A RETRIEVED-CONTEXT THEORY OF FINANCIAL DECISIONS

Jessica A. Wachter
Michael Jacob Kahana

A Retrieved-Context Theory Of Financial Decisions
Jessica A. Wachter and Michael Jacob Kahana
NBER Working Paper No. 26200
August 2019, Revised December 2022
JEL No. D91,E71,G11,G12,G41

## <u>ABSTRACT</u>

Studies of human memory indicate that features of an event evoke memories of prior associated contextual states, which in turn become associated with the current event's features. This mechanism allows the remote past to influence the present, even as agents gradually update their beliefs about their environment. We apply a version of retrieved context theory, drawn from the literature on human memory, to explain three types of evidence in the financial economics literature: the role of early life experience in shaping investment choices, occurrence of financial crises, and the impact of fear on asset allocation. These applications suggest a recasting of neoclassical rational expectations in terms of beliefs as governed by principles of human memory.

Jessica A. Wachter
Department of Finance
2300 SH-DH
The Wharton School
University of Pennsylvania
3620 Locust Walk
Philadelphia, PA 19104
and NBER
jwachter@wharton.upenn.edu

Michael Jacob Kahana
University of Pennsylvania
425 South University Avenue
Philadelphia, PA 19104
kahana@psych.upenn.edu

# 1 Introduction

Imagine returning to a childhood home after many years. As you ride along, familiar scenery drifts into view. Your thoughts, formerly so occupied with the latest work deadlines and obligations, or pressing matters of health and family, drift back to the last time you saw these sights. Pretty soon this drift is an unstoppable force, taking you by surprise. You simply cannot stop thinking about your childhood, and about the last time you saw these sights. You are surprised by the strength of the memories that have come back. Where have they been all this time? It is as if there were a different person hiding inside you.

What does this have to do with economics? The standard economic model of decision-making seems to leave little role for this phenomenon. A typical modeling device is to consider states of the world and probabilities assigned to each of them; the agent then maximizes an expectation of a utility function Savage (1954). The probabilities come from data, perhaps, and the utility function from an unknown source, but in any case the process seems far removed from the description above.

The subject of this paper is the role of encoding and retrieval of information in decision-making. Linking the encoding and the retrieval event is the idea of context. In the above example of a trip back to a childhood home, context is drawn from the physical environment, and drives memory retrieval, in an effect known in the psychology literature by the evocative term, "jump-back-in-time." In important prior work, Bordalo et al. (2020b) use a model of context to explain the formation of norms that drive decisions. Theirs is the first paper in economics to make an explicit link between the important idea of contextually-driven retrieval and economic decisions.

In the model of Bordalo et al. (2020b), context is static and is inherited directly from the environment. Consider, for example, the decision to purchase a coat from a catalogue. The price of the coat is in the catalogue, but what is the quality, namely the utility it will bring to the agent? If it is cold outside, the agent will bring to mind other instances of coat-wearing when it is cold, and, as in the empirical literature, is

2

more likely to order a sweater (Conlin et al., 2007). Bordalo et al. (2020b) consider many such examples, and link the idea to such notions as the tendency to notice price changes that may explain the affect of monetary policy on the real economy.

Bordalo et al. (2020b) introduce contextual retrieval into economic modeling, and show how it cannot be ignored in models of decision making. And yet, perhaps, there is more to the story. Consider, for example, the "jump-back-in-time" effect, which is not only well-known to most of us, but is well-studied and comes with its own distinct neural signature (Howard et al., 2012). The jump-back-in-time does not fit easily into a static framework. For a "jump" requires a before and an after, and therefore a stochastic process. Moreover, what happens is more than "remember better" – it is an entire set of memories that simply might not have occurred to one before. It is dynamic, and yet somehow beyond the reach of rational expectations as usually defined. Key to this notion is a persistent mental state that reacts to the outside world, and yet stands outside it – a context "layer" (or code, or representation) that differs from the features representing the direct experience.[1] This persistent mental state is simple and convincing. Most people "know" such a state exists. It is a piece missing from static context models.

Examining empirical findings in financial economics, one is immediately confronted with evidence that could hardly be explained *without* a model of contextually driven recall. For example, Cuculiza et al. (2020) show that analyst earnings forecasts become more negative upon anniversaries of terrorist attacks (as well as following the attacks themselves). Ramadorai et al. (2020) show that investors who were randomized into an IPO in which they received shares in companies that subsequently performed well not only purchase similar stocks, but trade more in general. Bassi et al. (2013) show that sunshine and good weather promote risk-taking behavior. The literature has identified both a January effect (Keim, 1983) and a Monday effect (Keim and Stambaugh, 1984),

---

[1]In contrast, the model of Bordalo et al. (2020b), is based in psychological studies such as Godden and Baddeley (1975) where a static context context is embedded in the environment. In this type of model the feature layer of the environment and the context layer are the same.

January and Monday both having an important role in context reinstatement. And, despite there being no apparent news, the Monday following Black Monday, October 27th, 1987 featured a one-day negative return of -8% and sharply elevated implied volatility.

Are these one-off findings, or do they represent the tip of a context-driven iceberg? Two papers making the explicit connection are Cohn et al. (2015) and Guiso et al. (2018). The former examines risk-taking by professional investors. When cued with a stock market decline, these investors are less likely to choose a risky gamble. Guiso et al. (2018) show that following the financial crisis, investors were also less likely to take a gamble. To draw a direct connection to risk aversion, they showed that when college students were shown scenes from a horror movie, they were also less likely to take a gamble. These studies show that a recent "experience" – even a simulated one – affects decision-making, raising the question of how many decisions that we make, including the ones that drive asset prices, are driven by experience.

Below the tip of the iceberg is the vast and growing field evidence of the important role of experience in shaping life decisions. Malmendier and Nagel (2011) showed that early life experience affected stock market participation; besides rigorously showing a role for experience effects in the important decision about whether to participate in the stock market, Malmendier and Nagel (2011) provide suggestive evidence for their experience-driven variable explaining fluctuations in asset valuation over the last half-century; a finding that Nagel and Xu (2018) make more precise. The literature has shows similar effects in domains ranging from everyday spending decisions to the setting of interest rate policy (Malmendier and Nagel, 2011, 2016; Malmendier et al., 2017; Malmendier and Shen, 2018). These models are very different from the above in that there is no cue; yet in the background is very much a role for memory.

While not obvious, there is a strong connection between these early-life experience-based accounts and contextually-driven retrieval. There are two parts to the memory: one is encoding, the other retrieval. Key to these descriptions is that early experience matters more than it "should," given the vast array of data were are confronted with.

4

What memory models allow us to understand is that previous experience – often early in life – shapes how we interpret the data. In some cases, this effect is so extreme that we cannot ignore it, as in post-traumatic stress disorder (Rubin and Berntsen, 2009). In the context-based model in this paper, early experience matters because it is encoded, and, when it is recalled, it is encoded again. Agents' thoughts become the very data on which they rely for future decision-making.

In this paper, we propose a memory-based model of decision-making under uncertainty. A wealth of data support the idea of a human memory system that maintains a record of associations between experiential features of the environment, and underlying contextual states (Kahana, 2012). This record of associations, together with inference about the current contextual state, constitutes a belief system that could potentially affect any kind of choice under uncertainty. This belief system responds to the current environment through retrieved context. The mechanism of retrieved context is how memory "knows" what information is most relevant to bring forward to our attention at any given time, and for this reason offers the potential for an account of optimality of memory (Azeredo da Silveira and Woodford, 2019). At the same time, any new experience, and the context itself, is then stored again in the memory system (Howard and Kahana, 2002). We apply this model to three illustrative examples: early-life decision-making, asset pricing around the financial crisis, and specifically the experiment of Guiso et al. (2018).

Whereas many applications of psychological principles to economic decision making have focused on cognitive biases such as loss aversion and narrow framing (see Barberis (2013)), or on limited attention (see Gabaix (2019)) the literature on human learning and memory offers a different perspective. Three major laws govern the human memory system: similarity, contiguity, and recency. Similarity refers to the priority accorded to information that is similar to the presently active features, contiguity refers to the priority given to features that share a history of co-occurrence with the presently active features, and recency refers to priority given to recently experienced features. All three "laws" exhibit universality across agents, feature types, and memory tasks and thus

provide a strong basis for a theory of economic decision-making.

While few economic models explicitly incorporate these laws, there are exceptions. Mullainathan (2002) proposes a model in which agents tend to remember those past events which resemble current events, and in that way incorporates associative memory. Perhaps most relevant to the present model is the fact that Mullainathan directly incorporates the fact that memory for a just-remembered item is greater than otherwise. This effect, sometimes termed *rehearsal*, is incorporated in our model through the idea that thoughts become data. Another means of modeling associative memory is that of Gilboa and Schmeidler (1995), who replace axiomatic expected utility with utility computed using probabilities that incorporate the similarity of the current situation to past situations. Other models directly model the principle of recency using a constant-gain model (Nagel and Xu, 2018), extrapolative expectations (Barberis et al., 2015), or other non-Bayesian expectations-formation mechanisms (Fuster et al., 2010).

The remainder of the paper is organized as follows. Section 2 describes the model and derives general properties. Section 3 describes the psychological and neural basis for the model. Section 4 discusses applications to problems in economics and finance. Section 5 briefly describes alternative approaches. Section 6 concludes.

## 2    Integrating Memory into Decision Making

We consider an agent who develops memories by experiencing events, and then makes decisions. We represent one occurrence by $Y_t$ and the time series of these events using $\{Y_t\}$. We assume $Y_t$ can take on one of $n$ values, which we group together in $\mathcal{Y} = \{y_1, \ldots, y_n\}$. We also assume a persistent latent process $Z_t$ taking on values in $\mathcal{Z} = \{z_1, \ldots, z_m\}$. Let $p_{ik}^Z$ denote the probability of transition from state $i$ to state $k$: $p_{ik}^Z \equiv \text{Prob}(Z_{t+1} = z_k \,|\, Z_t = z_i)$. Let $p(y_j|z)$ denote the probability that $Y_t = y_j$ conditional on $Z_t = z$ for $j = 1, \ldots, n$. It is as if nature delivers a set of persistent states about which the agent can partially learn through observation. This is a standard set-up in macroeconomics and finance (Hamilton, 1994; Sims, 2003).

We consider a static decision problem under uncertainty. In each period, the agent makes a choice denoted by $\pi$, perhaps subject to constraints. The agent has a utility function $V$ that depends on $\pi$ and on the outcome of the state of the world next period. We make the standard assumption that only outcomes traced to observables matter to the agent.[2]

**Assumption 1.** *The agent solves*

$$\max_{\pi} \mathbb{E}_t^Y \left[ V(Y_{t+1}, \pi) \right] \tag{1}$$

*where $\mathbb{E}_t^Y$ is the time-t subjective expectation over $Y_{t+1}$.*

The expectation $\mathbb{E}_t^Y$ is the subject of our paper, and will be described more fully in the sections that follow.

## 2.1   Human memory

A standard approach to the problem outlined above is to endow an agent with a system of prior beliefs on the joint dynamics $\{Y_t, Z_t\}$. Based on this prior and on the data, which are assumed to be equally at the agent's fintertips at all points in time, the agent can infer a posterior distribution over unknown quantities of interest such as transition probabilities and latent states. Depending on how restrictive one makes the agent's prior distributions, the problem becomes more or less well-identified, with an unavoidable tradeoff between bias and precision. Thus the agent's inference problem is a difficult one.

The literature on human memory offers an alternative approach to the problem of decision-making under uncertainty. We draw on the vast literature describing the influence of past experience on present behavior, a topic that has occupied the attention of experimental psychologists for more than a century (Ebbinghaus, 1913; Müller

---

[2]Assumption 1 implies that the latent state, and hence context, does not play a direct role in utility. The agent is only influenced by, say, mood through what mood tells the agent about the distribution of future features.

and Pilzecker, 1900; Jost, 1897; Müller and Schumann, 1894; Ladd and Woodworth, 1911; Carr, 1931). Because memories of recent experiences readily come to mind, this early work sought to uncover the factors that lead to forgetting. Experimental findings quickly challenged the folk assertion that memories decay over time, eventually becoming completely erased. Rather, they revealed that removing a source of interference, or reinstating the "context" of original learning, readily restored these seemingly forgotten memories (McGeoch, 1932; Underwood, 1948; Estes, 1955).

The original view of context was that it consisted of latent, background information unrelated to the present stimulus. Experiments often treated context as an aspect of the external environment. In contrast to this early work, Howard and Kahana (2002) proposed a model in which context became a *mental representation* of what had formerly been thought of as the physical environment. Context went from being an external physical concept to the internal state of the agent. In the Howard and Kahana theory, which we term *retrieved context theory*, the set of psychological (or neural) features that represent a stimulus enter into association this with internal mental state. The database of such associations form the basis for performance in recall, recognition, and categorization tasks. Subsequent work, e.g. Polyn et al. (2009) Lohnas et al. (2015), has thus emphasized the view of context as an internal process, evolving endogenously based on the stimuli which the agent encounters.

