surely. We can write

$$\mathcal{B}_{X,-Z}(\hat{h}^t) = \frac{\Pr_{\xi}(\hat{h}_1^t | M_{X,-Z}) \Pr_{\xi}(\hat{h}_0^t | M_{X,-Z}, \hat{h}_1^t)}{\Pr_{\xi}(\hat{h}_1^t | M_{X,Z}) \Pr_{\xi}(\hat{h}_0^t | M_{X,Z})}, \quad (\text{C.9})$$

where, recall,

$$\begin{aligned} \hat{h}_1^t &\equiv (y_{\tau}, x_{\tau}, \hat{z}_{\tau})_{\tau < t: \hat{z}_{\tau} \neq \emptyset} \\ \hat{h}_0^t &\equiv (y_{\tau}, x_{\tau}, \hat{z}_{\tau})_{\tau < t: \hat{z}_{\tau} = \emptyset}. \end{aligned}$$

From (C.9), we can write the left hand side of (C.8) as

$$\begin{aligned} &\frac{1}{t-1} \left[ \log \left( \frac{\Pr_{\xi}(\hat{h}_1^t | M_{X,-Z})}{\prod_{k \in \mathcal{E}(t)} p_0^1(y_k, x_k, z_k)} \right) + \log \left( \frac{\Pr_{\xi}(\hat{h}_0^t | M_{X,-Z}, \hat{h}_1^t)}{\prod_{k \notin \mathcal{E}(t)} p_0^0(y_k, x_k, \emptyset)} \right) \right] \\ &- \frac{1}{t-1} \left[ \log \left( \frac{\Pr_{\xi}(\hat{h}_1^t | M_{X,Z})}{\prod_{k \in \mathcal{E}(t)} p_0^1(y_k, x_k, z_k)} \right) + \log \left( \frac{\Pr_{\xi}(\hat{h}_0^t | M_{X,Z})}{\prod_{k \notin \mathcal{E}(t)} p_0^0(y_k, x_k, \emptyset)} \right) \right] \end{aligned} \quad (\text{C.10})$$

We know that the second term of (C.10) tends almost surely towards 0 as  $t \rightarrow \infty$  by Lemma D2.

As a result, to establish (C.8) it remains to show that the first term tends almost surely towards  $-d$ . Rewrite this term as

$$\begin{aligned} &\frac{\#\mathcal{E}(t)}{t-1} \left[ \frac{1}{\#\mathcal{E}(t)} \log \left( \frac{\Pr_{\xi}(\hat{h}_1^t | M_{X,-Z})}{\prod_{k \in \mathcal{E}(t)} p_0^1(y_k, x_k, z_k)} \right) \right] \\ &+ \frac{t-1-\#\mathcal{E}(t)}{t-1} \left[ \frac{1}{t-1-\#\mathcal{E}(t)} \log \left( \frac{\Pr_{\xi}(\hat{h}_0^t | M_{X,-Z}, \hat{h}_1^t)}{\prod_{k \notin \mathcal{E}(t)} p_0^0(y_k, x_k, \emptyset)} \right) \right]. \end{aligned} \quad (\text{C.11})$$

By Lemmas D2, C5, and C6, (C.11) tends almost surely towards

$$1 \times -d + 0 \times 0 = -d$$

as  $t \rightarrow \infty$ , which completes the proof.  $\square$

LEMMA C8. *If the continuous attention assumptions hold, then*

$$\frac{\mathcal{B}_{-X,-Z}(\hat{h}^t)}{e^{-d'(t-1)}} \xrightarrow{a.s.} K$$

for some  $d' \geq d$  and  $K$  satisfying  $0 < K < \infty$ .

*Proof.* Let  $d' = \delta_{-X,-Z}(p_0^1)$ . Using analagous arguments to those in the proof of Lemma C7, we can show that

$$\frac{\mathcal{B}_{-X,-Z}(\hat{h}^t)}{e^{-\delta_{-X,-Z}(p_0^1)(t-1)}} \xrightarrow{a.s.} K$$

for some  $K$  satisfying  $0 < K < \infty$ . Since  $\delta_{-X,-Z}(p_0^1) > \delta_{X,-Z}(p_0^1)$  (from the fact that adding more constraints weakly increases the minimized Kullback-Leibler divergence) and  $\delta_{X,-Z}(p_0^1) = d$  (by Lemma C6), the result follows.  $\square$

LEMMA C9. *Suppose the continuous attention assumptions hold. If the asymptotic rate of convergence of  $\hat{\pi}_Z^t$  to 1 is  $V(t)$  then the asymptotic rate of convergence of  $\eta(\hat{\pi}_Z^t)$  to 1 is  $V(t)$ .*

*Proof.* Suppose the asymptotic rate of convergence of  $\hat{\pi}_Z^t$  to 1 is  $V(t)$ . Then, by definition, there must exist a strictly positive constant  $C < \infty$  such that

$$\frac{1 - \hat{\pi}_Z^t}{V(t)} \xrightarrow{a.s.} C. \quad (\text{C.12})$$

The goal is to show that (C.12) implies that there exists a strictly positive constant  $C' < \infty$  such that

$$\frac{1 - \eta(\hat{\pi}_Z^t)}{V(t)} \xrightarrow{a.s.} C'. \quad (\text{C.13})$$

But (C.13) follows from (C.12) so long as there exists a strictly positive constant  $K < \infty$  such that

$$\frac{1 - \eta(\hat{\pi}_Z^t)}{1 - \hat{\pi}_Z^t} \xrightarrow{a.s.} K, \quad (\text{C.14})$$

which can easily be verified using l'Hopital's rule, together with the assumptions that  $\eta' = f$  is continuous and strictly positive on  $[0, 1]$ .  $\square$

*Proof of Proposition 5.* From Lemma C9, it is enough to show that

$$\frac{1 - \hat{\pi}_Z^t}{e^{-d(t-1)}} \xrightarrow{a.s.} C \quad (\text{C.15})$$

for some strictly positive constant  $C < \infty$ .

