

Online Appendix

A Existence of Equilibrium

Proposition 12 *There exists a PBE.*

Proof: Fix an ex-ante action $a \in A$, and consider the ‘‘restricted game’’ G_a as the continuation game following such action. From Definition 1, a profile of PBEs of all restricted games $\{\phi_a^*, \phi_b^*, m_H^*(a), m_L^*(a), \mu(\cdot | \hat{\sigma}_a)\}_{a \in A}$ is a PBE of the original game considered in the text if and only if

$$\phi_a^* \in \arg \max_{a \in A} \left\{ \begin{array}{l} E_{\hat{\sigma}_a} [E_\mu [u(a, \phi_b^*(a, \hat{\sigma}_a), \theta, \sigma_a) | \hat{\sigma}_a] | m_L^*(a), m_H^*(a)] \\ -q\psi_H(m_H^*(a)) - (1-q)\psi_L(m_L^*(a)) \end{array} \right\}.$$

First, we note that the all continuation games have a PBE. Fix an action $a \in A$ (which, for notational simplicity, will be omitted from the expressions below). Let $b_H \in \arg \max_b u_H(b, H)$, $b_L \in \arg \max_b u_L(b, L)$, and note that both b_H and b_L are not functions of m_H, m_L , or α . Define the correspondence $T : [-\eta_H, 1 - \eta_H] \times [-\eta_L, 1 - \eta_L] \times [0, 1] \times \Delta(B) \rightarrow [-\eta_H, 1 - \eta_H] \times [-\eta_L, 1 - \eta_L] \times [0, 1] \times \Delta(B)$ by

$$T(m_H, m_L, \alpha, \phi_{b_\emptyset}) \equiv (m_H^*(m_L, \phi_{b_\emptyset}, \alpha), m_L^*(m_H, \phi_{b_\emptyset}, \alpha), \alpha(m_L, m_H, \phi_{b_\emptyset}), \phi_{b_\emptyset}(m_L, m_H, \alpha)),$$

where:

$$m_H^* \equiv \arg \max_{m_H} (\eta_H + m_H) u_H(b_H, H) + (1 - \eta_H - m_H) [\alpha^* u_H(\phi_{b_\emptyset}, H) + (1 - \alpha^*) u_L(\phi_{b_\emptyset}, H)] - \psi_H(m_H), \quad (23)$$

$$m_L^* \equiv \arg \max_{m_L} (\eta_L + m_L) u_L(b_L, L) + (1 - \eta_L - m_L) [\alpha^* u_H(\phi_{b_\emptyset}, L) + (1 - \alpha^*) u_L(\phi_{b_\emptyset}, L)] - \psi_L(m_L), \quad (24)$$

$$\alpha^* = \frac{q(1 - \eta_H - m_H^*)}{q(1 - \eta_H - m_H^*) + (1 - q)(1 - \eta_L - m_L^*)}, \quad (25)$$

$$\phi_{b_\emptyset} \equiv \arg \max_{\{\phi(b)\}_{b \in B}} \left\{ \sum_b \phi(b) [\alpha^* u_H(b, H) + (1 - \alpha^*) u_L(b, L)] \right\}. \quad (26)$$

A fixed point of T , combined with $b^*(H) = b_H$ and $b^*(L) = b_L$ constitute a PBE of the restricted game G_a . Equation (23) is a continuous and concave function of m_H . Therefore, m_H^* is a non-empty, convex, and compact set. Similarly, equation (24) is a continuous and concave function of m_L , implying that m_L^* is a non-empty, convex, and compact set. Note that equation (25) defines α^* as a function of m_L, m_H, b_\emptyset . Therefore, α^* is a singleton, which is trivially a non-empty, convex, and compact set. Moreover, since the equation in (26) is a linear function of $\phi(b)$, ϕ_{b_\emptyset} is a non-empty, convex, and compact set. Then, Kakutani’s theorem establishes that there exists a fixed-point of T , which corresponds to a PBE of the restricted game G_a .

Given a profile of PBEs of all restricted games $\{G_a : a \in A\}$, define

$$\alpha^* \in \arg \max_{a \in A} \left\{ \begin{array}{l} E_{\hat{\sigma}_a} [E_\mu [u(a, \phi_b^*(a, \hat{\sigma}_a), \theta, \sigma_a) | \hat{\sigma}_a] | m_L^*(a), m_H^*(a)] \\ -q\psi_H(m_H^*(a)) - (1-q)\psi_L(m_L^*(a)) \end{array} \right\}.$$

Since A is finite, such an α^* exists. The profile $\{\alpha^*, \phi_b^*, m_H^*(a), m_L^*(a), \mu(\cdot | \hat{\sigma}_a)\}$ constructed previously constitutes a PBE of the original game. ■

B Non-Bayesian Framework

Throughout the paper, I have maintained the assumption that decision makers understand that they engage in memory manipulation and, thus, interpret their recollections according to Bayes’ rule (sophisticated decision makers). In this section, I consider the case of naive individuals. As in Mullainathan (2002), naive individuals are unaware of their imperfect memory and interpret recollections as if they were the true outcomes.