In a free recall experiment, the researcher presents subjects with a list of of items, often words, which they can recall in any order. We draw an analogy from this list of items to the features that nature presents to the agent (an analogy that is implicit in the notion of the free recall experiment). Howard and Kahana (2002) model these features (words) as basis vectors in a large $n$-dimensional space $f_t$. Their key innovation is the idea of a mental context that links these features through time:

$$x_t = (1 - \zeta)x_{t-1} + \zeta x_t^{\text{in}}, \tag{2}$$

where $x^{\text{in}}$ (the "in" stands for input) is the aspect of context that arises (is retrieved)

8

from the current environment, and where $\zeta$ lies between 0 and 1.[3] Context is an $m$-dimensional vector. As we will see, the context vector contains the decision-maker's subjective probabilities of latent states.

A defining assumption of retrieved context theory is the manner in which the agent retrieves context from the environment. The agent has a prior record of associations between features and internal context, stored in a "memory matrix." Retrieved context arises from multiplying current features $f_t$ with the memory matrix:[4]

$$x_t^{\text{in}} \propto M_{t-1} f_t.$$

where $M_{t-1}$ is the value of the memory matrix as of time $t-1$. It is convenient to scale $x^{\text{in}}$ so its elements sum to one:

$$x_t^{\text{in}} \equiv \frac{M_{t-1} f_t}{||M_{t-1} f_t||}, \tag{3}$$

where, unless stated otherwise, $\|\cdot\|$ denotes the sum of the elements in the vector.[5] After observing $f_t$ and forming $x_t$ based on retrieved context $x_t^{\text{in}}$, the memory matrix updates based on the *outer product* between current context and features:

$$M_t = M_{t-1} + x_t f_t^\top. \tag{4}$$

The current state for the agent is then summarized by $m \times n$ memory matrix $M_t$ and $m \times 1$ vector $x_t$.[6]

---

[3]Formally, features $f_i$ are elements of $\mathcal{B}^n \subset \mathbb{R}^n$, where $\mathcal{B}^n$ is a set of basis vectors that spans $n$-dimensional space. Context is an element of $\mathcal{A}^m \subset \mathbb{R}^m$, for $m \leq n$. That is $\mathcal{A}^m = \{x_t = [x_{1t}, \ldots, x_{mt}]^\top \in \mathbb{R}^m \,|\, \iota^\top x_t = 1\}$, where $\iota$ denotes a conforming vector of ones. We will have use for features that are not basis vectors, in which case they are elements of the unit circle in $n$-dimensional space.

[4]The symbol $\propto$ denotes equality up to multiplication by a positive scalar. Its use implies that we only care about the magnitude of the elements of the vector relative to one another, not in absolute terms.

[5]Because all the vectors we consider have non-negative entries, $\|\cdot\|$ is a valid distance measure; in fact it is distance under the $L^1$-norm. The memory literature, e.g. Polyn et al. (2009), uses the $L^2$-norm, with $x_t = \rho_t x_{t-1} + \zeta x_t^{\text{in}}$ and $\rho_t \approx 1 - \zeta$, to maintain $x_t$ on the unit circle.

[6]To complete the model, the agent must possess initial associations $M_0$. How the agent comes by

To understand the implications of (3) and (4), consider what happens when the agent is cued with features $f_t$. This cuing will recover all the contexts previously associated with those features. Retrieved context equals:

$$
\begin{aligned}
x_t^{\text{in}} &= \frac{M_{t-1}f_t}{||M_{t-1}f_t||} \\
&\propto M_0 f_t + \sum_{s=1}^{t}(x_s f_s^{\top})f_t \\
&\propto M_0 f_t + \sum_{s=1}^{t} x_s (f_s^{\top} f_t).
\end{aligned} \tag{5}
$$

A benchmark case (Assumption 2 below) has $f_s^{\top} f_t$ equal either to zero or one. It is zero if $f_s \neq f_t$; it is one if $f_s = f_t$. Equation 5 shows that features evoke the past contexts under which they are experienced. A context $x_s$ appears in the sum in (5) if the corresponding $f_s$ equals $f_t$; otherwise it does not. The interpretation remains valid even when $f_s$ are not orthonormal basis vectors. Even when features are not orthonormal, (5) is an average of past contexts under which similar features were experienced.[7]

According to retrieved context theory, context determines what an agent is most likely to remember. Figure 1 illustrates the mechanism. The current state of context contains a component that overlaps with the contexts of recent experiences, and a retrieved context component that overlaps with items experienced close in time to the just-recalled item(s). The figure illustrates these two effects as spotlights shining down on memories arrayed on the stage of life. Memories are not truly forgotten, but just obscured when they fall outside of the spotlights. Section 3 discusses the implications of this model for memory, and in particular for the temporal contiguity property in the introduction.

The theory developed by Howard and Kahana (2002), Polyn et al. (2009), and

---

these associations is beyond the reach of this paper.

[7]The theory is silent on initial associations $M_0$. It is convenient to assume that $M_0$ is of full rank so that any observed features will retrieve a non-zero context.

Lohnas et al. (2015) among others treats memory as an outcome of a mechanistic process. There is no decision-maker maximizing an objective function, nor is there an underlying probability space on which such an agent would form beliefs. The first step in mapping their framework to decision-making under uncertainty is to connect features with observable aspects of the environment:

**Assumption 2.** *At each time $t$, an $n$-dimensional vector $f_t$ characterizes the physical environment, with*

$$f_t(j) = \begin{cases} 1 & if \quad Y_t = y_j \\ 0 & otherwise. \end{cases}$$

*That is, $f_t = e_j$, the $j$th standard basis vector in $n$-dimensional space, where $j$ corresponds to the state of $Y_t$.*

The assumption of basis vectors for physical features is analytically convenient and standard in the memory literature. It is less restrictive than it may appear, given that $n$ could be very large.[8] It will at times be useful to consider features arising, for example, in an experiment, and also as part of the agent's recollections. These may not be basis vectors.

The second step in linking memory to decision-making is to relate the contents of memory to subjective probabilities.[9] Context forms the probabilities over latent states. Then, (2) and (5) imply that beliefs about latent states today are a weighted average of such beliefs in the recent past and over the agent's lifetime with the weighting on the latter determined by similarity of features.

**Assumption 3.** *At time $t$, the agent assigns the probability $x_t(i)$ to state $i$, and acts as if this probability is permanent.*

---

[8]A realistic dimensionality $n$ for the features space is on the order of $10^7$, an estimate of the number of neurons in the brain involved in memory storage.

[9]Behind Assumption 3 is the idea that what memory brings to mind is the basis of probabilities that inform financial decisions. That there should be some link between memory and beliefs is confirmed by recent experimental work in finance (Gödker et al., 2019; Enke et al., 2020). The second part of the assumption, namely that beliefs are permanent, pertains to the agent's view of future self. It is necessarily more speculative, but it is plausible given the observed power of context.

To motivate how a model of memory could becomes model of beliefs and ultimately choice, it is helpful to understand what happens when latent states are observed. The key insight is that (**??**) counts co-occurences. If $x_t$ and $f_t$ equal the $i$th basis vector and $j$th basis vector respectively, then $x_t f_t^\top$ represents exactly one occurrence of state $Z_t = z_i$ and $Y_t = y_j$. Elements appear in $M_t$ in proportion to their relative occurrence in the agent's environment.

**Theorem 1.** *Assume context corresponds to the underlying state $Z_t$, and that it is fully observable. Then $M_t$ correctly encodes posterior probabilities of each state.*

Given that $M_t$ "counts" occurences, a direct implication of (**??**) is that retrieved context has an interpretation as a conditional probability:

**Theorem 2.** *Consider a Bayesian agent who assumes $\{Z_t\}$ is iid and who interprets $M_{t-1}$ as containing (known) probabilities of joint occurrences of $\{Z_t\}$ and $\{Y_t\}$. Then retrieved context (3) is the conditional probability of $Z_t$ given $Y_t$.*

These analogies to Bayesian updating are limited because they assume states are observed, and say nothing yet about dynamics. Before considering these aspects of the model, which will lead to departures from Bayesian updating, we turn to one more important basic property:

Now suppose context is latent and evolves endogenously. Equation 4 still implies that the agent stores joint occurrences of context and features. Supposing that $\zeta = 1$, the agent will retain the incorrect associations, regardless of how much data are observed:

**Theorem 3.** *Assuming $\zeta = 1$, $M_t$ evolves such that the relative magnitudes of the elements in each column remain invariant while the magnitudes across columns change. That is, the agent correctly estimates the occurence of features, but incorrectly associates features with context.*

In the special case of $\zeta = 1$,

$$x_{t+1} = x_{t+1}^{\text{in}} = \alpha M_t f_{t+1},$$

12

where $\alpha$ is a positive scalar of proportionality. Consider the formation of the new association:

$$
\begin{aligned}
M_{t+1} &= M_t + x_{t+1} f_{t+1}^\top \\
&= M_t + \alpha M_t \left( f_{t+1} f_{t+1}^\top \right).
\end{aligned}
\tag{6}
$$

Suppose, for example, that $f_{t+1} = e_1$, the first basis vector. Then $f_{t+1} f_{t+1}^\top$ is the matrix with 1 in the first diagonal element and zero otherwise. Equation 6 takes the first column of $M_t$ and multiplies it by $1 + \alpha$, leaving the rest of the matrix unchanged.

Consider the intuition behind Theorem 1: the agent experiences events and stores them in memory. Similar reasoning is at work in Theorem 3: the agent "experiences" joint occurrences of $x_t$ and $f_t$ and stores them in memory (the agent's thoughts become data). However, $x_t$ is not reality; rather it is context that is retrieved based on prior associations from features. The agent nonetheless stores it as an event in memory. The agent correctly learns the frequency of outcomes of $\{Y_t\}$, but matches this frequency incorrectly to the underlying states.

This intuition extends to $\zeta < 1$. Note that columns of $M$ translate directly into which context is retrieved; for this reason, we describe the more general result directly in terms of retrieved context. Rather than being equal at all times, retrieved context decays at a rate that is slower than exponential, and at an ever slower rate as time goes by. The decay takes place in "event time," not in calendar time. For example, consider, as we will in Section 4, the experience of stock market losses. Suppose these take place at a sequence of times $t_1, t_2, \ldots$. For the first loss that the agent experiences, retrieved context is determined by the initial matrix $M_0$. For the second loss,

$$
x_{t_2}^{\text{in}} = \left( 1 - \frac{1 - \zeta}{2} \right) x_{t_1}^{\text{in}} + \frac{1 - \zeta}{2} x_{t_2 - 1},
$$

where $x_{t_2 - 1}$ is the context just before the loss. We generalize this to an arbitrary number of losses in the following theorem:

13

**Theorem 4.** *Consider events occurring at a subsequence of times* $\{t_1, t_2, \ldots, t_\ell, \ldots\}$. *Assumption 2 implies that retrieved context follows the vector process*

$$x_{t_\ell}^{\text{in}} = \left(1 - \frac{1 - \zeta}{\ell}\right) x_{t_{\ell-1}}^{\text{in}} + \frac{1 - \zeta}{\ell} x_{t_\ell-1} \tag{7}$$

*for* $\ell > 1$. *For the first occurrence of the event* $(\ell = 1)$, $x_{t_1}^{\text{in}} \propto M_0 e_i$, *where* $e_i$ *is the basis vector corresponding to the event.*

Theorems 3 and 4 give two key properties of the model: thoughts become data and recency that is independent of scale. Theorem 5 describes another key property: similarity, namely better memory for an item that is similar to the just-remembered item.

**Theorem 5** (Similarity). *Assume that* $\hat{f}_t$ *is geometrically close to* $f_t$ *in a vector space. Then the context retrieved by* $\hat{f}_t$ *is close that that retrieved by* $f_t$.

Similar features evoke similar contexts; but similar features are not always relevant. For example, similarity can arise because of anniversary effects, or through, say witnessing a scene from a movie. These events powerfully bring to mind a context similar to the original. The model shares the basic similarity mechanism with Bordalo et al. (2020b), who model context as an exogenous cue.

So far, we have shown similar features evoke similar contexts (a step necessarily absent in a model with exogenous context). Section 2.3 will show how context then determines what is remembered, and hence what the agent believes.

## 2.2   Implications of autoregressive context

Prior to the full specification of beliefs, we briefly discuss properties implied by autoregressive context. Many of these ideas have antecedents in the behavioral literature, and they will be useful in the applications.

Gennaioli and Shleifer (2018) posit that agents have a tendency to neglect risk. One possible source for this effect is contextual drift.