Since

$$1 - \hat{\pi}_Z^t = \frac{1}{1 + \frac{1 + \frac{1-\pi_X}{\pi_X} \mathcal{B}_{-X,Z}(\hat{h}^t)}{\frac{1-\pi_Z}{\pi_Z} \mathcal{B}_{X,-Z}(\hat{h}^t) + \frac{1-\pi_X}{\pi_X} \frac{1-\pi_Z}{\pi_Z} \mathcal{B}_{-X,-Z}(\hat{h}^t)}},$$

to demonstrate (C.15) it suffices to show that

$$\frac{e^{-d(t-1)}}{\frac{1-\pi_Z}{\pi_Z} \mathcal{B}_{X,-Z}(\hat{h}^t) + \frac{1-\pi_X}{\pi_X} \frac{1-\pi_Z}{\pi_Z} \mathcal{B}_{-X,-Z}(\hat{h}^t)} + \frac{\frac{1-\pi_X}{\pi_X} \mathcal{B}_{-X,Z}(\hat{h}^t) e^{-d(t-1)}}{\frac{1-\pi_Z}{\pi_Z} \mathcal{B}_{X,-Z}(\hat{h}^t) + \frac{1-\pi_X}{\pi_X} \frac{1-\pi_Z}{\pi_Z} \mathcal{B}_{-X,-Z}(\hat{h}^t)} \xrightarrow{a.s.} c' \quad (\text{C.16})$$

for some constant  $c'$  satisfying  $0 < c' < \infty$ .

The first term on the left hand side of (C.16) converges almost surely to a positive finite constant by Lemmas C7 and C8. The second term on the left hand side of (C.16) converges almost surely to 0 whenever  $\pi_X = 1$  (trivially) or  $x$  is important to predicting  $y$  (by Lemmas D3, C7, and C8).  $\square$

### C.5. Proofs of Results from Appendix B: Sophistication

*Proof of Proposition B1.* I prove by example. We will consider the case where  $\pi_X = 1$ , so

$$\hat{\pi}_Z^t = \pi_Z \frac{\Pr(\hat{h}^t | M_{X,Z})}{\Pr(\hat{h}^t | M_{X,-Z})}.$$

It is enough to find values of primitives such that there exists some history of the form  $\hat{h}^t = (y_{t-1}, x_{t-1}, \emptyset, \dots, y_1, x_1, \emptyset)$ , where, for the sophisticated agent,

$$\begin{aligned} \rho &= \frac{\Pr(\hat{h}^t | M_{X,Z})}{\Pr(\hat{h}^t | M_{X,-Z})} \\ &= \frac{\int_{\theta} \prod_{\tau=1}^{t-1} p_{\theta}(y_{\tau} | x_{\tau}) d\mu^{X,Z}}{\int_{\theta} \prod_{\tau=1}^{t-1} p_{\theta}(y_{\tau} | x_{\tau}) d\mu^{X,-Z}} \\ &= \frac{\int_{\theta} \prod_{x'} (\sum_{z'} \theta(x', z') g(z' | x'))^{\kappa'(x')} (1 - \sum_{z'} \theta(x', z') g(z' | x'))^{N'(x') - \kappa'(x')} d\mu^{X,Z}}{\int_{\theta} \prod_{x'} \theta(x')^{\kappa'(x')} (1 - \theta(x'))^{N'(x') - \kappa'(x')} d\mu^{X,-Z}} \\ &> 1. \end{aligned}$$

To this end, consider an example with  $Z = \{1, 2\}$ ,  $g(1|x) = 1/2$  for all  $x$ , and  $\psi(\theta) \sim U[0, 1]$ . For any length-2 history of the form  $(y_2, x, \emptyset, y_1, x, \emptyset)$ , we have

$$\rho = \frac{\int (.5(\tau_1 + \tau_2))^{y_1+y_2} (1 - .5(\tau_1 + \tau_2))^{2-y_1-y_2} d\tau_1 d\tau_2}{\int \tau^{y_1+y_2} (1 - \tau)^{2-y_1-y_2} d\tau}.$$

Calculating  $\rho$  in this example for the case where  $y_1 + y_2 = 1$  yields  $\rho = 1.25$ .  $\square$

*Proof of Lemma B1.* Note the inequality

$$\Lambda(\hat{h}^t) \leq \frac{1}{\mathcal{B}_{X,-Z}(\hat{h}^t)} + \frac{1 - \pi_X}{\pi_X} \frac{\mathcal{B}_{-X,Z}(\hat{h}^t)}{\mathcal{B}_{X,-Z}(\hat{h}^t)}. \quad (\text{C.17})$$

We will proceed by bounding the two terms on the right hand side of (C.17) above by a term that does not depend on  $t, \hat{h}^t$ .

When  $\hat{h}^t = (y_{t-1}, x_{t-1}, \emptyset, \dots, y_1, x_1, \emptyset)$ , we can write

$$\begin{aligned} \mathcal{B}_{X,-Z}(\hat{h}^t) &= \frac{\int \prod_{x'} \theta^X(x')^{\kappa^t(x')} (1 - \theta^X(x'))^{N^t(x') - \kappa^t(x')} \psi_{X,-Z}^X(\theta^X) d\theta^X}{\int \prod_{x'} \theta^X(x')^{\kappa^t(x')} (1 - \theta^X(x'))^{N^t(x') - \kappa^t(x')} \psi_{X,Z}^X(\theta^X) d\theta^X} \\ &\geq a_1, \end{aligned}$$

where  $a_1 \in (0, \infty)$  is a constant that does not depend on  $t, \hat{h}^t$ . The inequality holds because  $\psi_{X,-Z}^X(\theta^X) = \prod_{x'} \psi(\theta^X(x'))$  is bounded below by a positive constant and  $\psi_{X,Z}^X(\theta^X)$  is bounded above by a positive constant. As a result, the first term on the right hand side of (C.17) is bounded above by a term that does not depend on  $t, \hat{h}^t$ .