Two new features emerge under naiveté. First, unlike in the sophisticated case, the equilibrium is unique. Second, individuals may prefer to observe a signal even if it has no objective value. As a consequence, they may display ambiguity seeking behavior even under additive separability between attributes and money. Moreover, they may exhibit zeroth-order risk seeking behavior.

Consider a naive decision maker (NDM), who is unaware of her memory manipulation efforts. Unawareness implies that she applies Bayes' rule as if her recollections were generated by the case where she does not engage in memory manipulation, i.e. $m_L = m_H = 0$. When $\hat{\sigma} \in \{H, L\}$, she correctly infers that outcome $\sigma = \hat{\sigma}$ has been observed in period 1. However, when an outcome is forgotten, she attributes weight

$$\rho \equiv \frac{q(1 - \eta_H)}{q(1 - \eta_H) + (1 - q)(1 - \eta_L)} \quad (27)$$

to a high outcome and $(1 - \rho)$ to a low outcome. I refer to this rule as a *naive updating rule*.⁶⁰ The following definition proposes an adaptation of the PBE concept to naive decision makers:

Definition 3 A Perfect Naive Equilibrium (PNE) of the game is a strategy profile $(a^*, b^*, m_H^*(a), m_L^*(a))$ and posterior beliefs $\mu(\cdot | \hat{\sigma}_a)$ such that:

1. $a^* \in \arg \max_{a \in A} \left\{ \begin{array}{l} E_{\hat{\sigma}_a} [E_\mu [u(a, b_a^*(\hat{\sigma}_a), \theta, \sigma_a) | \hat{\sigma}_a] | m_L^*(a), m_H^*(a)] \\ - q\psi_H(m_H^*(a)) - (1 - q)\psi_L(m_L^*(a)) \end{array} \right\};$
2. $m_\sigma^*(a) \in \arg \max_{m_\sigma} \left\{ \begin{array}{l} (\eta_\sigma + m_\sigma) E_\mu [u(a, b_a^*(\hat{\sigma}_a), \theta, s) | \hat{\sigma}_a = \sigma] \\ + (1 - \eta_\sigma - m_\sigma) E_\mu [u(a, b_a^*(\hat{\sigma}_a), \theta, \sigma) | \hat{\sigma}_a = \emptyset] - \psi_\sigma(m_\sigma) \end{array} \right\},$
 $\sigma \in \{H, L\};$
3. $b_a^*(\hat{\sigma}) \in \arg \max_{b \in B} \{E_\mu [u(a, b, \theta, \sigma_a) | \hat{\sigma}_a = \hat{\sigma}]\};$
4. $\mu(\theta | \hat{\sigma}_a = \hat{\sigma})$ is obtained by the naive updating rule if $\Pr(\hat{\sigma}_a = \hat{\sigma} | m_L = m_H = 0, a) > 0$, $\hat{\sigma} \in \{L, H, \emptyset\}$, $a \in A$.

Conditions 1 – 3 are the same as in the PBE concept. Condition 4 modifies the standard Bayesian condition by requiring agents to follow the naive updating rule instead. An immediate consequence of the naive updating rule in this model is that it leads to posterior beliefs that first-order stochastically dominate beliefs obtained by Bayes' rule. Therefore, naive individuals are optimistic about their type θ .

An interesting special case of this naive framework is obtained when we take the forgetfulness memory system of Example 1. Recall that if the state \emptyset is interpreted as a recollection of a high outcome, then the model from Example 1 becomes one where the agent is able to convince herself that a low outcome was a high outcome by engaging in memory manipulation. Suppose an individual recollects a high outcome (i.e., $\hat{\sigma} = \emptyset$). If this individual is sophisticated, she then corrects for her memory imperfection and attributes some (Bayesian) weight to the possibility that she has observed a low outcome but managed to convince herself that the outcome was high instead. On the other hand, a naive individual believes her recollection is correct and attributes full weight to a high outcome ($\rho = 1$).

B1 Equilibrium Uniqueness

This subsection establishes that a PNE exists and, under mild conditions, is unique. The naive updating rule implies that the NDM's expected utility given $\hat{\sigma} = \emptyset$ is

$$u_\emptyset(a, b, \sigma) = \rho u_H(a, b, \sigma) + (1 - \rho) u_L(a, b, \sigma). \quad (28)$$

Upon observing an outcome $s \in \{H, L\}$, self 1 maximizes:

$$(\eta_s + m_s) u_s(a, b_a(s), s) + (1 - \eta_s - m_s) u_\emptyset(a, b_a(\emptyset), s) - \psi_s(m_s). \quad (29)$$

The key feature of the naive updating rule is that it is not a function of the amount of memory manipulation employed by self 1. This greatly simplifies the computation of the PNE of the model since, unlike in the sophisticated case, there is no feedback between self 2's expectation of the manipulation exerted by self 1 and self 1's manipulation choice. Then, the equilibrium amount of manipulation is determined by the maximum of expression (29).