**Theorem 6** (Neglected risk). *Fix a state $i$ and assume that between time $t$ and $t'$, the agent fails to experience features associated with state $i$. Further assume mutually orthogonality of features experienced between $t$ and $t'$.[10] Then the probability the agent places on state $i$ decays exponentially: $x_{t'}(i) = (1 - \zeta)^{t-t'} x_t(i)$.*

More generally, (2) implies that recently-experienced features are over-represented in the agent's context $x_t$, providing a mechanism for extrapolative beliefs (Barberis et al., 2015). It is, however, also the case that agents' beliefs also can adjust quickly based on new information arising from retrieved context. Beliefs can jump if features are sufficiently novel. The following theorem analyzes when this "jump" effect dominates, and when more traditional recency (neglected risk, extrapolative beliefs) dominates: [11]

**Theorem 7** (Short-run under-reaction; long-term reversal). *Consider a state $i$ that is uniquely associated with a latent state, and assume correct associations at $t$ (see Appendix B). Assume that $t$ is large, or that $\mathcal{Y}_i$ contains many elements. Then*

1. *If features associated with $i$ are novel, and if states are persistent, upward revisions to beliefs about $i$ tend to be followed by further upward revisions.*

2. *If the agent has repeatedly experienced features associated with $i$, then beliefs tend to reverse.*

There is another sense in which the agent over-reacts. If states are sufficiently transitory, the agent confuses a conditional probability with an unconditional one:

---

[10]The assumption of independent features allows us to focus on the exponential decay of context, suppressing the creation of associations in $M$ (explored at length later in the paper). Orthogonal features, are used in the memory laboratory to reset context. These features consist of "distractor tasks," such as asking subjects to solve arithmetic problems, or to view outdoor scenes (Howard and Kahana, 1999; Manning et al., 2016).

[11]For models and for a discussion of the evidence on under-reaction and reversal, see Daniel et al. (1998), Barberis et al. (1998), and Hong and Stein (1999). The price momentum effect is the finding that stocks with the highest price appreciation measured over the past 12 months outperform those with the lowest price appreciation (Jegadeesh and Titman, 1993). These gains partially reverse one year later. Macroeconomic expectations appear to under-react and then overshoot Angeletos et al. (2020), as do earnings expectations Bordalo et al. (2020a).

15

**Theorem 8** (Over-reaction)**.** *When context is sufficiently persistent and states transitory, upward revisions to beliefs about state $i$ tend to reverse.*

It is useful to compare the context updating equation to the physical updating process. Let $\mathbf{E}$ denote the expectation calculated by the econometrician. Revisions in beliefs equal:

$$
\begin{aligned}
\mathbf{E}_t[\Delta x_{t+1}(i)|\Delta x_t(i) > 0] &\approx \zeta(p_{ii}^Z - x_t(i)) \\
&\approx \zeta\left(p_{ii}^Z - ((1-\zeta)x_{t-1}(i) + \zeta)\right),
\end{aligned} \tag{8}
$$

because $x_{t+1}^{\text{in}}(i) = 1$. Consider first the case with $p_{ii}^Z$ high and $x_{t-1}(i)$ low (think of it as zero). This is the setting of the first statement of Theorem 7. Beliefs under-react to news because the persistence of state $i$ is greater than the weight of the new information is context. Slow adjustment of context implies that the agent cannot "take in" all the information at once. Now consider $p_{ii}^Z$ high, $x_{t-1}(i)$ high (think of it as one). This is the setting of Theorem 6 and of the second statement of Theorem 7. If the agent has seen sufficient observations consistent with a state, the agent forgets that other states are possible. Beliefs predictably reverse because of mean-reversion in the state. Finally, consider the case of $p_{ii}^Z$ low. This is the setting of Theorem 8. If the agent puts relatively high weight on new information, or if the state is transitory, then then the agent over-reacts. In the last two cases, the agent forgets that the world may be different next period. Note that our model generates over-reaction from contextual dynamics; this is contrast to Mullainathan (2002) and Bordalo et al. (2020b), in which over-reaction through a mechanism akin to similarity, namely that irrelevant information acts as a cue. Our model also incorporates this latter effect, but over-reaction can occur even when the agent's associations are correct.

While optimality of memory lies outside the scope of this study, one might conjecture that slow context evolution protects against over-reaction. If the state is transitory, then it is indeed inappropriate to react strongly. On the other hand, times might call for an instantaneous shift in perspective, as captured by Theorem 8. Understanding

16

the constraints on storage and retrieval is a prerequisite to a study of optimality of memory; Azeredo da Silveira and Woodford (2019) offer a theory of such constraints. A constrained-optimal view of memory might shed light on when there is sufficient information for an agent to act as Bayesian models specify, and when the agent instead must rely on the mechanisms that we emphasize here.

## 2.3   Features retrieval

The previous analysis still leaves unanswered the question: how do agents go from probabilities on underlying states in $\mathcal{Z}$ (given by context) to outcomes in $\mathcal{Y}$? While nature supplies the true $p(y|z)$, agents must infer this distribution. Consistent with the literature on human memory, we assume that agents draw on their memories in forming this distribution, and that the same associative principles apply. Besides being a necessary step in implementing our approach to decision-making, temporal contiguity (which the next section describes) requires a features retrieval step.

Consider the free recall experimental framework from Section 2.1 used to motivate internal context. In this experimental paradigm, subjects recall words in any order. The step from context to features is known as features retrieval. This features retrieval step (what features a current context calls to mind) will give us the probabilities that underly the decision problem (1).[12]

First define features that are retrieved by a vector putting 100% weight on context $i$: [13]

$$f_{i,t}^{\text{in}} \equiv \frac{M_{t-1}^{\top}\hat{e}_i}{||M_{t-1}^{\top}\hat{e}_i||}. \tag{9}$$

This equals the conditional probability $p(y \,|\, z_i)$. Averaging over these vectors give us the subjective probabilities underlying $\mathbb{E}_t^Y$.[14]

---

[12]While we model features retrieval in a way that supports Bayesian updating as a special case, the memory literature tends to assume a retrieval rule that is closer to winner-take-all ("stronger" features, i.e. those with a higher weight in (10), inhibit the recall of weaker ones), e.g. a "leaky accumulator" model.

[13]The notation $\hat{e}_i$ denotes features in $m$-dimensional context space, to distinguishing them from features basis vectors $e_j$.

[14]Implicit in Assumption 4 is an assumption on timing. The agent enters period $t$ with memory

**Assumption 4.** *Define $f_{i,t}^{\text{in}}$ as (9) and*

$$f_t^{\text{in}} \equiv \sum_{i=1}^{m} x_t(i) f_{i,t}^{\text{in}}. \tag{10}$$

*Then $f_t^{\text{in}}$ are the probabilities that define $\mathbb{E}_t^Y$,*

**Theorem 9.** *Under the conditions of Theorem 2, and under Assumption 3, features (10) represent the probability of distribution over $\mathcal{Y}$ given $Z_t$.*

The existence of retrieved features that are not the same as physical features raises the question: which are encoded in (4)? A natural assumption is that it is the physical features, However, while economic models assume agents see reality unfiltered (Woodford (2020) discusses exceptions), neurobiology supports a notion of filtering based on memory-driven expectations (Reynolds and Chelazzi, 2004; Makino and Komiyama, 2015). What is retained in each re-remembering is not exactly what occurred, but rather a distorted copy of the original event (Rubin et al., 2008). Experimental evidence supports encoding of retrieved features at study (Greene, 1992; Siegel and Kahana, 2014), and at test (Zaromb et al., 2006; Howard et al., 2009; Miller et al., 2012; Kuhn et al., 2018). Encoding of retrieved features is a form of rehearsal, as discussed by Mullainathan (2002). It is realistic to assume that sometimes the agent encodes retrieved features and sometimes physical features. When one or the other occurs is outside the scope of the model.

Features retrieval implies similarity: better memory for an item similar to a just-remembered item. Features retrieval also implies an intuitive property of the model: similar contexts imply similar beliefs. In fact, beliefs are a weighted average of beliefs held under similar contexts. The process of features retrieval and encoding is how disparate features become "glued" together. As we show in the following section, it is

---

matrix $M_{t-1}$ and context $x_{t-1}$. The agent then experiences features $f_t$, and forms retrieved context $x_t^{\text{in}}$ and using $x_{t-1}$, context $x_t$. The agent then uses $M_{t-1}$ and $x_t$ to form $f_t^{\text{in}}$. Only after this is $M_{t-1}$ updated to $M_t$. This assumption mirrors that in the memory literature, and implies a memory "cycle" through with $M_{t-1}$ remains constant.

key to understanding temporal contiguity and the jump back in time.

## 2.4 Temporal Contiguity

Temporal contiguity serves as a fundamental organizing principle of the memory system (Healey et al., 2019). In this section we show, by means of an example, how the model accounts for temporal contiguity effects. We consider a stylized model of the Great Depression, in which an economic collapse follows a financial crisis. Subsequent to the Great Depression, many narratives emphasized the importance of the stock market crash of 1929, runs on financial institutions, and the connection between these purely financial events and the subsequent period of high unemployment and sharply decreasing output and consumption from which the Great Depression received its name. But from where did these narratives arise in the first place?

Let $x_{1929}$ denote context as of 1929, which we refer to as time $t-1$. Let $f_{\text{crisis}}$ be features associated with the failure of a financial institution, also occurring at $t-1$. Equation 4 implies that this combination of features and context become associated in memory:

$$M_{t-1} = M_{t-2} + x_{1929} f_{\text{crisis}}^{\top}. \tag{11}$$

Let $f_{\text{depression}}$ be features associated with depression occurring in the next period $t$. Call $x_{\text{depression}}^{\text{in}}$ the context retrieved by features $f_{\text{depression}}$:

$$x_{\text{depression}}^{\text{in}} \propto M_{t-1} f_{\text{depression}}. \tag{12}$$

Retrieved context $x_{\text{depression}}^{\text{in}}$ is the agent's state of mind when confronted with the observable features of the Great Depression, such as mass unemployment. Calling this state of mind $x_{\text{depression}}^{\text{in}}$ is simply terminology; there need not be anything "depression-like" about this context. Features $f_{\text{crisis}}$ and $f_{\text{depression}}$ are orthogonal, so that even though $M_{t-1}$ appears in both (12) and (11), the occurrence of the crisis nothing to do

19

with the retrieval of $x_{\text{depression}}^{\text{in}}$.[15]

Context evolution (2) implies that the retrieved depression context combines with previous context, to create the current context:

$$x_{1930} \equiv x_t = (1 - \zeta)x_{1929} + \zeta x_{\text{depression}}^{\text{in}}, \tag{13}$$

where $x_{1929}$ is the time-$(t-1)$ context and $x_{1930}$ is the time-$t$ context.[16] Crucially, $x_{1930}$ is a weighted average between $x_{\text{depression}}^{\text{in}}$ and $x_{1929}$, even if the events leading to the retrieval of $x_{\text{depression}}^{\text{in}}$ had nothing to do with the events of $x_{1929}$. From (4), it follows that

$$M_t = M_{t-1} + x_{1930}f_{\text{depression}}^{\top}. \tag{14}$$

Thus far, it is clear that $f_{\text{depression}}$ and $x_{1930}$ are associated, as are $f_{\text{crisis}}$ and $x_{1929}$. What is not yet clear is how $f_{\text{crisis}}$ relates to $f_{\text{depression}}$.

Suppose $f_{\text{crisis}}$ appears at some future time $t' > t$. As described in Theorem 8, there is a jump back in time – the agent retrieves $x_{1929}$:

$$
\begin{aligned}
x_{t'}^{\text{in}} \quad &\propto \quad M_{t'-1}f_{\text{crisis}} \\
&\propto \quad \left( M_0 + \sum_{s=1}^{t-2} x_s f_s^{\top} + x_{1929}f_{\text{crisis}}^{\top} + \sum_{s=t}^{t'-1} x_s f_s^{\top} \right) f_{\text{crisis}} \tag{15} \\
&\propto \quad x_{1929}. \tag{16}
\end{aligned}
$$

To obtain (16) from (15), we assume $f_{\text{crisis}}$ only appears once – in 1929, and that there are no other prior associations with $f_{\text{crisis}}$. These assumptions simplify the algebra without changing our conclusions, and we relax them in Appendix E. In any case, a financial crisis retrieves $x_{1929}$ – a crisis retrieves the context of the previous crisis. It does not retrieve $x_{1930}$. And yet the agent believes a depression is imminent. Why is this?

---

[15]More precisely, $x_{\text{depression}}^{\text{in}}$ would be retrieved whether or not the crisis had occurred.

[16]For simplicity, we refer to this as context in 1930, even though unemployment rose throughout the early 1930s.