Also, when  $\hat{h}^t = (y_{t-1}, x_{t-1}, \emptyset, \dots, y_1, x_1, \emptyset)$ , we can write

$$\begin{aligned} \frac{\mathcal{B}_{-X,Z}(\hat{h}^t)}{\mathcal{B}_{X,-Z}(\hat{h}^t)} &= \frac{\int \prod_{x'} \theta^X(x')^{\kappa^t(x')} (1 - \theta^X(x'))^{N^t(x') - \kappa^t(x')} \psi_{-X,Z}^X(\theta^X) d\theta^X}{\int \prod_{x'} \theta^X(x')^{\kappa^t(x')} (1 - \theta^X(x'))^{N^t(x') - \kappa^t(x')} \psi_{X,-Z}^X(\theta^X) d\theta^X} \\ &\leq a_2, \end{aligned}$$

where  $a_2 \in (0, \infty)$  is a constant that does not depend on  $t, \hat{h}^t$ . The inequality follows from the fact that  $\psi_{X,-Z}^X$  is bounded below by a positive constant and  $\psi_{-X,Z}^X$  is bounded above by a constant. We conclude that the second term on the right hand side of (C.17) is bounded above by a term that does not depend on  $t, \hat{h}^t$ , which completes the proof.  $\square$

*Proof of Proposition B2.* It is the case that

$$\hat{\pi}_Z^t \leq b \iff \Lambda(\hat{h}^t) \leq \frac{1 - \pi_Z}{\pi_Z} \frac{b}{1 - b},$$

so the agent will never encode  $z$  if

$$\Lambda(\hat{h}^t) \leq \frac{1 - \pi_Z}{\pi_Z} \frac{b}{1 - b} \quad (\text{C.18})$$

for all  $t, \hat{h}^t$ .

From Lemma B1, we have that  $\Lambda(\hat{h}^t) \leq K$  for some  $K \in (0, \infty)$ , where this bound holds for all  $t, \hat{h}^t$ . Thus, Condition (C.18) holds so long as

$$K \leq \frac{1 - \pi_Z}{\pi_Z} \frac{b}{1 - b} \iff \pi_Z \leq \frac{b/(1 - b)}{K + b/(1 - b)} \equiv \hat{\pi}_Z. \quad (\text{C.19})$$

$\square$

*Proof of Proposition B3.* The statement regarding the naive agent was established in Proposition 4, so we restrict attention to the sophisticate. We can write

$$\begin{aligned}\hat{\pi}_X^t &= \frac{1 + \alpha_{X,-Z} \mathcal{B}_{X,-Z}(\hat{h}^t)}{1 + \alpha_{X,-Z} \mathcal{B}_{X,-Z}(\hat{h}^t) + \alpha_{-X,Z} \mathcal{B}_{-X,Z}(\hat{h}^t) + \alpha_{-X,-Z} \mathcal{B}_{-X,-Z}(\hat{h}^t)} \\ &\leq \frac{1 + \alpha_{X,-Z} \mathcal{B}_{X,-Z}(\hat{h}^t)}{1 + \alpha_{X,-Z} \mathcal{B}_{X,-Z}(\hat{h}^t) + \alpha_{-X,Z} \mathcal{B}_{-X,Z}(\hat{h}^t)}.\end{aligned}$$

From the inequality, it suffices to show that, almost surely,  $\mathcal{B}_{-X,Z}(\hat{h}^t)$  does not vanish as  $t$  gets large conditional on  $\pi_Z < \bar{\pi}_Z$ .<sup>7</sup>

When the agent never encodes  $z$ , we can write

$$\begin{aligned}\mathcal{B}_{-X,Z}(\hat{h}^t) &= \frac{\int \prod_{x'} \theta^X(x')^{\kappa^t(x')} (1 - \theta^X(x'))^{N^t(x') - \kappa^t(x')} \psi_{-X,Z}^X(\theta^X) d\theta^X}{\int \prod_{x'} \theta^X(x')^{\kappa^t(x')} (1 - \theta^X(x'))^{N^t(x') - \kappa^t(x')} \psi_{X,Z}^X(\theta^X) d\theta^X} \\ &= \frac{\left( \frac{\int \prod_{x'} \theta^X(x')^{\kappa^t(x')} (1 - \theta^X(x'))^{N^t(x') - \kappa^t(x')} \psi_{-X,Z}^X(\theta^X) d\theta^X}{\int \prod_{x'} \theta^X(x')^{\kappa^t(x')} (1 - \theta^X(x'))^{N^t(x') - \kappa^t(x')} d\theta^X} \right)}{\left( \frac{\int \prod_{x'} \theta^X(x')^{\kappa^t(x')} (1 - \theta^X(x'))^{N^t(x') - \kappa^t(x')} \psi_{X,Z}^X(\theta^X) d\theta^X}{\int \prod_{x'} \theta^X(x')^{\kappa^t(x')} (1 - \theta^X(x'))^{N^t(x') - \kappa^t(x')} d\theta^X} \right)} \\ &= \frac{E^\lambda[\psi_{-X,Z}^X(\theta^X)|\hat{h}^t]}{E^\lambda[\psi_{X,Z}^X(\theta^X)|\hat{h}^t]} \\ &\geq \frac{E^\lambda[\psi_{-X,Z}^X(\theta^X)|\hat{h}^t]}{\bar{g}}\end{aligned}$$

for some constant  $\bar{g} \in (0, \infty)$ , where  $E^\lambda$  is the expectation operator taken with respect to the Lebesgue prior over  $\theta^X$ . The inequality follows from the fact that  $\psi_{X,Z}^X$  is bounded above by some finite positive constant. From this inequality, it suffices to show that, almost surely,  $E^\lambda[\psi_{-X,Z}^X(\theta^X)|\hat{h}^t]$  does not vanish.

From Freedman (1963),  $\lambda_t \rightarrow \delta_{\theta_0^X}$ , almost surely (weak-star), where  $\lambda_t$  denotes the posterior over  $\theta^X$  given prior  $\lambda$  and history  $\hat{h}^t$ , and  $\delta_{\theta_0^X}$  denotes a point mass at  $\theta_0^X$ . Since  $\psi_{-X,Z}^X$  is continuous in a neighborhood of  $\theta_0^X$  (by assumption), this implies that  $E^\lambda[\psi_{-X,Z}^X(\theta^X)|\hat{h}^t] \rightarrow \psi_{-X,Z}^X(\theta_0^X)$  almost surely. Since  $\theta_0^X$  is in the support of  $\mu_{-X,Z}$  (because  $x$  is unimportant to predicting  $y$ ), we have  $\psi_{-X,Z}^X(\theta_0^X) > 0$  (by assumption).  $\square$

7. Recall that statements regarding almost sure convergence are with respect to the measure over infinite-horizon mentally represented histories as generated by  $\theta_0, g$ , and the agent's encoding rule. When  $\pi_Z < \bar{\pi}_Z$ , the relevant measure is the distribution over  $(y_\tau, x_\tau, \mathcal{O})_{\tau=1}^\infty$ , as implied by  $\theta_0, g$ .