⁶⁰Benabou and Tirole (2006a) present a slightly different formulation of naivety. In their model, agents follow an updating rule that underweights bad news by a fixed proportion. Unlike under our notion of naivety, they still obtain multiple equilibria when agents are naive. However, as in our model, naivety also tends to reduce the welfare costs of self-deception.

Proposition 13 *There exists a PNE. Furthermore, if ψ_σ is strictly convex and $b_a(\hat{\sigma})$ is a (single-valued) function where $\sigma \in \{H, L, \emptyset\}$ and $\hat{\sigma} \in \{H, L, \emptyset\}$, the PNE is essentially unique.⁶¹*

Proof. Existence of PNE follows an argument analogous to the existence of a PBE. Note that Condition 3 from Definition 3 implies that $b_a(\hat{s})$ is not a function of self 1's memory manipulation. Strict convexity of ψ_s implies that expression (29) is strictly concave in m_s . Then, the equilibrium amounts of memory manipulation m_L^* and m_H^* are unique. Condition 4 implies that beliefs must also coincide in all recollections such that $\Pr(\hat{\sigma}_a = \hat{\sigma} | m_L^* = m_H^* = 0) > 0$. ■

Corollary 2 *The PNE is essentially unique when either: (i) $u(a, \cdot, \theta, \sigma_a) : B \rightarrow \mathbb{R}$ is a strictly concave function, or (ii) B is a singleton (i.e., the individual does not take ex-post actions).*

Remark 7 *Suppose that B is finite and fix the natural rates of remembering an outcome η_L and η_H . Since the set of utility functions $u : \Theta \times A \times B \times \mathbb{R} \rightarrow \mathbb{R}$ such that $\arg \max_{b \in B} \{E_\mu[u(a, b, \theta, \sigma_a) | \hat{\sigma}_a = \hat{\sigma}]\}$ contains more than one element is nowhere dense, the PNE is essentially unique for generic utility functions when the set of ex-post actions B is finite.*

The generic uniqueness of the PNE contrasts with the multiplicity of the PBE. Multiplicity arises from the fact that self 1 affects self 2's equilibrium inference when the individual is sophisticated. In the naive case, because there is no effect from memory manipulation on self 2's inference, uniqueness is obtained.

B2 Ambiguity-Seeking Behavior

For simplicity, consider the forgetfulness memory system of Example 1 and, as in the Section A, assume that the utility is additively separable between attributes and money. Then, the equilibrium amount of memory manipulation is $m_L^* = \min \left\{ \psi_L'^{-1}(\Delta u) ; 1 \right\}$. The ex-ante expected utility of the NDM is

$$\begin{aligned} U(\Sigma) &= (1 - q)(1 + m_L^*)u_L + [q - (1 - q)m_L^*]u_\emptyset - (1 - q)\psi_L(m_L^*) \\ &= (1 - q)(1 + m_L^*)u_L + [q - (1 - q)m_L^*]u_H - (1 - q)\psi_L(m_L^*), \end{aligned}$$

where the second inequality uses the fact that $u_\emptyset = u_H$. The NDM prefers to observe the signal Σ if and only if the expected improvement in self-image $|m_L^*| \Delta u$ is greater than the cost of memory manipulation $\psi_L(m_L^*)$.⁶² Thus, naive individuals may prefer to observe signals even if the objective value of information (which in this case is zero) is lower than the expected costs of manipulation.

Remark 8 *Proceeding as in Appendix A, the NDM's expected utility from the monetary lottery can be represented by*

$$U(\Sigma) = w(q)u_H + [1 - w(q)]u_L,$$

where $w(q) = q - (1 - q) \left[m_L^* + \frac{\psi_L(m_L^*)}{\Delta u} \right]$, $w(0) = 0$, and $w(1) = 1$. Thus, the NDM is ambiguity averse if $|m_L^*| \Delta u < \psi_L(m_L^*)$ and ambiguity seeking if the reverse inequality is satisfied. Hence, a naive individual may be ambiguity seeking even when the utility function is additively separable between attributes and money.

B3 Zeroth-Order Risk Seeking Behavior

This subsection shows that the NDM may be zeroth-order risk seeking. As in Subsection A.3, consider a lottery that pays $x = \varepsilon s$, $s \in \{H, L\}$, where $qH + (1 - q)L = 0$. Let $m_s^*(\varepsilon)$ denote the equilibrium amount of memory manipulation as a function of ε . The certainty equivalent of this lottery is:

$$\begin{aligned} \int u(\theta, CE(\varepsilon)) dF(\theta) &= (1 - q)(1 + m_L^*(\varepsilon))u_L(L) \\ &\quad + [q - (1 - q)m_L^*(\varepsilon)]u_H(H) - (1 - q)\psi_L(m_L^*(\varepsilon)) - q\psi_H(m_H^*(\varepsilon)), \end{aligned}$$

⁶¹The PNE is essentially unique in the sense that, all equilibria feature the same choices of actions a and b , manipulation efforts m_L and m_H , and beliefs given recollections that are believed to be reached with positive probability (i.e., $(\hat{\sigma}_a = \hat{\sigma} | m_L = m_H = 0) > 0$). Equilibria may diverge only with respect to beliefs at recollections that are not believed to be reached with positive probability. Obviously, one can ensure uniqueness of beliefs in all recollections by assuming that the NDM believes that all recollections are reached with positive probability: $0 < \min \{\eta_H, \eta_L\} < 1$.