The reason lies with the features retrieved by $x_{1929}$:[17]

$$
\begin{aligned}
f_{t'}^{\text{in}} \quad &\propto \quad M_{t'-1}^{\top} x_{1929} \\
&\propto \quad M_0^{\top} x_{1929} + \sum_{s=1}^{t'-1}(f_s x_s^{\top}) x_{1929} \\
&\propto \quad \underbrace{\left(M_0^{\top} + \sum_{s=1}^{t-2} f_s x_s^{\top}\right) x_{1929}}_{\text{prior associations}} + f_{\text{crisis}} + \underbrace{\sum_{l=1}^{t'-t}(1-\zeta)^l f_{t+l-1}}_{\text{depression features}} + \cdots \qquad (17)
\end{aligned}
$$

The re-appearance of the crisis retrieves features associated with $x_{1929}$ prior to the actual crisis (the first term in (17)). They retrieve $f_{\text{crisis}}$, because $x_{1929}$ is the very context under which the crisis was experienced. Most importantly, they retrieve $f_{\text{depression}}$ because $x_{1929}$ was part of context (13) at the time of depression. The time–$t$ term in (17) equals

$$
\begin{aligned}
f_t x_t^{\top} x_{1929} \quad &= \quad f_{\text{depression}} x_{1930}^{\top} x_{1929} \\
&= \quad f_{\text{depression}} \left((1-\zeta)x_{1929} + \zeta x_{\text{depression}}^{\text{in}}\right)^{\top} x_{1929} \\
&= \quad (1-\zeta) f_{\text{depression}}
\end{aligned}
$$

Thus re-appearance of a crisis retrieves the depression. If the depression features had previously continued for more multiple periods, as in fact occurred, the crisis context re-instates these as well, with geometrically declining weights.

Because we do not take a stand on whether context $x_{1929}$ is retrieved at some future time, there is a potential for additional terms in (17). That is, $x_s^{\text{in}}$, for $s = t+1, \ldots, t'-1$ might be correlated with $x_{1929}$. If, for example, a financial crisis occurred again while the agent was already in context $x_{1930}$, that would lead to additional terms in (17) containing $f_{\text{depression}}$ and would strengthen the associations between depressions and crises. If, on the other hand, a crisis occurred during a context very different from

---

[17]To simplify the algebra, we assume that $x_{t'} \approx x_{t'}^{\text{in}}$ as will be the case if crisis features are persistent. We also assume $x_{t'}$ is a basis vector. We relax these assumptions in Appendix E.

$x_{1929}$, the additional terms in (17) would be orthogonal to $f_{\text{depression}}$, and the association would weaken.

Note that retrieved features (17) form the probability distribution for the agent's expectations $\mathbb{E}^Y$, after having witnessed a financial crisis. We have seen that temporal contiguity leads the agent to place additional weight on the depression outcome. How much weight? One way to answer this question is to compare the probabilities to those retrieved by the depression: $f_{\text{depression}}$ retrieved by $x_{1929}$ versus $x_{1930}$. Substituting in the latter, we find:

$$f_{\text{depression}}^{\text{in}} \propto \underbrace{\left(M_0^\top + \sum_{s=1}^{t-2} f_s x_s^\top\right) x_{1930}}_{\text{prior associations}} + (1-\zeta)f_{\text{crisis}} + \underbrace{\sum_{l=1}^{t'-t}(1-\zeta)^{l-1}f_{t+l-1} + \cdots}_{\text{depression features}} \quad (18)$$

Both retrieve depression features, and with similar weights. This is purely because they occurred close in time.

It is useful to compare the role of associations in this model to a Bayesian one. Consider the reasoning that lies behind the formation of subjective Bayesian probabilities. The agent conceptualizes states of the world. The agent knows features that occur in each state of the world. There is a concept of events occurring "at the same time" – and yet in reality hardly anything occurs truly contemporaneously. If the agent's conceptualization of the underlying states happens to be correct, then arrival of one of the features is correctly taken as a signal that others will soon occur, and the agent's model of the world moves closer to the truth. The Bayesian set-up is convenient, but fragile. If the agent's prior did not allow for the correct features to signal a change in states (and there are a great many potential signals), there would be no way for the agent to learn.

Indeed, the "signal" argument requires that, in 1929, the agent foresaw the Great Depression, and then prior to 2008, believed that a signal such as one received in 1929 might recur, foreshadowing the next Great Depression. It seems more likely that in 1929, agents did not place some probability on the Great Depression being imminent;

22

had they, events might have unfolded differently. Thus they did not start with priors that the Bayesian analysis would require. Rather, after the Great Depression occurred, they associated the events (note that our analysis greatly reduces the sensitivity of the prior; $M_0$ is part of the analysis, but does not prevent the agent from forming other associations). The crisis followed by the depression *created an association that simply was not present before.* Then, after years went by and context shifted, the assumption became that great depressions, and for that matter financial crises, were a thing of the past (accorded zero probability), until events sufficiently similar to 1929 made individuals feel that the Great Depression was about to occur all over again. One might argue that the fact that individuals explicitly considered this possibility, whereas they had not in 1929 is what *prevented* the Great Depression from recurring in 2009.

One might object that this temporal association is unrealistic. Clearly agents should be able to distinguish between reminders and actual shifts in distributions. Simply because an event recurs, or a context changes should not change an agent's perspective. The well-known phenomenon of post-traumatic stress disorder, suggests that this idea is not so easily dismissed. In the next section we discuss the psychological and neural basis for contextual retrieval.

# 3 Psychological and Neural Basis for Contextual Retrieval

Before turning to specific applications, we summarize the psychological and neural evidence for context as an internal state.

In the memory laboratory, researchers create experiences by presenting subjects with lists of easily identifiable items, such as common words or recognizable pictures. Subjects attempt to remember these items under varying retrieval conditions: these include *free recall*, in which subjects recall as many items as they can in any order,

*cued recall*, in which subjects attempt to recall a particular target item in response to a cue, and *recognition* in which subjects judge whether or not they encountered a test item on a study list. In each of these experimental paradigms, memory obeys the classic "Laws of Association" which appear first in the work of Aristotle, and later in Hume (1748). The first of these is *recency*: human subjects exhibit better memory for recent experiences, *semantic similarity*: we remember experiences that are most similar in meaning to those we are currently experiencing, and finally, *temporal contiguity*: we remember items that occurred contiguously in time to recently-recalled items. Although quantified in the memory laboratory, each of these phenomena appears robustly in real-world settings, as described below.

A longstanding and persistently active research agenda in experimental psychology seeks to uncover the cognitive and neural mechanisms that could give rise to these regularities. Students of memory have proposed many hypotheses which they have tested in the laboratory. Some striking findings include the fact that recency and contiguity appear regardless of whether you measure memory for list items presented seconds apart or many minutes apart, or for autobiographical memories separated by days or weeks.

What gives rise to the recency and contiguity effects that appear ubiquitously in both laboratory memory experiments and in our daily lives? One influential class of explanations posits the existence of a fixed-capacity memory buffer, better known as "short-term memory." In such models, retrieval involves two stages: first, subjects report the items maintained in the short-term store; next, they search through long-term memory guided by interitem and context-to-item associations, and subject to interference from similar memories. This model accounts for contiguity owing to the strengthening of interitem associations among items that share time in the short-term store (Kahana, 1996). The classic work of Mullainathan (2002) builds on a short-term memory model.[18]

---

[18]Mullainathan's model is fundamentally one of short-term memory in that it has two states $R = 0, 1$ (in the notation of that paper). If an event has $R = 1$, it enters short-term memory, which it then

Although the short-term memory model produces recency and contiguity in immediate recall, it cannot readily explain why similar recency and contiguity effects appear for experiences that are widely separated in time, and thus neither likely to be present in short-term store at the time of recall, or to have occurred together in short-term store (Howard and Kahana, 1999; Healey et al., 2019).[19] Related to accounts based on short-term memory, neurobiological models of association posit that patterns of brain activity can associate with one another when they co-occur within a short time window governing synaptic plasticity (Abbott and Blum, 1996; Kempter et al., 1999). These models, however, also struggle to explain why robust temporal contiguity appear for temporally separated events.

These findings suggested the alternative retrieved context framework that we extend to economic choice behavior. According to this view, context evolves recursively by adding the retrieved past contexts associated with an item, remembered or experienced, to the prior state of context. The retrieved context will bear similarity to contiguously experienced items, generating the contiguity effect. Because retrieval depends on the *relative* similarities among competing items, strong contiguity effects can appear even for items separated by very long intervals. The same is true for recency effects, in both the model and in the data.

Figure 2 illustrates the temporal contiguity effect (TCE) and how it has provided empirical support for the idea of context retrieval. To measure the effect of contiguity on memory retrieval, researchers examine subjects' tendency to successively recall items experienced in proximate list positions. In free recall, this tendency appears as decreasing probability of successively recalling items $f_t$ and $f_{t+lag}$ as a function of $lag$, conditional on the availability of that transition (Kahana, 1996). This TCE reaches its

---

falls out of with an exponentially increasing probability. If $R = 0$, it is the generic memory store.

[19]Within economics, Mullainathan (2002) could be understood as a model of short-term memory. Items are available to be recalled (in the short-term store) or are not. They enter and exit the short-term store with probabilities determined by associations with current events. However, there is no mechanism by which an item from the distant past can evoke another item simply by virtue of being proximate in time. Nagel and Xu (2018) invoke long-term recency to explain economic phenomena, but do not otherwise employ associativeness.

maximum at $lag = \pm 1$, but also exhibits a forward asymmetry in the form of higher probability for positive as compared with negative lags. Equations 2–4 generate a forward asymmetry in the contiguity effect because recalling an item reinstates both its associated study-list context and its associated pre-experimental context. Whereas the study-list context became associated, symmetrically, with both prior and subsequent list items, the pre-experimental context became associated only with subsequently encoded list items, leading to a forward asymmetric contiguity effect, as seen in the data.[20]

Figure 1A shows that interitem distraction does not disrupt the TCE. Figure 1B-D shows that the TCE appears robustly for both younger and older adults, for subjects of varying intellectual ability, and for both naïve and highly practiced subjects. Figure 1E shows that the TCE appears even for transitions between items studied on distinct lists, despite these items being separated by many other item presentations. Figure 1F-H shows that the TCE also predicts confusions between different study pairs in a cued recall task, in errors made when subjects attempt to recall an individual list item in response to a sequential cue, and in tasks that do not depend on inter-item associations at all, such as picture recognition (see caption for details). Finally, long-range contiguity appears in many real-life memory tasks, such as recalling autobiographical

---

[20]Consider the example in Section 2.4, and treat the 1929, 1930 episodes as a "study phase." The model predicts that crises better recall economic contractions than the reverse. Let

$$x_{1929} = (1 - \zeta)x_{\text{prior}} + \zeta x^{\text{in}}_{\text{crisis}},$$

where $x^{\text{in}}_{\text{crisis}}$ is comprised of the 1929 stock market crash *and any previous financial crises*. These previous financial crises form the "pre-experimental" context associated with a crisis. The depression then becomes associated with this pre-experimental crisis context as well as with the 1929 crash because they are both part of $x^{\text{in}}_{\text{crisis}}$, $x^{\text{in}}_{\text{crisis}}$ is part of $x_{1929}$, and $x_{1929}$ is a part of $x_{1930}$ which co-occurs with $f_{\text{depression}}$.

In contrast, there are no means by which previous economic contractions can become associated with crises. Consider (13), where $x^{\text{in}}_{\text{depression}}$ consists of the Great Depression and previous economic contractions. Like crises, these previous economic contractions are part of "pre-experimental" context. However, whereas prior crises become associated with depressions, prior depressions do not become associated with crises. Why? Though these depressions are part of $x^{\text{in}}_{\text{depression}}$, they are orthogonal to $x_{1929}$, which is the means by which crises and depressions are associated. Thus there are two routes in memory by which a crisis recalls a depression: the 1929 crisis itself, and any previous crisis. However, there is only one route by which a large contraction recalls a crisis, and that is the Great Depression.

memories (Moreton and Ward, 2010) and remembering news events (Uitvlugt and Healey, 2019). These findings argue for a general associative memory mechanism, like context retrieval, that requires neither strict temporal proximity, nor specialized mnemonic strategies.

A second source of data in favor of retrieved context arises from neurobiology. The theory implies that brain states representing the context of an original experience reactivate or replay during the subsequent remembering of that experience. Several studies tested this idea using neural recordings. These studies found that in free recall (Manning et al., 2011), cued recall (Yaffe et al., 2014) and recognition memory (Howard et al., 2012; Folkerts et al., 2018) brain activity during memory retrieval resembles not only the activity of the original studied item, but also the brain states associated with neighboring items in the study list. Thus, one observes contiguity both at the behavioral and at the neural level, with these effects being strongly correlated (Manning et al., 2011). Finally, this recursive nature of the contextual retrieval process offers a unified account of many other psychological phenomena including the spacing effect (Lohnas and Kahana, 2014b), the compound cueing effect (Lohnas and Kahana, 2014a), and the phenomena of memory consolidation and reconsolidation (Sederberg et al., 2011).

Memory theory thus indicates that remembering an item involves a jump-back-in-time to the state of mind that obtained when the item was previously experienced. This reinstatement, in turn, becomes encoded with the new experience and also persists to flavor the encoding of subsequently experienced items. The persistence of the previously retrieved contextual states enables memory to carry the distant past into the future, allowing the contextual states associated with an old memory to re-enter one's life following a salient cue and associate with subsequent "neutral" memories. While the original memory is retained in association with its encoding context, the retrieval and re-experiencing of that memory forms a new memory in association with the mixture of the prior and retrieved context.