## Appendix D: Bayes' Factors

This appendix is broken up into four parts. First, I establish several results which are useful in establishing asymptotic properties of Bayes' factors generally, and in turn aid in characterizing the agent's asymptotic forecasts and beliefs. Next, I consider the evolution of Bayes' factors for three special cases: when the agent always encodes  $z$ , when the agent never encodes  $z$ , and when the agent encodes  $z$  infinitely often.

Let  $p_0(y, x, \hat{z})$  and  $\hat{p}(y, x, \hat{z})$  denote probability mass functions over  $(y, x, \hat{z}) \in \{0, 1\} \times X \times \hat{Z}$ . Define the Kullback Leibler distance between  $\hat{p}(y, x, \hat{z})$  and  $p_0(y, x, \hat{z})$  as

$$d_K(\hat{p}, p_0) = \sum_{y, x, \hat{z}} p_0(y, x, \hat{z}) \log \left( \frac{p_0(y, x, \hat{z})}{\hat{p}(y, x, \hat{z})} \right) \quad (\text{D.1})$$

with the convention that  $0 \log(0/\hat{p}) = 0$  for  $\hat{p} \geq 0$  and  $p_0 \log(p_0/0) = \infty$  for  $p_0 > 0$  (see, e.g., Cover and Thomas 2006).

For all  $(y, x, \hat{z})$ , assume that  $\hat{p}(y, x, \hat{z})$  can be written as  $\hat{p}(y, x, \hat{z}|\theta) = \theta(x, \hat{z})^y (1 - \theta(x, \hat{z}))^{1-y} p_0(x, \hat{z})$  (sometimes abbreviated as  $\hat{p}_\theta(y, x, \hat{z})$ ), where  $p_0(x, \hat{z}) = \sum_{y' \in \{0, 1\}} p_0(y', x, \hat{z})$ . Define  $\hat{p}(y, x, \hat{z}|\theta(M)) = p_{\theta(M)}(y, x, \hat{z})$  in the obvious manner ( $\theta(M)$  is defined as in Subsection A.1.1) and let  $\underline{\theta}(M) = \arg \min_{\theta(M) \in \Theta(M)} d_K(\hat{p}_{\theta(M)}, p_0)$  denote a minimizer of the Kullback-Leibler distance between  $\hat{p}_{\theta(M)}(\cdot)$  and  $p_0(\cdot)$  among parameter values in the support of  $\mu^M(\cdot)$ . Finally, define  $\delta_M = \delta_M(p_0) = d_K(\hat{p}_{\underline{\theta}(M)}, p_0)$ .

LEMMA D1. For all  $M \in \mathcal{M}$ ,  $p_0$ , and  $c^M \in C^M$ ,  $\underline{\theta}(c^M) = p_0(y = 1|c^M)$ .

*Proof.* Fix some  $p_0(\cdot)$ ,  $M$ , and  $c^M$ . We have

$$-d_K(\hat{p}_{\theta(M)}, p_0) = \sum_{y, x, \hat{z}} p_0(y|x, \hat{z}) p_0(x, \hat{z}) \log \left( \frac{\theta(c^M(x, \hat{z}))^y (1 - \theta(c^M(x, \hat{z})))^{1-y}}{p_0(y|x, \hat{z})} \right).$$

The right-hand-side of this expression can be re-written as

$$\sum_{c^M \in C^M} [p_0(y = 1|c^M) p_0(c^M) \log(\theta(c^M)) + p_0(y = 0|c^M) p_0(c^M) \log(1 - \theta(c^M))] - K, \quad (\text{D.2})$$

where  $K$  does not depend on  $\theta(M)$ . It is routine to show that each term in the sum of (D.2) is maximized when  $\theta(c^M) = p_0(y = 1|c^M)$ , which concludes the proof.  $\square$

Let  $\hat{h}^t = (y_{t-1}, x_{t-1}, \hat{z}_{t-1}, \dots, y_1, x_1, \hat{z}_1)$  be some random sample from  $p_0(y, x, \hat{z})$ . Define

$$\mathcal{I}_t(M) = \mathcal{I}(M|\hat{h}^t) = \int \frac{\prod_{k=1}^{t-1} \hat{p}(y_k, x_k, \hat{z}_k|\theta)}{\prod_{k=1}^{t-1} p_0(y_k, x_k, \hat{z}_k)} \mu^M(d\theta) \quad (\text{D.3})$$

as well as the predictive distribution

$$\hat{p}_t^M(y, x, \hat{z}) = \hat{p}^M(y, x, \hat{z}|\hat{h}^t) = \int \hat{p}(y, x, \hat{z}|\theta) \mu^M(d\theta|\hat{h}^t). \quad (\text{D.4})$$

Note that, while not explicit in the notation, both  $\mathcal{S}_t(M)$  and  $\hat{p}_t^M(\cdot)$  depend on  $p_0$ . To avoid confusion, I will sometimes make this dependence explicit by writing  $\mathcal{S}_t(M|p_0)$  and  $\hat{p}_t^M(\cdot|p_0)$ .

It will be useful to establish some lemmas with priors which are slightly more general than what has been assumed.

**DEFINITION D1.**  $\mu^M$  is *uniformly non-doctrinaire* if it makes each  $\theta(c^M)$  independent with non-doctrinaire prior  $\psi_{c^M}$ .

Note that it is possible for  $\psi_{c^M}$  to vary with  $c^M$  when  $\mu^M$  is uniformly non-doctrinaire.

**LEMMA D2.** For all  $M \in \mathcal{M}$ ,  $p_0$ , and uniformly non-doctrinaire  $\mu^M$ ,

$$\frac{1}{t-1} \log \mathcal{S}_t(M|p_0) \rightarrow -\delta_M(p_0), \quad (\text{D.5})$$

$p_0^\infty$  almost surely.