⁶²As in Subsection 4.1, the NDM's surplus from observing a signal is decreasing in the favorableness of her prior distribution over her attributes under Assumption 3. However, unlike Conjecture 1, this surplus may be positive when the individual is naive.

Recall that $u_\sigma(s) \equiv \int u(\theta, s) dF(\theta|\sigma)$ $u_s(\sigma) = v_\sigma + \tau(s)$, where the last equality follows from additive separability. Then, taking the limit as $\varepsilon \rightarrow 0_+$, we obtain:

$$\tau(CE(0)) = -m_L^*(1-q)\Delta v - (1-q)\psi_L(m_L^*) - q\psi_H(m_H^*).$$

Hence, $CE(0) > 0$ if $|m_L^*(0)|(v_H - v_L) > \psi_L(m_L^*(0)) + \frac{q}{1-q}\psi_H(m_L^*(0))$ and the NDM is zeroth-order risk seeking. In the opposite case, the NDM is zeroth-order risk averse. Thus, we have established the following result:

Proposition 14 *The NDM is:*

- zeroth-order risk averse if $|m_L^*(0)|(v_H - v_L) < \psi_L(m_L^*(0)) + \frac{q}{1-q}\psi_H(m_L^*(0))$, and
- zeroth-order risk seeking if $|m_L^*(0)|(v_H - v_L) > \psi_L(m_L^*(0)) + \frac{q}{1-q}\psi_H(m_L^*(0))$.

B4 Convergence

This subsection shows that, under certain regularity conditions, the beliefs of a NDM converge to the truth despite her departure from Bayes' rule. Let $\tilde{\theta}_n$ denote the NDM's expected type.

Proposition 15 *Let $N \rightarrow \infty$ and fix a NPE. There exists a random variable z such that $\tilde{\theta}_n \rightarrow z$ for almost all histories.*

Proof. Let \tilde{q}_n denote the expected value of q according to the naive updating rule. It is straightforward to verify that $1 - \tilde{q}_n$ follows a supermartingale and, by Doob's theorem, it converges in distribution. Then using the fact that $q(\theta)$ is continuous and strictly increasing concludes the proof. ■

The proposition above does not imply that manipulation converges to zero since z may be nondegenerate and, therefore, Δu may not converge to zero. In order to ensure that this is the case, we will use the following regularity condition:

Regularity Condition. The prior distribution $f(q)$ satisfies

$$\lim_{n \rightarrow \infty} \left\{ \begin{array}{l} \frac{\int_0^x q^{\alpha n+1}(1-q)^{\beta n} [q(1-\eta_H) + (1-q)(1-\eta_L)]^{(1-\alpha-\beta)n} f(q) dq}{\int_0^1 q^{\alpha n+1}(1-q)^{\beta n} [q(1-\eta_H) + (1-q)(1-\eta_L)]^{(1-\alpha-\beta)n} f(q) dq} \\ - \frac{\int_0^x q^{\alpha n}(1-q)^{\beta n+1} [q(1-\eta_H) + (1-q)(1-\eta_L)]^{(1-\alpha-\beta)n} f(q) dq}{\int_0^1 q^{\alpha n}(1-q)^{\beta n+1} [q(1-\eta_H) + (1-q)(1-\eta_L)]^{(1-\alpha-\beta)n} f(q) dq} \end{array} \right\} = 0,$$

for all $(\alpha, \beta, x) \in [0, 1]^3$.

The regularity condition states that the impact of a single observation is small when the number of observations is large. It is satisfied, for example, in the forgetfulness model of Example 1 and in the limited memory model of Example 2 when the prior distribution is uniform. The following proposition shows that although the NDM is optimistic in any finite period, her beliefs converge to the truth and her memory manipulation converges to zero.