One might argue that, while retrieved context theory offers a persuasive account of

human memory phenomena, memories need not affect behavior, and still less, conscious decisions such as how much to invest in the stock market. While evidence on the role of experience in economic decision-making suggests otherwise, one might still argue that experience operates through a conscious process of attaining knowledge, rather than memory per se. Such a purely rational account would, however, miss important memory phenomena.[21] Evidence shows that agents re-live events from both the remote and recent past, often involuntarily (Rubin and Berntsen, 2009). If memory is a process of knowledge accumulation, it appears to be one that is outside of conscious control. An extreme example of the power of involuntary memory is post-traumatic stress disorder (PTSD), in which a traumatic event is not only "persistently reexperienced" but "causes clinically significant distress or impairment in social, occupational, or other important areas of functioning."[22] Studies show that PTSD is diminished when brain injury, childhood amnesia, or pharmacologically-induced amnesia blunts encoding, indicating that it is primarily a memory disorder (Rubin et al., 2008). As such, it exhibits patterns that are well-accounted for by retrieved context theory (Cohen and Kahana, 2020). Overall, evidence on PTSD suggests that it is best understood in terms of principles that govern "normal" memory functioning (Rubin et al.). In other words, there is no clear line separating trauma-induced and normal memories. It appears that people relive the past involuntarily and unawares, to the extent that they base their behavior on a biased representation of the external environment. In what follows, we show how this idea can account for economic phenomena that are difficult to explain otherwise.

---

[21]Though such an account would also have to explain why agents base their decisions on their *particular* experience.

[22]See the *Diagnostic and Statistical Manual of Mental Disorders* (4th ed., text revision.; APA, 2000, pp. 467–468).

# 4    Applications

This section describes three applications of our theory. Section 4.1 describes an application to portfolio allocation, illustrating how long-run stock-market experience might influence portfolio choice. Section 4.2 shows how memory dynamics might effect stock prices and interest rates in an otherwise standard macro-finance model in which circumstances lead an agent to recall a rare event such as an economic depression. Lastly, Section 4.3 shows how the model can account for the effects of changes in short-term context, thus explaining observed experimental effects on portfolio choice.

In each of these sections, rather than modeling the full lifetime of an agent's memories, we assume a self-contained decision problem. We follow the memory literature in making an assumption on the matrix $M_t$ prior to the decision problem at hand. Associations represented by $M_t$ are motivated by the temporal contiguity property (Section 2.4). That said, we do not generate the prior $M_t$ within the model. In particular, we require a memory matrix that is sparser than lifetime simulations of the process (2), (4) and (9) would imply. Augmenting the model with costly storage (Azeredo da Silveira and Woodford, 2019) could endogenously generate such sparsity while maintaining the model's ability to account for temporal contiguity. Encoding of retrieved features, generated using a process that downweights low probability items also "organizes" memory, leading to a sparser $M_t$ as noted by Polyn et al. (2009). Wachter and Kahana (2020) show that winner-take-all retrieval and encoding leads to a more organized $M_t$.[23]

## 4.1    Retrieved-context theory and the persistence of beliefs

A basic account of the behavior of financial markets and the macroeconomy requires that agents disagree (Lucas, 1975; Grossman and Stiglitz, 1980). Such disagreement poses a problem for standard Bayesian models in which agents begin with possibly

---

[23]Wachter and Kahana (2020) also shows that winner-take-all can account for the representative heuristic (Bordalo et al., 2021), and is therefore related to diagnostic expectations (Bordalo et al., 2018).

different priors, but nonetheless see the same data and take a rational view of others' beliefs beliefs. Recent evidence linking economic decisions to lifetime experience suggests that experience may be the place to look for understanding how disagreement arises and what causes it to persist (Malmendier and Nagel, 2011, 2016). The explanatory power of experience immediately suggests a role for memory.

We focus on the results of Malmendier and Nagel (2011), who show that experienced stock returns affect portfolio decisions, because the departure from rationality is particularly striking. The size of the equity premium – the expected return on stocks over Treasury bills – is discussed in finance textbooks, the media, and in popular books. Provided that the equity premium is positive, participation in the stock market is optimal (Arrow, 1971). Yet a large percentage of households do not participate in equity markets (Campbell, 2016). In what follows, we show how the theory in Section 2 can generate, based on life experience, permanent pessimism regarding stock returns. While survey evidence suggests agents are on average pessimistic (Goetzmann et al., 2017), the aim in this section is not to explain average pessimism, but why some agents remain pessimistic in the face of contrary data. This is what is required to account for non-participation.

### 4.1.1 The portfolio choice problem

Assume latent states $\mathcal{Z} = \{z_1, z_2\}$, and let $p < 1/2$ denote the unconditional probability of the adverse state $z_2$. The investor allocates wealth between a risky asset with net return $\tilde{r}$, and a riskfree bond with zero net return. The agent also receives risky labor income $\tilde{\ell}$. The set $\mathcal{Y}$ thus consists of joint outcomes of labor income and stock returns. Assume that the value of labor income is known given the state: $\ell(z_1) > \ell(z_2) = 0$. Assume that the stock return $\tilde{r}$ takes on values $r_g$ (gain) and $r_l < r_g$ (loss). Assume $\text{Prob}(r_g|z_1) = \frac{1}{2}(1-p)^{-1}$ and $\text{Prob}(r_g|z_2) = 0$ so that the marginal distribution of $\tilde{r}$ is 50% gains and 50% losses. Because gains do not occur in state 2, their probability is slightly elevated under the normal state $z_1$.

For simplicity, we assume mean-variance preferences:

$$\max_{\pi} \mathbb{E}[(1 + \pi\tilde{r} + \tilde{\ell})] - \frac{1}{2}\text{Var}(1 + \pi\tilde{r} + \tilde{\ell}), \tag{19}$$

where $\pi$ is the percent allocation to the risky asset. The expectation and the variance in (19) are with respect to the agents' subjective expectation in (1). Setting the derivative of the objective function with respect to $\pi$ equal to zero leads to

$$\pi = \frac{\mathbb{E}\tilde{r} - \text{Cov}(\tilde{r}, \tilde{\ell})}{\text{Var}(\tilde{r})}. \tag{20}$$

Let $\mu$ equal the mean of $\tilde{r}$ and $\sigma$ the standard deviation, so that $r_g = \mu + \sigma$, $r_\ell = \mu - \sigma$. Let $\tilde{\ell}(z_1) = \ell > 0$ and $\tilde{\ell}(z_2) = 0$, so that the optimal allocation (20) equals

$$\pi(p) = \frac{\mu - p\sigma\ell}{\sigma^2}. \tag{21}$$

The greater the probability that the agent assigns to the adverse state, the less he or she allocates to the risky asset.

### 4.1.2   Memory for stock market gains and losses

In the Bayesian benchmark (Theorem 1) the agent retains a perfect memory of gains and losses, and their associations with depressions. We consider the implications of contextual retrieval and encoding for the agent's beliefs and portfolio allocation. To highlight what is novel about our mechanism, we temporarily turn of the persistent feature of context. In the next subsection, we show how incorporating persistence matters, and how it can combine with contextual retrieval to generate a slowly-decaying experience effect.

Because both labor income states and stock returns influence utility, Assumption 1 implies that they both must be features of the environment.[24] Table 1 summarizes the

---

[24]Under our assumptions, gains cannot occur with a negative outcome of $\ell$. There are thus three possible features.

features space.

Table 1: Features corresponding to gains, losses, and depressions

| Basis Vector | Features | Outcome for wealth |
|---|---|---|
| $e_1$ | gain | $1 + \pi r_g + \ell$ |
| $e_2$ | loss | $1 + \pi r_l + \ell$ |
| $e_3$ | depression | $1 + \pi r_l$ |

The following matrix represents the agent's prior associations:

$$M_{t-1} \propto \begin{bmatrix} \frac{1}{2} & 0 & 0 \\ 0 & \frac{1}{2} - p^* & p^* \end{bmatrix} \tag{22}$$

The second row of $M_{t-1}$ has nonzero entries corresponding to columns for both loss and for depression. This means that loss and depression occur in the same context: a key assumption for the results that follow. Section 2.3 shows how associations of the form (22) arise, simply due to the fact that stock market losses preceded a depression.[25]

Prior to exploring the implications of the full model dynamics, we follow the format of Section 2.1 and consider the special case of $\zeta = 1$. While stark, this most directly contrasts the implications of the memory model with the standard Bayesian approach. We first explore retrieved context, and then retrieved features. If the agent experiences a gain, namely $f_t = e_1$, retrieved context is the first basis vector:

$$x_t = x_t^{\text{in}} \propto M_{t-1}e_1 = \begin{bmatrix} 1 \\ 0 \end{bmatrix}.$$

---

[25]More precisely, Section 2.4 shows these events would have contexts that are close to one another because they occur close in time. Having contexts that are close is sufficient to have context retrieval imply similar features retrieval. In order to focus on the main mechanism in this and in the following examples, we assume that the agent experiences depressions and stock market losses under the same context. It is reasonable to imagine memory consolidates contexts that are sufficiently close into a single context.

It follows from (9) that context in turn retrieves features corresponding to a gain:

$$f_t^{in} \propto M_{t-1}^\top \begin{bmatrix} 1 \\ 0 \end{bmatrix} = \frac{1}{2} e_1 \propto e_1$$

Now consider retrieved context in response to a loss ($f_t = e_2$):

$$x_t = x_t^{in} \propto M_{t-1} e_2 = \begin{bmatrix} 0 \\ \frac{1}{2} - p^* \end{bmatrix} \propto \begin{bmatrix} 0 \\ 1 \end{bmatrix}. \tag{23}$$

It follows from (9), that context (23) retrieves some probability on depression as well as on loss:

$$f_t^{in} \propto M_{t-1}^\top \begin{bmatrix} 0 \\ 1 \end{bmatrix} = \begin{bmatrix} 0 \\ \frac{1}{2} - p^* \\ p^* \end{bmatrix} \propto \begin{bmatrix} 0 \\ 1 - 2p^* \\ 2p^* \end{bmatrix} \tag{24}$$

All that happened was a stock market loss, but the loss brought to mind the depression, simply because the two shared a context.

Having recalled features (24), the agent encodes them with context (23). Memory $M_t$ evolves according to (4), with

$$x_t f_t^\top = \begin{cases} \begin{bmatrix} 1 \\ 0 \end{bmatrix} e_1^\top = \begin{bmatrix} 1 & 0 & 0 \\ 0 & 0 & 0 \end{bmatrix} & \text{if } gain \\ \begin{bmatrix} 0 \\ 1 \end{bmatrix} [0\ 1 - 2p^*\ 2p^*] = \begin{bmatrix} 0 & 0 & 0 \\ 0 & 1 - 2p^* & 2p^* \end{bmatrix} & \text{if } loss \end{cases} \tag{25}$$

After $\tau$ periods of which $k$ are gains:

$$M_{t+\tau} \propto \begin{bmatrix} \frac{1}{2}t + k & 0 & 0 \\ 0 & (\frac{1}{2} - p^*)t + (1 - 2p^*)(\tau - k) & p^* t + 2p^*(\tau - k) \end{bmatrix} \tag{26}$$

Note that the relative probability of losses and depressions, $M(2,2)/M(2,3)$, remains

33

the same (as in Theorem 3), regardless of how much experience the agent accumulates.[26] The agent does learn, in particular, about the relative probabilities of stock market losses and gains. In fact, the agent's beliefs will converge to the truth in this regard. However, the agent still overestimates the probability of a depression.[27]

In this stark example, the agent fails to update probabilities entirely. A loss makes the agent think of a depression by reinstating a context. This act of recalling the depression context is similar to experiencing the depression. Thus, a high probability of depression remains associated with losses in the agent's mind – thoughts extrapolative data.

We assumed $\zeta = 1$ as an purely for illustration. In practice, the speed at which the association weakens (if indeed it does weaken), depends on many factors. Below, we re-consider this example using the full model, and show that learning takes place, though it can be very slow. Figure 3 illustrates this case in a line with 'x' marks; note that after 20 years of data the probability still fails to converge to the Bayesian case, and as a result, the optimal portfolio choice remains significantly lower.[28] In contrast, the Bayesian agent quickly learns that depressions are unusual, though he or she will take a long time to learn the precise value.[29]

### 4.1.3 Generalizing to autoregressive context

This section generalizes the conclusions of the previous section to $\zeta < 1$. Consider a richer features space, but continue to assign $e_1$ to gains $e_2$ to losses, and $e_3$ to losses combined with a depression. For simplicity, assume the following form for initial

---

[26]Note that $((\frac{1}{2} - p^*)t + (1 - 2p^*)(\tau - k))/(p^*t + 2p^*(\tau - k)) = (\frac{1}{2} - p^*)/p^*$ regardless of $\tau$ or $k$.