*Proof.* Fix some  $M \in \mathcal{M}$ ,  $p_0$ , and uniformly non-doctrinaire  $\mu^M$ . From Walker (2004, Theorem 2), it is sufficient to show that the following conditions hold:

1.  $\mu^M(\{\theta : d_K(\hat{p}_\theta, p_0) < d\}) > 0$  only for, and for all,  $d > \delta_M$
2.  $\lim_t \inf d_K(\hat{p}_t^M, p_0) \geq \delta_M$ ,  $p_0^\infty$  almost surely
3.  $\sup_t \text{Var}(\log(\mathcal{S}_{t+1}(M)/\mathcal{S}_t(M))) < \infty$

The “only for” part of the first condition holds trivially from the definition of  $\delta_M$  and the “for all” part follows from the fact that  $d_K(\hat{p}_{\theta(M)}, p_0)$  is continuous in a neighborhood of  $\theta(M)$  (since  $\hat{p}_{\theta(M)}(\cdot)$  is continuous in  $\theta(M)$ ) and  $\mu^M(\cdot)$  places positive probability on all open neighborhoods in  $\Theta(M)$ . The second condition also holds trivially since  $d_K(\hat{p}_t^M, p_0) \geq \min_{\theta(M) \in \Theta(M)} d_K(\hat{p}_{\theta(M)}, p_0) = \delta_M$  for all  $t, \hat{h}^t$ .

The third condition requires a bit more work to verify. Note that

$$\mathcal{S}_{t+1}(M) = \frac{\hat{p}_t^M(y_t, x_t, \hat{z}_t)}{p_0(y_t, x_t, \hat{z}_t)} \mathcal{S}_t(M) \Rightarrow \log(\mathcal{S}_{t+1}(M)/\mathcal{S}_t(M)) = \log\left(\frac{\hat{p}_t^M(y_t, x_t, \hat{z}_t)}{p_0(y_t, x_t, \hat{z}_t)}\right),$$

so condition (3) is equivalent to

$$\sup_t \text{Var}\left[\log\left(\frac{\hat{p}_t^M(y_t, x_t, \hat{z}_t)}{p_0(y_t, x_t, \hat{z}_t)}\right)\right] < \infty \quad (\text{D.6})$$

which can easily be shown to hold so long as

$$\sup_t E\left\{\sum_{y,x,\hat{z}} p_0(y,x,\hat{z}) \log\left(\frac{\hat{p}_t^M(y|x,\hat{z})}{p_0(y|x,\hat{z})}\right)^2\right\} < \infty \quad (\text{D.7})$$

or

$$\sup_t E\left[\log\left(\hat{p}_t^M(y|x,\hat{z})\right)^2\right] < \infty \quad (\text{D.8})$$

for all  $(y, x, \hat{z})$  which satisfy  $p_0(y, x, \hat{z}) > 0$ .

To verify (D.8), fix some  $(y, x, \hat{z})$  with  $p_0(y, x, \hat{z}) > 0$  and let  $N(c^M(x, \hat{z})|\hat{h}^t) = N_t$  denote the number of times the covariates have taken on some value  $(x', \hat{z}') \in c^M(x, \hat{z})$  along history  $\hat{h}^t$  and  $\kappa(c^M(x, \hat{z})|\hat{h}^t) = \kappa_t$  denote the number of times both that the covariates have taken on such a value and  $y = 1$ . Then

$$q_t = \frac{\kappa_t + 1}{N_t + 2} \quad (\text{D.9})$$

roughly equals the empirical frequency of  $y = 1$  conditional on  $(x', \hat{z}') \in c^M(x, \hat{z})$  up to period  $t$ .

An implication of the Theorem in Diaconis and Freedman (1990) is that

$$|\hat{p}_t^M(y|x, \hat{z}) - q_t^y(1 - q_t)^{1-y}| \rightarrow 0 \quad (\text{D.10})$$

at a uniform rate across histories since the marginal prior density over  $\theta(c^M(x, \hat{z}))$  is non-doctrinaire. Consequently, fixing an  $\varepsilon > 0$  there exists an  $n > 0$  such that, independent of the history,

$$|\log(\hat{p}_t^M(y|x, \hat{z}))^2 - \log(q_t^y(1 - q_t)^{1-y})^2| < \varepsilon$$

for all  $t \geq n^8$  which implies that

$$E[|\log(\hat{p}_t^M(y|x, \hat{z}))^2 - \log(q_t^y(1 - q_t)^{1-y})^2|] < \varepsilon \quad (\text{D.11})$$

for all  $t \geq n$ . Since  $E[\log(\hat{p}_t^M(y|x, \hat{z}))^2] < \infty$  for all finite  $t$ , to verify (D.8) it is sufficient to show that

$$\sup_t E[\log(q_t^y(1 - q_t)^{1-y})^2] < \infty \quad (\text{D.12})$$

by (D.11).

By symmetry, it is without loss of generality to verify (D.12) for the case where  $y = 1$ . To this end,

$$\begin{aligned} E[\log(q_t)^2] &= E[E[\log(q_t)^2|N_t]] \\ &= E\left[\left((1 + N_t)(1 - \tilde{\theta})^{N_t} \log\left(\frac{1}{2 + N_t}\right)\right)^2\right] \end{aligned}$$

where

$$\tilde{\theta} \equiv p_0(y = 1|c^M(x, \hat{z})).$$

---

8. One can show that this statement follows from Diaconis and Freedman's (1990) result using an argument similar to Fudenberg and Levine (1993, Proof of Lemma B.1).

Now, since  $\lim_{N \rightarrow \infty} (1+N)(1-\tilde{\theta})^N \log\left(\frac{1}{2+N}\right)^2 = 0$ , there exists a constant  $M < \infty$  such that

$$(1+N)(1-\tilde{\theta})^N \log\left(\frac{1}{2+N}\right)^2 < M$$

for all  $N$ . As a result,

$$E \left[ (1+N_t)(1-\tilde{\theta})^{N_t} \log\left(\frac{1}{2+N_t}\right)^2 \right] < M < \infty$$

for all  $t$  which verifies (D.12) and concludes the proof.  $\square$

Define the Bayes' factor conditional on  $p_0$  as

$$\mathcal{B}_{i,j}(\hat{h}^t | p_0) = \mathcal{B}_{i,j}^t(p_0) = \frac{\mathcal{J}_t(M_{i,j} | p_0)}{\mathcal{J}_t(M_{X,Z} | p_0)} \quad (\text{D.13})$$

Note that  $\mathcal{B}_{i,j}(\hat{h}^t) = \mathcal{B}_{i,j}(\hat{h}^t | p_0)$  for some  $p_0$  whenever we can write  $\Pr_{\xi}(\hat{h}^t | \theta) = \prod_{k=1}^{t-1} \hat{p}(y_k, x_k, \hat{z}_k | \theta) = \prod_{k=1}^{t-1} \theta(x_k, \hat{z}_k)^{y_k} (1-\theta(x_k, \hat{z}_k))^{1-y_k} p_0(x_k, \hat{z}_k)$  for some  $p_0$ .