Proposition 16 *Suppose the regularity condition is satisfied and let $N \rightarrow \infty$. Then, $m_\sigma \xrightarrow{D} 0$ and $\hat{\theta}_n \xrightarrow{D} \theta$.*

Proof. Let $\hat{q}(h^n)$ denote the NDM's beliefs about q given h_n . The regularity condition implies that $(\hat{q}_t|h^n, H) - (\hat{q}_t|h^n, L) \xrightarrow{D} 0$. Let \tilde{E} be the expectation under the NDM's beliefs. Then, $\tilde{E}[q|h^n, H] - \tilde{E}[q|h^n, L] \rightarrow 0$. Then, the amount of memory manipulation converges to zero. But $m_s \xrightarrow{D} 0$ implies that $\hat{\theta}_n \xrightarrow{D} \theta$. ■

C Non-Separable Preferences

In Section A, the DM's preferences additively separable between attributes and money. In this appendix, I consider general utility functions. A key ingredient of the general model is the degree of complementarity between attributes and money. Since a DM is not as affected by monetary outcomes when she is uninformed about her attributes when attributes and money are complementary, complementarity can be interpreted as providing "psychological insurance." Therefore, the DM may prefer a lottery whose outcomes are informative about her attributes if the complementarity effect is greater than the costs of self-deception. Moreover, the resulting probability weighting function may have an "inverted S-shape" as in Tversky and Kahneman (1992) and Prelec (1998).

Let $u_\sigma(s) \equiv \int u(\theta, s) dF(\theta|\sigma)$ denote the expected utility from a monetary amount equal to s conditional on outcome σ . From Lemma 2, $m_H^* > 0 \geq m_L^*$. Define the degree of complementarity between θ and money by

$$\chi(H, L) \equiv u_H(H) + u_L(L) - u_L(H) - u_H(L). \quad (30)$$

Note that $\chi(H, L) \geq 0$ if u has increasing differences and $\chi(H, L) \leq 0$ if u has decreasing differences. The additively separable case presented in the text corresponds to the case where $\chi(H, L) = 0$. The ex-ante expected utility from the lottery is

$$\begin{aligned} U(\Sigma) &= q(\eta_H + m_H^*)u_H(H) + (1-q)(\eta_L + m_L^*)u_L(L) \\ &\quad + q(1-\eta_H - m_H^*)u_\emptyset(H) + (1-q)(1-\eta_L - m_L^*)u_\emptyset(L) - MC, \end{aligned}$$

where $MC = q\psi_H(m_H^*) + (1-q)\psi_H(m_L^*)$ is the expected memory cost. Then, long but tedious algebraic manipulations yield

$$U(\Sigma) = qu_H(H) + (1-q)u_L(L) + z\chi(H, L) - MC. \quad (31)$$

where $z = \frac{q(1-q)(1-\eta_L - m_L^*)(1-\eta_H - m_H^*)}{q(1-\eta_H - m_H^*) + (1-q)(1-\eta_L - m_L^*)} > 0$ and $MC = q\psi_H(m_H^*) + (1-q)\psi_H(m_L^*)$.

The utility of a monetary lottery can be decomposed in three terms: First, the expected utility $qu_H(H) + (1-q)u_L(L)$ of the lottery when memory is perfect. Second, the expected manipulation costs MC . These two effects are precisely the same as in the additively separable case (see equation 11). The third effect, which is not present when the utility is additively separable, is the degree of complementarity between attributes and money. When signals are forgotten, there is probability $\alpha(m_H^*, m_L^*)$ that a high signal was observed and the complementary probability that a low signal was observed. Thus, forgetting a signal can be seen as providing “psychological insurance” to the agent. This raises her expected utility if θ and money are complementary ($\chi > 0$) and decreases her expected utility if they are substitutes ($\chi < 0$).

Proceeding as in Subsection A.1, the DM’s expected utility can be represented by

$$U(\Sigma) = w(q) \times u_H(H) + [1 - w(q)] \times u_L(L),$$

where $w(q) = q + \frac{z\chi(H, L) - MC}{u_H(H) - u_L(L)}$. Moreover, it is straightforward to show that $w(0) = 0$, and $w(1) = 1$. Therefore, when attributes and money are complementary, the DM may exhibit ambiguity loving behavior. In particular, the following example shows that the model may lead to an inverted S-shaped probability weighting function:

Example 7 (Inverted S-shaped Probability Weights) Consider the limited memory model of Example 2 and suppose that the manipulation effort is a binary variable: $m_H \in \{0, \frac{3}{4}\}$, where $\psi_H(\frac{3}{4}) = \frac{1}{5}$. Let $\chi(H, L) = u_H(H) - u_H(L) = 1$. Then, self 1 chooses to engage in memory manipulation if $q \in (0, \frac{11}{12})$. It is straightforward to show that, for values of q such that the DM engages in memory manipulation, the probability weighting function has an inverted S-shape:

$$w(q) \begin{cases} > q & \text{if } q \in (0, \frac{1}{2}) \\ < q & \text{if } q \in (\frac{1}{2}, \frac{11}{12}) \end{cases}.$$

As in Section A, denote by U^I the utility of a lottery with the same distribution over monetary outcomes as the one above but whose monetary outcomes are uninformative about θ . Rearranging equation (31), we obtain

$$U(\Sigma) = U^I + y\chi(H, L) - MC, \quad (32)$$

where $y = q(1-q) \left[1 + \frac{(1-\eta_L - m_L^*)(1-\eta_H - m_H^*)}{q(1-\eta_H - m_H^*) + (1-q)(1-\eta_L - m_L^*)} \right] > 0$. Consider the choice between the lottery Σ and another lottery with the same distribution over monetary outcomes but whose monetary outcomes are uninformative about θ . Equation (32) implies that the DM will prefer lottery Σ if the degree of complementarity is high enough or if the expected memory cost is low enough: $y\chi(H, L) \geq MC$. Therefore, when attributes and money are complementary, the DM may prefer the attribute-dependent lottery.