[27]We can directly map the entries of $M_t$, which have the interpretation of unconditional probabilities, into the agent's decision problem by assuming a neutral "decision context", $x_t = \frac{1}{2}\hat{e}_1 + \frac{1}{2}\hat{e}_2$. This context implies retrieved features $f_t^{\text{in}} = [\frac{1}{2}, \frac{1}{2} - p^*, p^*]$, which implies (21), with $p^*$ substituted in for $p$.

[28]In calculating portfolio choice, we use (21), with $p^*$ taken from (26). Alternatively, we could use the idea of a neutral decision context as defined in the section below.

[29]For the purposes of the figure, $p = 0.02$, $p^* = 0.50$, $\sigma = 1$, and $\ell = 2$.

associations:

$$
M_0 \propto \begin{bmatrix} 1 & 0 & 0 & 0 & \cdots \\ 0 & 1-2p^* & 2p^* & 0 & \cdots \\ 0 & 0 & 0 & \hat{M}_0 & \\ \vdots & \vdots & \vdots & & \end{bmatrix},
\tag{27}
$$

where $\hat{M}_0$ represents associations other than gains, losses, and depressions.[30] Many paths of experience could influence the agent's memory for losses. Here, we will focus on one such path. Appendix F gives details of the arguments below.

Assume the agent experiences losses at times $t_1, t_2, \ldots, t_\ell, \ldots$. Also assume that the agent does not experience actual depressions; relaxing this assumption will strengthen our results. Context at the time of the first loss equals $x_{t_1} = (1-\zeta)x_{t_1-1} + \zeta\hat{e}_2$, where $x_{t_1-1}^\top \hat{e}_2 = 0$, and where features retrieved by $x_{t_1-1}$ put no weight on depressions. It follows from (10) that

$$
f_{t_1}^{\text{in}} = \zeta[0,\ 1-2p^*,\ 2p^*,\ 0,\ \ldots]^\top + (1-\zeta)f_{t_1-1}^{\text{in}}.
\tag{28}
$$

The subjective probability of depression (the third element of (28)) equals $2\zeta p^*$. Retrieved features are encoded with context $x_{t_1}$. Contextual drift will lead this agent to dis-associate the context loss with losses and depressions, thus putting a progressively lower weight on depressions.[31] If, for example, the agent experienced a single loss and saw no other features to remind him or her of a loss or a depression, then the weight on the depression would decay exponentially to zero. However, should another loss arise, even after an arbitrarily long time delay, the agent will recall the depression as if no

---

[30]The column sums in $M_0$ can be interpreted as the number of observations of each feature in a prior sample. See Appendix A, Lemma A.1. Our assumptions on these do not qualitatively effect the results.

[31]This conclusion would not necessarily hold under winner-take-all features retrieval. See Section 2.3 and the introduction to Section 4 for further discussion. Winner-take-all would lead retrieved features to continue to place high weight on the depression, assuming that agents' beliefs in a depression were above a threshold. In this respect, it would make it easier for us to derive our main result. However, incorporating winner-take-all features retrieval would complicate the analysis, which is why we do not assume it here.

time has passed.

Theorem 4 shows that retrieved context follows a recursion, assuming that encoding takes place with basis vector features. Appendix F proves an analogous result when encoding is with non-basis features. Let $\tilde{q}_t$ denote the agent's subjective probability of $z_2 \in \mathcal{Z}$, namely, $\tilde{q}_t = x_t(2)$. Assume at least one loss has taken place ($t > t_1$). Appendix F shows

$$f_{t,2}^{\text{in}} = \left(1 - \frac{\tilde{q}_{t-1}(1 - \tilde{q}_{t-1})}{1 + \sum_{s=1}^{t-1} \tilde{q}_s}\right) f_{t-1,2}^{\text{in}} + \frac{\tilde{q}_{t-1}(1 - \tilde{q}_{t-1})}{1 + \sum_{s=1}^{t-1} \tilde{q}_s} f_{t-1,2}^{\text{in},\perp}, \tag{29}$$

where

$$f_{t-1,2}^{\text{in},\perp} = (1 - \tilde{q}_{t-1})^{-1} \sum_{j \neq 2} x_{t-1}(j) f_{t-1,j}^{\text{in}}.$$

The third element of (29) gives the subjective probability of a depression.[32] Iterating on (29) shows how the subjective probability decays from its initial value of $2\zeta p^*$. As of time $T$, the subjective depression belief $\tilde{p}_T = f_{T,2}^{\text{in}}(3)$ equals

$$\tilde{p}_T = 2\zeta p^* \prod_{t=t_1}^{T} \left(1 - \frac{\tilde{q}_{t-1}(1 - \tilde{q}_{t-1})}{1 + \sum_{s=1}^{t-1} \tilde{q}_s}\right). \tag{30}$$

Note that the rate of decay depends on $\tilde{q}_{t-1}(1 - \tilde{q}_{t-1}) = (x_{t-1}^\top \hat{e}_2)(1 - x_{t-1}^\top \hat{e}_2)$. If $x_{t-1}$ is either orthogonal or proportional to $\hat{e}_2$, then $f_{t,2}^{\text{in}} = f_{t-1,2}^{\text{in}}$. In neither case do beliefs decay, because associations are unchanged. If $x_{t-1}$ has some, but not perfect, overlap with $\hat{e}_2$, then the agent learns new associations with the current context, causing beliefs to decay more quickly. Over time, the rate of decay slows, as captured by $1 + \sum_{s=1}^{t-1} \tilde{q}_s$. As the sum increases, the term multiplying $f_{t-1,2}^{\text{in}}$ also tends to increase, causing the process to become more persistent.

We now assume that, following each loss, context decays to a neutral value (one that is not associated with depression features). The most recent loss occurred at

---

[32] We make the conservative assumption that $f_{t-1,2}^{\text{in},\perp}$ places zero weight on the depression; if not, beliefs in depressions will be greater.

$t_\ell = \mathrm{argmin}_{t_j}\{t - t_j; t_j < t\}$, and $\tau = t - t_\ell$, the time elapsed since the loss event. Then

$$x_{t_\ell + \tau} = (1 - \zeta)^\tau \zeta \hat{e}_2 + (1 - (1 - \zeta)^\tau \zeta)\bar{x}, \tag{31}$$

is the required time path of context.[33]

We continue to assume that stock market gains and losses are equally likely.[34] We also assume that gains and losses occur one out of every $J$ periods, capturing the fact that the agent experiences other types of features (greater values of $J$ correspond to slower decay of probabilities). We assume that portfolio choice takes place when the agent is in a neutral decision context, namely a context that is $\frac{1}{2}(\hat{e}_1 + \hat{e}_2)$. This neutral decision context implies a depression probability of $p^* \prod_{t=t_1}^{T}\left(1 - \frac{\tilde{q}_{t-1}(1 - \tilde{q}_{t-1})}{1 + \sum_{s=1}^{t-1}\tilde{q}_s}\right)$. Figure 3 shows these values, and the resulting portfolio choice when we average over 1000 individuals and assume $J = 4$ and $p^* = 1/2$. The figure shows that slow decay of beliefs is not special to $\zeta = 1$; incorporating contextual drift allows for the more realistic conclusion that early memories still exercise influence, but that they fade gradually over time. This is key to explaining results such as Malmendier and Nagel (2011), in which early-life experience, while mattering less on average than recent experience, continues to have a statistically significant effect on choice.[35]

---

[33] We continue to assume that features that are losses retrieve context $\hat{e}_2$. Appendix F explains why this assumption is valid for $p^*$ close to $\frac{1}{2}$. In short: for $p^* \approx \frac{1}{2}$, retrieved features are more like depressions than losses. Thus actual losses retrieve the context associated with the initial loss.

[34] Strictly speaking, because depressions do not occur in the sample we are considering, losses should be slightly less likely; we ignore this effect here.

[35] In contrast, a Bayesian model would produce a null result in regressions on experience, because a relatively short amount of data suffices to determine optimal choice as described above. While more complicated Bayesian models might produce other patterns, it is hard to imagine one in which the early period takes on particular importance

## 4.2 Context and the jump back in time: Application to the financial crisis

The failure of Lehman Brothers is widely recognized as a point of inflection in the 2008 financial crisis.[36]

An open question is: why was the failure of Lehman Brothers so pivotal? A growing line of research answers this question by focusing on the importance of financial intermediation to the overall the economy. Brunnermeier and Sannikov (2014) and He and Krishnamurthy (2013) develop models in which the balance sheets of intermediaries contribute to business cycle fluctuations. However, while it may be necessary to have specialized institutions trade certain complicated investments, it is not clear why the failure of a financial institution should be followed by a broad-based stock market decline. Common stocks are not intermediated assets: trading costs for common stocks, already quite low for the past half-century, have only gotten lower (Jones, 2002). Another possibility is that Lehman represented a sunspot that caused a run on other intermediaries, and other forms of debt (Allen and Gale, 2009; Gorton and Metrick, 2012). Unanswered is why this should cause the stock market to crash, as it did in the fall of 2008, when most companies have very low leverage and can fund themselves through retained earnings?[37]

Gennaioli and Shleifer (2018) emphasize a third possibility: individuals and financial institutions took on too much debt because they incorrectly extrapolated from a recent low-risk environment. This debt created unstable conditions. The Lehman bankruptcy caused a sudden shift in beliefs by reminding agents of the risks they had forgotten. This account is most in spirit of the discussion here. Indeed, our hypothesis is that the financial crisis was a psychological event caused by the failure of Lehman Brothers.

---

[36]See, for example, French et al. (2010).

[37]Kahle and Stulz (2013) argue that firms dependent on bank-lending were not unduly affected by the crisis. Gomes et al. (2019) argue that fluctuations in borrowing conditions are more likely to be affected by investment opportunities than the other way around.

In the model described below, the failure of an important financial institution in the absence of insurance reminds investors of the Great Depression.[38] Some felt that they had – literally – returned to the Great Depression. Once this feeling entered the discourse, it proved hard to shake. Subsequent events showed that in fact there was no Great Depression; yet the association continued through a greatly renewed interest in financial crises and the macroeconomy, and through even the very name Great Recession.[39]

### 4.2.1  Endowment and preferences

Assume an endowment economy with identical log utility agents, each with time-discount factor $\beta$. Let $W_t$ denote an agent's wealth at time $t$, and let $C_t$ be consumption. Each agent solves

$$\max_{\pi_t, C_t} (1 - \beta) \log C_t + \beta \mathbb{E}_t[\log W_{t+1}] \tag{32}$$

subject to

$$W_{t+1} = (W_t - C_t) \left[ R_{f,t+1} + \pi_t (R_{t+1} - R_{f,t+1}) \right], \tag{33}$$

where, $\pi_t$ denotes the percent allocated to the risky asset and where, because agents are identical, we omit an agent subscript. Each agent trades in a risky asset with gross return $R_{t+1}$ and a riskless asset with gross return $R_{f,t+1}$ (known at time $t$).[40] We specify the aggregate endowment; asset prices will then equilibrate so that it is optimal for the agent to consume this endowment (Lucas, 1978). Let $g$ denote the growth rate in consumption during normal times, and $\delta \in [0, 1)$ the decline in consumption, should

---

[38]See, for example, the reporting of *The Guardian* on the day's events: https://www.theguardian.com/business/2008/sep/15/marketturmoil.stockmarkets.

[39]The model below is stylized; it cannot capture many interesting features of the Great Recession. One line of research in particular focuses on a channel from asset valuations to the real economy, either through real-business-cycle mechanisms (Gourio, 2012) or New Keynesian mechanisms (Caballero and Simsek, 2020). This literature can be viewed as taking beliefs in a changed regime as given and deriving joint implications for real outcomes and asset markets. Our model is about endogenizing the beliefs.

[40]Under time-consistent beliefs, the assumption of log utility implies that (32) is a recursive formulation of the multiperiod consumption and savings problem (Samuelson, 1969).

a depression occur:

$$\frac{C_{t+1}}{C_t} = \begin{cases} 1 + g & \text{with probability } 1 - p \\ (1+g)(1-\delta) & \text{with probability } p \end{cases} \tag{34}$$

The aggregate market is a claim to dividends $D_t$ satisfying $D_{t+1}/D_t = (C_{t+1}/C_t)^\phi$, with $\phi > 1$. The assumption of $\phi > 1$ captures the fact that payouts to shareholders fell by far more than consumption during the Great Depression (Longstaff and Piazzesi, 2004).[41] Equilibrium will require that $\pi_t = 1$, and that the dividend claim and the riskfree asset are in zero supply.

We briefly describe equilibrium under full information. The wealth-to-consumption ratio $W_t/C_t = 1/(1 - \beta)$. In equilibrium, the riskfree rate equals a constant given by

$$\begin{aligned} R_f &= \mathbb{E}\left[\beta \frac{C_t}{C_{t+1}}\right]^{-1} \\ &= \beta^{-1}(1+g)\left(1 + p\left(\frac{1}{1-\delta} - 1\right)\right)^{-1}. \end{aligned} \tag{35}$$

From (35), it follows that (in a comparative statics sense) an increase in the depression probability $p$ lowers the interest rate. An increase in the depression probability leads the investor to want to save today. Bond prices rise, and riskfree rates fall.