LEMMA D3. For all  $M_{i,j} \in \mathcal{M}$ ,  $p_0$ , and uniformly non-doctrinaire  $\mu^{i,j}, \mu^{X,Z}$ ,

$$\frac{1}{t-1} \log \mathcal{B}_{i,j}^t(p_0) \rightarrow \delta_{X,Z}(p_0) - \delta_{i,j}(p_0), \quad (\text{D.14})$$

$p_0^\infty$  almost surely.

*Proof.* Note that

$$\frac{1}{t-1} \log \mathcal{B}_{i,j}^t = \frac{1}{t-1} \log(\mathcal{J}_t(M_{i,j})) - \frac{1}{t-1} \log(\mathcal{J}_t(M_{X,Z})) \quad (\text{D.15})$$

so the result follows immediately from Lemma D2.  $\square$

REMARK D1. An immediate implication of Lemma D3 is that  $\delta_{i,j}(p_0) > \delta_{X,Z}(p_0)$  implies  $\mathcal{B}_{i,j}^t(p_0) \rightarrow 0$ ,  $p_0^\infty$  almost surely.

Remark D1 applies when  $\delta_{i,j}(p_0) > \delta_{X,Z}(p_0)$ ; what does the Bayes' factor  $\mathcal{B}_{i,j}^t(p_0)$  tend towards asymptotically when  $\delta_{i,j}(p_0) = \delta_{X,Z}(p_0)$ ? I now present a Lemma (due to Diaconis and Freedman 1992) that will aid in estimating the Bayes' factor in this case and establishing asymptotic results. First some definitions. Let  $H(q)$  be the entropy function  $q \log(q) + (1-q) \log(1-q)$  (set at 0 for  $q = 0$  or 1) and define

the following

$$\begin{aligned}\varphi(\kappa, N, \psi) &= \int_0^1 \theta^\kappa (1-\theta)^{N-\kappa} \psi(\theta) d\theta \\ \varphi(\kappa, N) &= \int_0^1 \theta^\kappa (1-\theta)^{N-\kappa} d\theta = \text{Beta}(\kappa+1, N-\kappa+1) \\ \hat{q} &= \frac{\kappa}{N} \\ \varphi^*(\kappa, N) &= \begin{cases} \frac{e^{N H(\hat{q})}}{\sqrt{N}} \sqrt{2\pi} \sqrt{\hat{q}(1-\hat{q})} & \text{for } 0 < \kappa < N \\ \frac{1}{N} & \text{for } \kappa = 0 \text{ or } N \end{cases}\end{aligned}$$

LEMMA D4. For any non-doctrinaire  $\psi(\cdot)$  there are  $0 < a < A < \infty$  such that for all  $N = 1, 2, \dots$  and  $\kappa = 0, 1, \dots, N$ ,  $a\varphi^*(\kappa, N) < \varphi(\kappa, N, \psi) < A\varphi^*(\kappa, N)$ .

*Proof.* Note that for any non-doctrinaire  $\psi$  there exist constants  $b, B$  such that  $0 < b \leq \psi(\theta) \leq B < \infty$  for all  $\theta \in (0, 1)$ . The result then follows from Lemma 3.3(a) in Diaconis and Freedman (1992). For a brief sketch, note that  $b\varphi(\kappa, N) \leq \varphi(\kappa, N, \psi) \leq B\varphi(\kappa, N)$ . Now use Stirling's formula on  $\varphi(\kappa, N)$  for  $\kappa$  and  $N - \kappa$  large.  $\square$

### D.1. Bayes' Factors When the Agent Always Encodes $z$

LEMMA D5.  $\mathcal{B}_{X,-Z}(h^t) \rightarrow 0$  and  $\mathcal{B}_{-X,-Z}(h^t) \rightarrow 0$ ,  $P_{\theta_0}$  almost surely. Additionally, if  $x$  is important to predicting  $y$ , then  $\mathcal{B}_{-X,Z}(h^t) \rightarrow 0$ ,  $P_{\theta_0}$  almost surely; otherwise,  $\mathcal{B}_{-X,Z}(h^t) \rightarrow \infty$ ,  $P_{\theta_0}$  almost surely.

*Proof.* First, I will establish that  $\mathcal{B}_{X,-Z}^t$  and  $\mathcal{B}_{-X,-Z}^t$  converge almost surely to zero, and that  $\mathcal{B}_{-X,Z}^t$  converges almost surely to zero whenever  $x$  is important to predicting  $y$ . When the agent always encodes  $z$ , each period's observation is independently drawn from  $p_0^1(y, x, z) = \theta_0(x, z)^y (1 - \theta_0(x, z))^{1-y} g(x, z)$  for all  $(y, x, z)$ . Then, Lemma D3 implies that it is sufficient to show that  $\delta_{X,-Z}(p_0^1) > \delta_{X,Z}(p_0^1)$ ,  $\delta_{-X,-Z}(p_0^1) > \delta_{X,Z}(p_0^1)$  and, whenever  $x$  is important to predicting  $y$ ,  $\delta_{-X,Z}(p_0^1) > \delta_{X,Z}(p_0^1)$ . We can easily establish these inequalities by applying Lemma D1 for each  $M \in \mathcal{M}$ .