However, when the monetary lottery is “small” (i.e., when the lottery pays $x = \varepsilon s$ for ε low), the complementarity effect vanishes. Since the memory cost converges to a strictly positive number as ε converges to zero, it follows that the certainty equivalent of the lottery converges to $CE(0) < 0$. Therefore,

$$\lim_{\varepsilon \rightarrow 0^+} \frac{\pi(\varepsilon)}{\varepsilon} = - \lim_{\varepsilon \rightarrow 0^+} \frac{CE(\varepsilon)}{\varepsilon} = +\infty$$

and, for any degree of complementarity between attributes and money, the DM always exhibits zeroth-order risk aversion. This is formally stated in the following proposition:

Proposition 17 *In any PBE, the DM exhibits zeroth-order risk aversion.*

Proof. For any PBE, define the expected manipulation cost as $MC(\varepsilon) = q\psi_H(m_H^*(\varepsilon)) + (1-q)\psi_H(m_L^*(\varepsilon))$. Note that $\lim_{\varepsilon \rightarrow 0^+} \chi(\varepsilon H, \varepsilon L) = 0$. Therefore, for small ε , equation (15) becomes:

$$\int u(\theta, CE(\varepsilon)) dF(\theta) = qu_H(\varepsilon H) + (1-q)u_L(\varepsilon L) - MC(\varepsilon). \quad (33)$$

Since $MC(0) > 0$ and, by the Theorem of the Maximum, $MC(\varepsilon)$ is continuous, it follows that $MC(\varepsilon) > 0$ for small ε . Hence, $\lim_{\varepsilon \rightarrow 0^+} MC(\varepsilon) > 0$. Then, equation (33) yields:

$$\lim_{\varepsilon \rightarrow 0^+} \int u(\theta, CE(\varepsilon)) dF(\theta) > qu_H(0) + (1-q)u_L(0) = \int u(\theta, 0) dF(\theta),$$

where the last equality follows from Bayes' rule. Since u is continuous and increasing in money, this implies that $\lim_{\varepsilon \rightarrow 0^+} CE(\varepsilon) > 0$. Hence, $\lim_{\varepsilon \rightarrow 0^+} \pi(\varepsilon)/\varepsilon = -\lim_{\varepsilon \rightarrow 0^+} CE(\varepsilon)/\varepsilon < 0$. ■

It is interesting to contrast the general model with the a model from the following example where the DM does not face memory costs:

Example 8 (Exogenous Memory Model) *Take $\psi_s(m_s) = +\infty$ for all $m_s \neq 0$. Thus, the agent cannot engage in endogenous memory manipulation. Let $\eta_s < 1$ so that the agent forgets outcome with (exogenous) probabilities $1 - \eta_s > 0$. If $\eta_H > \eta_L$, memory is selective in the sense that good news is more likely to be remembered than bad news.*

When memory manipulation is endogenous (and differentiable at $m_s = 0$), the effect from memory manipulation always dominates the complementarity effects and the DM displays zeroth-order risk aversion. When memory manipulation is exogenous, the order of risk aversion is determined by the degree of complementarity between attributes and money. Note that for small ε , attributes and money are complementary if $u'_H(0) < u'_L(0)$ and substitutes if $u'_H(0) > u'_L(0)$.

Proposition 18 *In the exogenous memory model: (i) the DM is first-order risk averse if $u'_L(0) > u'_H(0)$; (ii) the DM is first-order risk seeking if $u'_L(0) < u'_H(0)$; and (iii) the DM has second-order risk preferences if $u'_L(0) = u'_H(0)$.*

Proof. Since $MC(\varepsilon) = 0$ for all ε , equation (15) becomes

$$\int u(\theta, CE(0)) dF(\theta) = qu_H(\varepsilon H) + (1-q)u_L(\varepsilon L) + z\chi(\varepsilon H, \varepsilon L). \quad (34)$$

Substituting $\chi(0, 0) = 0$, yields

$$\int u(\theta, CE(0)) dF(\theta) = qu_H(0) + (1-q)u_L(0).$$

Therefore, Bayes' rule implies that $\int u(\theta, CE(0)) dF(\theta) = \int u(\theta, 0) dF(\theta)$ and, because u is strictly increasing in money, $\pi(0) = -CE(0) = 0$. Differentiating equation (34), it follows that

$$CE'(0) = \frac{qu'_H(0)H + (1-q)u'_L(0)L + z(H-L)[u'_H(0) - u'_L(0)]}{qu'_H(CE) + (1-q)u'_L(CE)}.$$

Substituting $qH + (1-q)L = 0$, we obtain

$$CE'(0) = K[u'_H(0) - u'_L(0)],$$

where $K = \frac{H}{qu'_H(0) + (1-q)u'_L(0)} \left(q + \frac{z}{1-q} \right) > 0$. Then, by L'Hospital's rule, it follows that

$$\lim_{\varepsilon \rightarrow 0^+} \pi(\varepsilon)/\varepsilon = -CE'(0) = -K[u'_H(0) - u'_L(0)],$$

which concludes the proof. ■

Therefore, the DM may display risk preferences of first order when there are no manipulation costs. Unlike when memory manipulation is endogenous, the DM may be first-order risk seeking or have risk preferences of second order.