Let $S_t$ equal the value of the aggregate stock market. In equilibrium,

$$S_t = \mathbb{E}_t\left[\beta \frac{C_t}{C_{t+1}}(S_{t+1} + D_{t+1})\right], \tag{36}$$

---

[41]This specification implies dividends and consumption are perfectly conditionally correlated. In the data, dividends have greater normal-times volatility than consumption, and they are imperfectly correlated with consumption. Both facts could be introduced into the model by assuming that dividends also are subject to independent shocks. Because these shocks are unpriced, and assuming that we abstract from dividends as features about which the agent learns, they would have a negligible effect on the results of interest. The assumption that normal-times growth in dividends is $(1 + g)^\phi$ could similarly be relaxed without affecting the results.

Solving for a fixed point yields:

$$\frac{S_t}{D_t} = \frac{\beta(1 + p\left((1-\delta)^{\phi-1} - 1\right))}{(1+g)^{1-\phi} - \beta(1 + p((1-\delta)^{\phi-1} - 1))} \tag{37}$$

for all $t$. For $\phi > 1$, an increase in $p$ lowers the stock price. Realized returns on the stock market equal

$$R^S_{t+1} = \begin{cases} \bar{R}^S & \text{with probability } 1 - p \\ \bar{R}^s(1-\delta)^\phi & \text{with probability } p \end{cases} \tag{38}$$

where $\bar{R}^S$ is the stock market return during normal times:[42]

$$\bar{R}^S = \beta^{-1}(1+g)\left(1 + p\left(\frac{1}{(1-\delta)^{1-\phi}} - 1\right)\right)^{-1}. \tag{39}$$

The expected return for the stock market equals:

$$\mathbb{E}_t R^S_{t+1} = \bar{R}^S(1 + p((1-\delta)^\phi - 1)). \tag{40}$$

Subtracting the log of the riskfree rate (35) from the log of the expected return (40) give the equity premium. For small $p$ (in the continuous-time limit), the equity premium is well-approximated by

$$\log \mathbb{E}_t \left[R^S_{t+1}/R_f\right] \approx p\left((1-\delta)^{-1} - 1\right)\left(1 - (1-\delta)^\phi\right).$$

The right hand side is is the negative change in marginal utility multiplied by the change in stock price during depressions.

---

[42]Let $\Phi = \beta(1+g)^{\phi-1}(1 + p\left((1-\delta)^{\phi-1} - 1\right))$. Then

$$\bar{R}^S = \frac{S_{t+1}/D_{t+1} + 1}{S_t/D_t}(1+g)^\phi = \frac{\Phi/(1-\Phi) + 1}{\Phi/(1-\Phi)}(1+g)^\phi = \Phi^{-1}(1+g)^\phi$$

### 4.2.2 Features, context, and memory

Now assume that agents update beliefs according to retrieved-context theory. An application to asset pricing immediately raises the question of agent's beliefs about others' beliefs. Here we assume agents have identical initial associations $M_0$, identical experiences (and thus the same $M_t$), and that beliefs are common knowledge. Relaxing these assumptions would be desirable and would have interesting implications.

There are three possible outcomes in $\mathcal{Y}$, captured by basis features vectors: $e_1$ (normal), $e_2$ (crisis), and $e_3$ (depression).[43] At each time $t$, the agent observes $W_t$ and $f_t$ (we abstract from wealth as a feature). We conjecture an equilibrium in which the agent possesses the correct mapping between $f_t$ and $r_{f,t+1}$ and $f_t$ and $r_t$. The agent solves (32–33) under the subjective expectation, formed from memory as Section 2.3 describes.[44] The agent's problem satisfies Assumption 1 in that it depends only on the distribution of future features.[45]

The agent's optimal consumption continues to equal $1/(1-\beta)$ times wealth. Consistent with (10), the subjective probability of depression (which is next period's feature) is $p_t = e_3^\top f_t^{\text{in}}$. The agent will allocate $\pi_t = 1$ to the consumption claim provided that $r_{f,t+1}$ satisfies (35), and $r_{t+1}$ satisfies (38) and (39) with $\phi = 1$. The assumption that $p_t$ is permanent (agents do not foresee a change in their own, or others' beliefs) implies that (37) is also satisfied in equilibrium. Table 2 summarizes features and their effects on outcomes of interest.

Equilibrium returns on the consumption claim equal growth in consumption, scaled

---

[43] We assume two underlying states in $\mathcal{Z}$, one in which crises take place with greater probability. While it is not necessary to take a stand on whether there is a physical association between a depression and a financial crisis, for ease of comparison with the full-information case above, we assume that the true correlation between disasters and crises equals zero. That is, disasters are equally likely in the two states.

[44] The agent retrieves context $x_t^{\text{in}}$ from features $f_t$ and memory $M_{t-1}$. Retrieved context and lagged context $x_{t-1}$ combine to determine $x_t$ through (2). Using $M_{t-1}^\top$, the agent retrieves $f_t^{\text{in}}$, which gives the probability distribution over future features.

[45] We assume the agent is capable of conceiving of the equilibrium and calculating quantities that are directly implied, through algebraic equations, to features (this is required for optimization as well as equilibrium). This is a strong assumption, which we make to focus on one departure from the standard model at a time.

Table 2: Features corresponding to normal times, financial crises, and depressions

| Basis Vector | Features | Consumption claim return | Crisis? |
|---|---|---|---|
| $e_1$ | normal | $\beta^{-1}(1+g)$ | No |
| $e_2$ | crisis | $\beta^{-1}(1+g)$ | Yes |
| $e_3$ | depression | $\beta^{-1}(1+g)(1-\delta)$ | Yes |

by $\beta^{-1}$, and thus depend only on the realization of a depression. While the agent uses current features (such as crises) to form beliefs as in (1), and indeed the value function depends in equilibrium on the riskfree rate and hence on current features, the only *future* feature of interest to the agent is whether or not there will be a depression.

Prior to the crisis, we assume associations take the form:

$$
M_{t-1} \propto \begin{array}{ccc} \text{normal} & \text{crisis} & \text{depression} \\ \left[ \begin{array}{ccc} 1-p^c & 0 & 0 \\ 0 & p^c(1-q) & p^c q \end{array} \right] \end{array} \tag{41}
$$

As in the previous example, two sets of features share a context. In this case, it is a financial crisis and a depression. We see this from the fact that $M_{t-1}$ has nonzero entries in its second row. As in that example, the motivation is as in Section 2.3: associations of the form (22) arise, simply due to the fact that stock market losses preceded a depression.

Assume that there has been a sufficiently long period of normal features $e_1$. so that the agent exhibits neglected risk: namely context $x_{t-1} = [1,0]^\top$ (Theorem 6).[46] Though agents neglect the depression state, they have not forgotten it. Representing the failure of Lehman brothers is $f_t = e_2$, the well-publicized failure of a major financial institution. Retrieved context in response to crisis features equals

$$
x_t^{\text{in}} \propto M_{t-1}e_2 \propto [0,1]^\top. \tag{42}
$$

---

[46]Strictly speaking, this requires a large number of neutral features; it is simple to alter this example to include these but it also makes the notation more complicated.

Context therefore equals

$$x_t = (1 - \zeta) \begin{bmatrix} 1 \\ 0 \end{bmatrix} + \zeta \begin{bmatrix} 0 \\ 1 \end{bmatrix} = \begin{bmatrix} 1 - \zeta \\ \zeta \end{bmatrix}. \tag{43}$$

As in Section 2.3, we calculate features retrieved by each component of context:

$$f_{1,t}^{\text{in}} \quad \propto \quad M_{t-1}^{\top} \begin{bmatrix} 1 \\ 0 \end{bmatrix} \propto e_1$$

$$f_{2,t}^{\text{in}} \quad \propto \quad M_{t-1}^{\top} \begin{bmatrix} 0 \\ 1 \end{bmatrix} \propto \begin{bmatrix} 0 \\ 1 - q \\ q \end{bmatrix}$$

Therefore, applying (43) and (10),

$$f_t^{\text{in}} = \begin{bmatrix} 1 - \zeta \\ \zeta(1 - q) \\ \zeta q \end{bmatrix}. \tag{44}$$

Equation 44 represents reinstatement of the depression context – the financial crisis reminds the agent of the depression. This same expression also gives the probabilities that the agent uses in the decision problem; that is, they are the probabilities underlying $\mathbb{E}_t^Y$. The subjective probability of a depression rises from 0 to $\zeta q$, causing an immediate decline in stock prices (37) and in the riskfree rate (35).[47] The sharp drop between $t = 0$ and $t = 1$ of Figure 4 illustrates this effect on the price-dividend ratio and on the riskfree rate.[48] Figure 4 also shows a sharply negative return corresponding to this

---

[47]Note that Assumption 3 implies that the agent views this change as permanent. Equation 36 should thus be understood as the price of a future stream of dividends assuming this increased probability. It is equivalent to substitute the probability directly into the right-hand side of (37).

[48]When we report the riskfree rate in the figure, we assume a zero lower bound. That is, we assume that for institutional reasons, the observed riskfree rate cannot fall below zero, whereas the true riskfree rate might

event.[49] Note that while both the recent past of economic calm and the depression cue are part of recent experience, we might well say that the "jump-back-in-time" effect dominates the recency effect that is the focus of prior work on memory and economics. Pure recency would suggest positive economic news dominates (because the bank failure, by construction, does not affect consumption). The novel and unique cue from bank failures, however, powerfully reinstates the context of the depression.

What does the model say about the time path of context, and hence that of prices, interest rates, and stock returns, following the event? We discuss in detail one such possible path. Consistent with events in late 2008, we assume several (specifically, three) observations of crisis features, and then normal features. Assume a sufficiently long prior sample so that updating memory is not first-order.[50] If the agent continues to observe crisis features, context updates as follows:

$$x_{t+1} = (1 - \zeta)^2 \begin{bmatrix} 1 \\ 0 \end{bmatrix} + (\zeta(1 - \zeta) + \zeta) \begin{bmatrix} 0 \\ 1 \end{bmatrix} \tag{45}$$

Thus recall of the depression increases, the stock price declines further, and realized returns continue to be negative. Figure 4 shows this continued decline. Thus while the initial drop was an over-reaction relative to the correct iid benchmark (Theorem 8), there is also a sense in which it was, in the short run, an under-reaction because prices fall more before they stabilize (Theorem 7). Given sufficient crisis observations, the stock price would stop declining once context reached a steady state of $[1 - q, q]^\top$, implying zero weight on normal features. In this example, however, assume that normal features return after three periods, leading to a partial recovery. Returns are positive and high, because they represent good news that there has not been a depression.[51] Prices recover more slowly than they fall. In this model with autoregressive beliefs,

---

[49] We assume $q = 0.5$, $p^c = 0.05$, $\delta = 0.15$, $\beta = .98$, $g = 0$, $\phi = 2$.

[50] The solution shown in Figure 4 assumes a prior sample of 100 years and calculates the exact path of context.

[51] Due to log utility, the equity premium is small in this model, even with a high probability of a rare. Little of realized returns correspond to an equity premium.

this is mainly due to the effects of duration embedded in (37). An increase in the probability of disaster decreases the effective maturity of the stock return, and so any increase in $p$ from this low point has a smaller effect.

This duration explanation cannot, however, account for the fact that the price-dividend ratio asymptotes to a lower level. This is because as depression features are retrieved, not only are they encoded with the crisis context (as in Section 4.1), but they are also encoded with the normal context. The agent then associates this context with depression, so that even a return to a normal context implies a permanently elevated probability of depression. More precisely, the agent associates whatever was in context just before the crisis with depression, even though the depression did not occur. The continued retrieval of the depression context ensures that this is a permanent effect. To the extent that (previously normal features) reminiscent of 2009 re-appear, they will remind agents of the Great Recession, which in turn will recall the Great Depression.[52]

## 4.3 Fear and asset allocation

Psychological and neuroscientific research reveals a tight link between memory and emotion. For example, when people are sad or depressed, they tend to recall negative events (Matt et al., 1992; Teasdale and Fogarty, 1979). When people remember an emotionally-valent word (positive or negative) in a list of mixed-valence words, the next word they remember tends to be emotionally congruent (Long et al., 2015; Siddiqui and

---

[52]While Figure 4 represents a clear departure from a full-information benchmark, another benchmark of interest is one in which a Bayesian agent believes in the existence of two states and that the financial crisis serves as a signal for depression. That is, the Bayesian agent shares the associations $M_{t-1}$, and thus is as close as possible to the agent we consider. Wachter and Zhu (2019) model Bayesian learning in such a setting. The Bayesian agent does not neglect risk (or at least not to the same degree): there is always some probability on the depression state regardless of how long a period of normalcy occurs. The total decline in price relative to the Bayesian benchmark represents an overreaction; the Bayesian agent believes the state will revert. However, unlike the Bayesian agent, there is also under-reaction in the sense that prices do fully react immediately to the "news" — they take several periods to respond. The permanent adjustment in prices due to the new associations shown in Figure 4 is absent in the Bayesian model, nor does the Bayesian model explain how an agent came to associate crisis and depression in the first place. The initial slow adjustment and the permanent change in prices also differentiates the model from that of Gennaioli and Shleifer (2018). Note that this comparison also assumes the agent begins with conditional probabilities that the present model endogenizes through temporal contiguity.