It is left to show that  $\mathcal{B}_{-X,Z}^t \xrightarrow{a.s.} \infty$  whenever  $x$  is not important to predicting  $y$ . First, write out the Bayes' factor:

$$\mathcal{B}_{-X,Z}^t = \frac{\Pr(h^t | M_{-X,Z})}{\Pr(h^t | M_{X,Z})} \quad (\text{D.16})$$

$$= \prod_{z'} \frac{\int_0^1 \theta^{\kappa(z'|h^t)} (1-\theta)^{N(z'|h^t) - \kappa(z'|h^t)} \psi(\theta) d\theta}{\prod_{x'} \int_0^1 \theta^{\kappa(x',z'|h^t)} (1-\theta)^{N(x',z'|h^t) - \kappa(x',z'|h^t)} \psi(\theta) d\theta}. \quad (\text{D.17})$$

From (D.17) it is sufficient to show that

$$\frac{\int_0^1 \theta^{\kappa(z|h^t)} (1-\theta)^{N(z|h^t) - \kappa(z|h^t)} \psi(\theta) d\theta}{\prod_{x'} \int_0^1 \theta^{\kappa(x',z|h^t)} (1-\theta)^{N(x',z|h^t) - \kappa(x',z|h^t)} \psi(\theta) d\theta} \xrightarrow{a.s.} \infty \quad (\text{D.18})$$

for each  $z \in Z$ .

Fix some  $z$ . I will use Lemma D4 to estimate (D.18). Let  $\kappa_t = \kappa(z|h^t)$ ,  $N_t = N(z|h^t)$ ,  $\hat{q}_t = \kappa_t/N_t$ ,  $\kappa_t^{x'} = \kappa(x', z|h^t)$ ,  $N_t^{x'} = N(x', z|h^t)$ , and  $\hat{q}_t^{x'} = \kappa_t^{x'}/N_t^{x'}$ .

Applying Lemma D4,

$$\frac{\int_0^1 \theta^{\kappa_t} (1-\theta)^{N_t-\kappa_t} \psi(\theta) d\theta}{\prod_{x'} \int_0^1 \theta^{\kappa_t^{x'}} (1-\theta)^{N_t^{x'}-\kappa_t^{x'}} \psi(\theta) d\theta} \geq \frac{a\varphi^*(\kappa_t, N_t)}{A^{\#X} \prod_{x'} \varphi^*(\kappa_t^{x'}, N_t^{x'})} \quad (\text{D.19})$$

for some constants  $0 < a < A < \infty$ . By the strong law of large numbers, the right hand side of (D.19) tends almost surely towards

$$C \frac{\sqrt{\prod_{x'} N_t^{x'}}}{\sqrt{N_t}} \xrightarrow{a.s.} \infty$$

where  $C$  is some positive constant independent of  $t$ . □

## D.2. Bayes' Factors When the Agent Never Encodes $z$

Define

$$p_0^0(y, x, \hat{z}) = \begin{cases} \sum_{z' \in Z} \theta_0(x, z')^y (1 - \theta_0(x, z'))^{1-y} g(x, z') & \text{for each } y, x, \text{ and } \hat{z} = \emptyset \\ 0 & \text{for } \hat{z} \neq \emptyset \end{cases} \quad (\text{D.20})$$

to equal the distribution over  $(y, x, \hat{z})$  conditional on the agent not encoding  $z$ . Lemma D6 establishes the almost sure limit of several Bayes' factors when the agent never encodes  $z$ .

**LEMMA D6.** *Suppose  $E_{\theta_0}[y|x] \neq E_{\theta_0}[y]$  for some  $x \in X$ . Then, for all uniformly non-doctrinaire  $\mu^{-X, -Z}$ ,  $\mu^{-X, Z}$ , and  $\mu^{X, Z}$ ,  $\mathcal{B}_{-X, -Z}(\hat{h}^t | p_0^0) \rightarrow 0$  and  $\mathcal{B}_{-X, Z}(\hat{h}^t | p_0^0) \rightarrow 0$ ,  $(p_0^0)^\infty$  almost surely.*

*Proof.* Lemma D3 implies that it is sufficient to show that  $\delta_{-X, -Z}(p_0^0) > \delta_{X, Z}(p_0^0)$  and  $\delta_{-X, Z}(p_0^0) > \delta_{X, Z}(p_0^0)$  whenever  $E_{\theta_0}[y|x] \neq E_{\theta_0}[y]$  for some  $x \in X$ . We can easily verify these inequalities by applying Lemma D1 for each  $M \in \mathcal{M}$ . □

## D.3. Bayes' Factors When the Agent Encodes $z$ Infinitely Often

**LEMMA D7.** *Suppose that, with positive probability under  $P_{\theta_0, \xi}(\cdot)$ , the agent encodes  $z$  infinitely often. Conditional on the agent encoding  $z$  infinitely often,  $\mathcal{B}_{X, -Z}(\hat{h}^t) \rightarrow 0$  and  $\mathcal{B}_{-X, -Z}(\hat{h}^t) \rightarrow 0$  almost surely.*

*Proof.* I will establish that, almost surely,

$$\log(\mathcal{B}_{i, -Z}(\hat{h}^t)) \rightarrow -\infty \quad (\text{D.21})$$

for each  $i \in \{X, \neg X\}$ .

Defining

$$\begin{aligned}\hat{h}'_1 &\equiv (y_\tau, x_\tau, \hat{z}_\tau)_{\tau < t: \hat{z}_\tau \neq \emptyset} \\ \hat{h}'_0 &\equiv (y_\tau, x_\tau, \hat{z}_\tau)_{\tau < t: \hat{z}_\tau = \emptyset}\end{aligned}$$

we can write

$$\mathcal{B}_{i, -Z}^t = \frac{\Pr_\xi(\hat{h}'_0 | M_{i, -Z}) \Pr_\xi(\hat{h}'_1 | M_{i, -Z}, \hat{h}'_0)}{\Pr_\xi(\hat{h}'_0 | M_{X, Z}) \Pr_\xi(\hat{h}'_1 | M_{X, Z}, \hat{h}'_0)}$$

so the LHS of (D.21) can be expressed as

$$\log \left( \frac{\Pr_\xi(\hat{h}'_0 | M_{i, -Z})}{\Pr_\xi(\hat{h}'_0 | M_{X, Z})} \right) + \log \left( \frac{\Pr_\xi(\hat{h}'_1 | M_{i, -Z}, \hat{h}'_0)}{\Pr_\xi(\hat{h}'_1 | M_{X, Z}, \hat{h}'_0)} \right). \quad (\text{D.22})$$