D Proof of Claim 2

The p.d.f. of q conditional on h^n is

$$f(q|h^n) = \frac{\prod_{t:\sigma_t=\emptyset} [q(1-\eta_H - m_{H,t}^*) + (1-q)(1-\eta_L - m_{L,t}^*)] \times \prod_{t:\sigma_t=H} q(\eta_H + m_{H,t}^*) \times \prod_{t:\sigma_t=L} (1-q)(\eta_L + m_{L,t}^*) \times f(q)}{\int \left\{ \prod_{t:\sigma_t=\emptyset} [q(1-\eta_H - m_{H,t}^*) + (1-q)(1-\eta_L - m_{L,t}^*)] \times \prod_{t:\sigma_t=H} q(\eta_H + m_{H,t}^*) \times \prod_{t:\sigma_t=L} (1-q)(\eta_L + m_{L,t}^*) \times f(q) \right\} dq}.$$

Let $\#H$ denote the number of times that a signal $\hat{\sigma} = H$ was recollected: $\#\{t : \hat{\sigma}_t = H\}$. Similarly, define $\#L$ as $\#\{t : \hat{\sigma}_t = L\}$.⁶³ Then, after some algebraic manipulations, we can write:

$$f(q|h^n) = \frac{\prod_{t:\sigma_t=\emptyset} [q(1-\eta_H - m_{H,t}^*) + (1-q)(1-\eta_L - m_{L,t}^*)] \times q^{\#H} \times (1-q)^{\#L} \times f(q)}{\int \prod_{t:\sigma_t=\emptyset} [q(1-\eta_H - m_{H,t}^*) + (1-q)(1-\eta_L - m_{L,t}^*)] \times q^{\#H} \times (1-q)^{\#L} \times f(q) dq}.$$

Note that $f(q|h^n)$ is not a function of $m_{H,t}$ and $m_{L,t}$ for any history h^n such that $\hat{\sigma}_t \neq \emptyset$. This follows from the signals σ_t being i.i.d. and the fact that $\hat{\sigma}_t = \sigma_t$ when $\hat{\sigma}_t \neq \emptyset$. Integrating the equation above, we obtain

$$F(x|h^n) = \frac{\int_0^x \prod_{t:\sigma_t=\emptyset} [q(1-\eta_H - m_{H,t}^*) + (1-q)(1-\eta_L - m_{L,t}^*)] \times q^{\#H} \times (1-q)^{\#L} \times f(q) dq}{\int_0^1 \prod_{t:\sigma_t=\emptyset} [q(1-\eta_H - m_{H,t}^*) + (1-q)(1-\eta_L - m_{L,t}^*)] \times q^{\#H} \times (1-q)^{\#L} \times f(q) dq}.$$

We are now ready to prove the Claim:

Proof of Claim 2. We have to show that

$$\frac{\int_0^x \prod_{t:\sigma_t=\emptyset} [q(1-\eta_H - m_{H,t}^*) + (1-q)(1-\eta_L - m_{L,t}^*)] \times q^{\#H} \times (1-q)^{\#L} \times f(q) dq}{\int_0^1 \prod_{t:\sigma_t=\emptyset} [q(1-\eta_H - m_{H,t}^*) + (1-q)(1-\eta_L - m_{L,t}^*)] \times q^{\#H} \times (1-q)^{\#L} \times f(q) dq} \leq \frac{\int_0^x \prod_{t:\sigma_t=\emptyset} [q(1-\eta_H - m_{H,t}^*) + (1-q)(1-\eta_L - m_{L,t}^*)] \times q^{\#H-1} \times (1-q)^{\#L+1} \times f(q) dq}{\int_0^1 \prod_{t:\sigma_t=\emptyset} [q(1-\eta_H - m_{H,t}^*) + (1-q)(1-\eta_L - m_{L,t}^*)] \times q^{\#H-1} \times (1-q)^{\#L+1} \times f(q) dq}.$$

When $x = 0$, both sides become 0 and, when $x = 1$, both sides are equal to 1.