When the agent fails to encode  $z$  only a finite number of times along a history, we can ignore the first term of (D.22) because it tends towards a finite value as  $t \rightarrow \infty$ . Otherwise, Lemma C3 says that the first term of (D.22) is identically 0 for  $i = X$ , as well as for  $i = \neg X$  when  $X$  is a singleton; Lemma D6 (together with Assumption 2) says that the first term tends towards  $-\infty$  with probability 1 for  $i = \neg X$  when  $X$  contains at least two elements. As a result, no matter which case we are in it is sufficient to show that the second term of (D.22) tends towards  $-\infty$  with probability 1 in order to establish (D.21). This can be verified by showing that

$$\limsup_t \frac{1}{\#\mathcal{E}(t)} \log \left( \frac{\Pr_\xi(\hat{h}'_1 | M_{i, -Z}, \hat{h}'_0)}{\prod_{\tau \in \mathcal{E}(t)} p_0^1(y_\tau, x_\tau, z_\tau)} \right) - \frac{1}{\#\mathcal{E}(t)} \log \left( \frac{\Pr_\xi(\hat{h}'_1 | M_{X, Z}, \hat{h}'_0)}{\prod_{\tau \in \mathcal{E}(t)} p_0^1(y_\tau, x_\tau, z_\tau)} \right) < 0 \quad (\text{D.23})$$

with probability 1 for  $i \in \{X, \neg X\}$ , where

$$\begin{aligned}p_0^1(y, x, z) &= \theta_0(x, z)^y (1 - \theta_0(x, z))^{1-y} g(x, z) \\ \mathcal{E}(t) &= \{\tau < t : \hat{z}_\tau \neq \emptyset\}.\end{aligned}$$

The second term on the LHS of (D.23) tends towards 0 with probability 1 by Lemma D2.<sup>9</sup> To complete the proof, it then remains to show that the first term on the LHS of (D.23) remains bounded away from 0 as  $t \rightarrow \infty$  for  $i \in \{X, \neg X\}$ , or

$$\limsup_t \frac{1}{\#\mathcal{E}(t)} \log \left( \frac{\Pr_\xi(\hat{h}'_1 | M_{i, -Z}, \hat{h}'_0)}{\prod_{\tau \in \mathcal{E}(t)} p_0^1(y_\tau, x_\tau, z_\tau)} \right) < 0. \quad (\text{D.24})$$

9. Note that

$$\frac{1}{\#\mathcal{E}(t)} \log \left( \frac{\Pr_\xi(\hat{h}'_1 | M_{X, Z}, \hat{h}'_0)}{\prod_{\tau \in \mathcal{E}(t)} p_0^1(y_\tau, x_\tau, z_\tau)} \right) = \frac{1}{\#\mathcal{E}(t)} \log \left( \frac{\Pr_\xi(\hat{h}'_1 | M_{X, Z})}{\prod_{\tau \in \mathcal{E}(t)} p_0^1(y_\tau, x_\tau, z_\tau)} \right)$$

since, under  $\mu^{X, Z}$ , subjective uncertainty regarding  $\theta(x, \emptyset)$  and  $\theta(x, z')$ ,  $z' \neq \emptyset$ , is independent.

We can re-write the LHS of (D.24) as

$$\frac{1}{\#\mathcal{E}(t)} \log \left( \frac{\prod_{x'} \prod_{z'} \int \theta(x', z')^{\kappa(x', z' | \hat{h}_1^t)} (1 - \theta(x', z'))^{N(x', z' | \hat{h}_1^t) - \kappa(x', z' | \hat{h}_1^t)} \mu^{i, -Z}(d\theta | \hat{h}_0^t)}{\prod_{x'} \prod_{z'} \theta_0(x', z')^{\kappa(x', z' | \hat{h}_1^t)} (1 - \theta_0(x', z'))^{N(x', z' | \hat{h}_1^t) - \kappa(x', z' | \hat{h}_1^t)}} \right). \quad (\text{D.25})$$

Since  $\mu^{i, -Z}(\cdot | \hat{h}_0^t)$  places full support on vectors of success probabilities ( $\theta$ ) with  $\theta(x, z) = \theta(x, z')$  for all  $x, z, z'$ , we can bound (D.25) by noting that

$$\begin{aligned} & \prod_{x'} \prod_{z'} \int \theta(x', z')^{\kappa(x', z' | \hat{h}_1^t)} (1 - \theta(x', z'))^{N(x', z' | \hat{h}_1^t) - \kappa(x', z' | \hat{h}_1^t)} \mu^{i, -Z}(d\theta | \hat{h}_0^t) \\ & \leq \max_{\theta(0), \theta(1)} \prod_{x'} \prod_{z'} \theta(x')^{\kappa(x', z' | \hat{h}_1^t)} (1 - \theta(x'))^{N(x', z' | \hat{h}_1^t) - \kappa(x', z' | \hat{h}_1^t)} \\ & = \prod_{x'} \prod_{z'} \frac{\kappa(x' | \hat{h}_1^t)}{N(x' | \hat{h}_1^t)} \left( 1 - \frac{\kappa(x' | \hat{h}_1^t)}{N(x' | \hat{h}_1^t)} \right)^{N(x', z' | \hat{h}_1^t) - \kappa(x', z' | \hat{h}_1^t)} \end{aligned}$$

which implies that (D.25) is bounded above by

$$\frac{1}{\#\mathcal{E}(t)} \log \left( \frac{\prod_{x'} \prod_{z'} \frac{\kappa(x' | \hat{h}_1^t)}{N(x' | \hat{h}_1^t)} \left( 1 - \frac{\kappa(x' | \hat{h}_1^t)}{N(x' | \hat{h}_1^t)} \right)^{N(x', z' | \hat{h}_1^t) - \kappa(x', z' | \hat{h}_1^t)}}{\prod_{x'} \prod_{z'} \theta_0(x', z')^{\kappa(x', z' | \hat{h}_1^t)} (1 - \theta_0(x', z'))^{N(x', z' | \hat{h}_1^t) - \kappa(x', z' | \hat{h}_1^t)}} \right). \quad (\text{D.26})$$

for all  $t, \hat{h}^t$ . By the strong law of large numbers, expression (D.26) can be shown to tend towards  $-d_K(\hat{p}_{\theta(M_{X, -Z})}, P_0^1) < 0$  with probability 1 conditional on the agent encoding  $z$  infinitely often; this establishes (D.24) and completes the proof.  $\square$