The derivative of the LHS with respect to x is

$$\frac{\prod_{t:\sigma_t=\emptyset} [x(1-\eta_H - m_{H,t}^*) + (1-x)(1-\eta_L - m_{L,t}^*)] \times x^{\#H} \times (1-x)^{\#L} \times f(x)}{\int_0^1 \prod_{t:\sigma_t=\emptyset} [q(1-\eta_H - m_{H,t}^*) + (1-q)(1-\eta_L - m_{L,t}^*)] \times q^{\#H} \times (1-q)^{\#L} \times f(q) dq},$$

and the derivative of the RHS with respect to x is

$$\frac{\prod_{t:\sigma_t=\emptyset} [x(1-\eta_H - m_{H,t}^*) + (1-x)(1-\eta_L - m_{L,t}^*)] \times x^{\#H-1} \times (1-x)^{\#L+1} \times f(x)}{\int_0^1 \prod_{t:\sigma_t=\emptyset} [q(1-\eta_H - m_{H,t}^*) + (1-q)(1-\eta_L - m_{L,t}^*)] \times q^{\#H-1} \times (1-q)^{\#L+1} \times f(q) dq}.$$

⁶³Obviously, $\#H$ and $\#L$ are functions of histories. We omit this dependence for notational clarity.

Note that $\frac{dRHS}{dq} > \frac{dLHS}{dq}$ if and only if

$$\begin{aligned} & \frac{\prod_{t:\sigma_t=\emptyset} [x(1-\eta_H - m_{H,t}^*) + (1-x)(1-\eta_L - m_{L,t}^*)] \times x^{\#H} \times (1-x)^{\#L} \times f(x)}{\int_0^1 \prod_{t:\sigma_t=\emptyset} [q(1-\eta_H - m_{H,t}^*) + (1-q)(1-\eta_L - m_{L,t}^*)] \times q^{\#H} \times (1-q)^{\#L} \times f(q) dq} \\ & \leq \frac{\prod_{t:\sigma_t=\emptyset} [x(1-\eta_H - m_{H,t}^*) + (1-x)(1-\eta_L - m_{L,t}^*)] \times x^{\#H-1} \times (1-x)^{\#L+1} \times f(x)}{\int_0^1 \prod_{t:\sigma_t=\emptyset} [q(1-\eta_H - m_{H,t}^*) + (1-q)(1-\eta_L - m_{L,t}^*)] \times q^{\#H-1} \times (1-q)^{\#L+1} \times f(q) dq}. \end{aligned}$$

Rearranging, we obtain:

$$\begin{aligned} & \frac{x}{\int_0^1 \prod_{t:\sigma_t=\emptyset} [q(1-\eta_H - m_{H,t}^*) + (1-q)(1-\eta_L - m_{L,t}^*)] \times q^{\#H} \times (1-q)^{\#L} \times f(q) dq} \\ & \leq \frac{(1-x)}{\int_0^1 \prod_{t:\sigma_t=\emptyset} [q(1-\eta_H - m_{H,t}^*) + (1-q)(1-\eta_L - m_{L,t}^*)] \times q^{\#H-1} \times (1-q)^{\#L+1} \times f(q) dq}. \end{aligned}$$

Thus, $\frac{dRHS}{dq} > \frac{dLHS}{dq}$ if and only if

$$\rho(x) > \frac{\int_0^1 \prod_{t:\sigma_t=\emptyset} [q(1-\eta_H - m_{H,t}^*) + (1-q)(1-\eta_L - m_{L,t}^*)] \times q^{\#H} \times (1-q)^{\#L} \times f(q) dq}{\int_0^1 \prod_{t:\sigma_t=\emptyset} [q(1-\eta_H - m_{H,t}^*) + (1-q)(1-\eta_L - m_{L,t}^*)] \times q^{\#H-1} \times (1-q)^{\#L+1} \times f(q) dq},$$

where $\rho(x) = \frac{x}{1-x}$. Since $\rho(0) = 0$, $\rho(1) = +\infty$, $\rho(x)$ is strictly increasing in x , and the term on the right is a positive constant, there exists a unique \bar{x} such that

$$\rho(x) > (<) \frac{\int_0^1 \prod_{t:\sigma_t=\emptyset} [q(1-\eta_H - m_{H,t}^*) + (1-q)(1-\eta_L - m_{L,t}^*)] \times q^{\#H} \times (1-q)^{\#L} \times f(q) dq}{\int_0^1 \prod_{t:\sigma_t=\emptyset} [q(1-\eta_H - m_{H,t}^*) + (1-q)(1-\eta_L - m_{L,t}^*)] \times q^{\#H-1} \times (1-q)^{\#L+1} \times f(q) dq},$$

if $x < (>) \bar{x}$.

Therefore, we have that $\frac{dRHS}{dq} > \frac{dLHS}{dq}$ if $x < \bar{x}$ and $\frac{dRHS}{dq} < \frac{dLHS}{dq}$ if $x > \bar{x}$. Thus, the inequality is satisfied for all q (it is satisfied with strict inequality whenever $q \in (0, 1)$ and with equality at $q \in \{0, 1\}$). ■

References

Prelec, D. (1998). "The probability weighting function," *Econometrica*, 66, 497–527.

Tversky, A. and D. Kahneman (1992). "Advances in prospect theory: cumulative representation of uncertainty," *Journal of Risk and Uncertainty*, 5, 297–323.