Unsworth, 2011). Emotion also affects memory for neutral items — subjects exhibit superior recall when tested under a state that is emotionally congruent to the study state (Eich, 1995).

The connection between emotion and memory extend to the financial domain. Guiso et al. (2018) found that after the 2008–209 financial crisis, professional investors required twice the premium to accept a risky bet rather than a sure payoff than before, suggesting a role for fear in decision-making. In this case, the finance professionals' response need not be due to fear; an alternative is that they were materially worse off following the crisis, and that they exhibit risk aversion that increases as wealth falls. To demonstrate that it is the memory and not the wealth change that influences risk aversion, Guiso et al. conduct an experiment in in which subjects are randomly assigned to view a scene from a horror movie. Subjects who viewed the scene required a 50% greater premium to accept the lottery as compared to those that did not. Similar results are found by Cohn et al. (2015), who conduct an experiment in which financial professionals are selected randomly to view a chart of a stock market boom versus a crash. Investors in the boom condition invested significantly more in the risky asset that those in the crash condition. The effect of emotion on financial decisions extends beyond the laboratory. Cuculiza et al. (2020) show that analyst earnings forecasts become more negative upon anniversaries of terrorist attacks (as well as following the attacks themselves). Ramadorai et al. (2020) examine the trading behavior of investors following the outcome of IPO lotteries. Investors who received shares in companies that subsequently performed well not only purchase similar stocks, but trade more in general. While Loewenstein (2000) proposes a model by which emotion influences utility, his model is silent on the connection between emotion and recent experience.

In what follows, we apply our framework to explain findings on fear and aversion to risk. We specifically consider the set-up of Guiso et al. (2018), in which an agent chooses between a risky asset (a lottery) and a sure investment. As in Section 4.1, the agent chooses an allocation $\pi$ to a risky investment with net return $\tilde{r}$, and $1 - \pi$ to a safe investment with net return of zero. Also as in Section 4.1, the agent is subject to

a second source of risk. While we focused on job loss in that example, here we would want to think more broadly in terms of anything that might severely and suddenly impair the human capital of the agent, analogous to physical danger. For convenience, we continue to refer to this as human capital. The risky investment (capturing the lottery in the experiment) experiences a gain or a loss, each with equal probability:

$$\tilde{r} = \begin{cases} \mu + \sigma & \text{prob. } 1/2 \\ \mu - \sigma & \text{prob. } 1/2. \end{cases}$$

Wealth equals $1 + \pi\tilde{r} - \tilde{\ell}$, where $\tilde{\ell}$ is the second source of risk mentioned above:

$$\tilde{\ell} = \begin{cases} 0 & \text{prob. } 1 - p \\ \delta & \text{prob. } p. \end{cases}$$

In what follows, we refer to $\tilde{\ell}$ as a human capital shock.

Because $\tilde{r}$ is the outcome of a lottery, $\tilde{\ell}$ and $\tilde{r}$ must be independent. Moreover, $\tilde{r}$ has a well-defined set of outcomes (it obeys the Savage (1954) model). Unlike in Section 4.1, we cannot rely on the agent's misperception of the correlation between the stock return and human capital. While it is possible that the agent has such a misperception, here we assume that agents understand that the outcome of the lottery does not bear on other events.

While the agent is told that the outcome of the lottery is 50/50, the probability $p$ is unknown. The agent determines the more ambiguous probabilities of $\tilde{\ell}$ based on memory. Assume that the agent has log utility over wealth, so that perceived probabilities of $\tilde{\ell}$ influence the utility of the lottery (any power function over wealth would have this property). The agent feels more fearful of a bad outcome, and thus is less likely to take on risk. Whereas the connection between an increased likelihood of a bad outcome and a willingness to take on risk is not irrational, the increased fear of the bad outcome is.

As in Section 4.1.3, we assume the first features vector corresponds to a normal

state, whereas the next two correspond to the negative outcomes. There are also a large number of "neutral" features, i.e. neither positive or negative. The matrix $M_{t-1}$ takes the form:

$$
M_{t-1} \propto
\begin{array}{cccc}
\text{normal} & \text{fin. crisis} & \text{danger} & \text{other associations} \\
\end{array}
\begin{bmatrix}
1 - p_1 - p_2 & 0 & 0 & 0 & \cdots \\
0 & p_1 & p_2 & 0 & \cdots \\
0 & 0 & 0 & \hat{M}_t & \\
\vdots & \vdots & \vdots & &
\end{bmatrix}
$$

The second context corresponds to a state in which negative events occur. Temporal contiguity could account for the associations between a financial crisis and danger in this context. The relative values of columns indicate how much of each feature the agent has experienced. Given our focus on context and feature retrieval, this will not be important in what follows. Features retrieved by each component of context are:

$$
\begin{aligned}
f_{1,t}^{\text{in}} &\propto M_{t-1}^{\top}\hat{e}_1 \propto e_1 \\
f_{2,t}^{\text{in}} &\propto M_{t-1}^{\top}\hat{e}_2 \propto (p_1 + p_2)^{-1}(p_1 e_2 + p_2 e_3)
\end{aligned}
$$

Let $\tilde{q}_t = x_t(2)$ denote the subjective probability of the second state (in which $\tilde{\ell} = \delta$), and assume that, prior to the experiment, $\tilde{q}_{t-1} = p_1 + p_2 \equiv p$.

The scene from the movie is a feature that is similar to, but not exactly the same as, danger. That is, $f_t \approx e_3$, with the difference representing movie-specific features. Assuming $f_t$ is sufficiently close to $e_3$, then Theorem 5 applies: retrieved context will be similar to what would occur under actual danger. That is $x_t^{\text{in}} \approx \hat{e}_2$, and $\tilde{q}_t \approx (1-\zeta)\tilde{q}_{t-1} + \zeta$. If the agent had watched this particular scene on a prior occasion, then the approximation $x_t^{\text{in}}$ would be further from $\hat{e}_2$ because movie-specific features would pull up the context from the last time the movie was viewed. Guiso et al. (2018) chose a scene likely to be unfamiliar to subjects.

Ultimately it is retrieved features that matter for decision-making. It follows from

(10) that

$$f_t^{\text{in}} \approx (1 - \tilde{q}_t)e_1 + \tilde{q}_t \propto (p_1 + p_2)^{-1}(p_1 e_2 + p_2 e_3), \tag{46}$$

so that the probability of danger or depression equals $\tilde{q}_t$. The agent chooses $\pi$ to maximize:

$$\mathbb{E}^\ell \left[ \frac{1}{2} \log(1 + \pi(\mu + \sigma) - \tilde{\ell}) + \frac{1}{2} \log(1 + \pi(\mu - \sigma) - \tilde{\ell}) \right], \tag{47}$$

where the outcome $\tilde{\ell} = \delta$ occurs with probability $\tilde{q}_t$ and is 0 otherwise.[53] We assume an excess return $\mu = 4\%$, a standard deviation $\sigma = 20\%$, a prior probability of the negative labor market outcome $p = 2\%$, and a percent decline $\delta = 0.8$, should the outcome occur. As elsewhere, $\zeta = 0.35$. The agent trades off the higher return arising from greater $\pi$ with greater risk. The higher is the probability of a bad realization, the less risk the agent can afford to take. The experiment reminds agents that such bad realizations can occur.

Figure 5 shows (47) as a function of the allocation $\pi$. At an initial level of zero, taking on a small value of risk is optimal (the function is increasing). After a certain level of $\pi$, the function begins to fall, representing the curvature in the utility function. When the agent is not fearful, this occurs at 70%. When the agent is fearful, it occurs at 30%. The response of the agent to the experiment cannot be Bayesian: a movie has not changed anything about the outside world. In that sense, the response of risk-taking to viewing a horror movie is a good test of our theory. This, and related results, show that it is possible to manipulate internal mental state (context) in a way that changes decision-making. Context alters agents beliefs about the outside world, and hence agents' decisions.

---

[53]In the notation of the general decision problem from Section 2, (Assumption 1), $V(Y_{t+1}, \pi) = \frac{1}{2} \log(1 + \pi(\mu + \sigma) - \tilde{\ell}(Y_{t+1})) + \frac{1}{2} \log(1 + \pi(\mu - \sigma) - \tilde{\ell}(Y_{t+1}))$, and $\mathbb{E}^Y$ arises from the probabilities (46).

# 5 Retrieved Context Theory and Alternative Memory Models

Our theoretical framework for modeling human memory—retrieved context theory—builds upon scholarship going back to the early 20th century. McGeoch's classic (1932) theory of forgetting assumed that memories do not wither away in time, but rather that the retrieval cues used to search memory may either reveal or occlude a particular experience.[54] McGeoch theorized that this retrieval-based interference depends on the state of context, the mental "set" of the rememberer, and the activation of similar "competitor" memories. Estes (1959) and Bower (1972) developed a mathematical foundation for these ideas, positing a VAR context representation that determined the retrievability of items from memory, accounting for the law of recency (Crowder, 1976). These models also allowed for explicit manipulations of context that could alter the memorability of particular experiences, as seen in experimental studies.

Unlike recency effects, which are easy to quantify, the influence of contiguity eluded careful measurement (Murdock, 1974). In studying the order and timing of recall sequences, (Kahana, 1996) introduced a conditional-probability measure of contiguity by computing the likelihood of recalling an item as a function of its contiguity to the just recalled item. Howard and Kahana (1999) theorized that contiguity could arise from retrieval of context, rather than direct interitem associations, and subsequent work provided support for this account (see, Healey et al., 2019, for a review).

Retrieved-context theory provides a unified account of recency and contiguity effects at short and long time scales. As a vector-based model of associative learning, this theory nests earlier work on similarity-based organization in memory and cue-dependent interference effects (Kahana, 2012). Two key aspects of this model make novel predictions not shared by most other memory theories: First, each encoding

---

[54]Applications 1 and 3 seem, at first glance, to be consistent with this notion of "withering away." However, these two examples show that any such withering must occur at vastly different time scales, already revealing that there is more than meets the eye.

event involves an internal retrieval whose output modifies memory. Thus, thoughts become memories that influence subsequent retrieval. Second, the recursive definition of context predicts a forward asymmetry in memory retrieval because the contextual states previously associated with an item recombine with the items context and associate with subsequent items. The recursive contextual dynamics imply that when a new event matches, or closely resembles, an earlier experience, the context of that earlier experience will re-embed in memory; analogously, new contexts that resemble old contexts will tend to retrieve old features, which will also re-embed in memory. In the present theory of economic choice, these re-embeddings help to explain persistent disagreement in the face of accumulating evidence.

Whereas modern memory theorists have adopted McGeoch's early emphasis on retrieval processes, other classic models saw memory-guided choice as primarily reflecting the strength of memories established during learning. Consider a repeated event, such that each occurrence supplies an agent with information about the world (e.g., ordering a cappuccino and learning about its price). According to strength theory, the association between a cappuccino and its price continuously updates, and a future reminder of a cappuccino retrieves a single sufficient statistic reflecting the distribution of experienced prices. In this model, one updates the summary statistic and discards the memories of each experienced event. A fundamental problem facing strength theories of memory is defining the unit of memory; at some point a new experience is sufficiently different from an existing "memory" to constitute a new memory, but the model is silent as to how that point is determined.

Unlike strength theory, similarity-based exemplar theory assumes that the memory system separately records each event (feature vector) in an ever growing matrix of memories (Estes, 1986; Hintzman, 1988; Nosofsky, 1992). The events retain their individuality, including their ordinal position within the series. Aggregation of memory happens at the time of retrieval, when the memory system compares a cue event against all previously stored examples, computing the summary statistic at the time of test (as in Bordalo et al. (2020b)). Although a wealth of data favor exemplar over

strength-based models, the latter have resisted extinction due to their parsimony and computational efficiency (Murdock, 1985; Wixted, 2007).

Exemplar and strength-based models share a major limitation: they lack any mechanism for associating events that co-occur in a spatiotemporal context (that is, they lack a mechanism to explain temporal contiguity). Such associations have long held a central role in philosophical conceptions of the "association of ideas" (Hume) and form the basis for the Aristotelian "Law of Contiguity". Continental philosophers (Herbart, 1834) developed theories based on chained associations that stimulated the earliest experimental work on human memory (Ebbinghaus, 1913). These theories did not posit a specific representation of time, but rather assumed that contiguously experienced events become associated and that the strength of this associations falls as a function of the temporal separation of the events (Solway et al., 2012). The repetition of an item, evoked by an external stimulus or an internal retrieval, triggers retrieval of the item's neighbors as a function of the strength of their association. Chaining theory encounters serious obstacles when two lists of sequentially presented items share an overlapping item, or an overlapping subsequence of items. In this case, one cannot recall either list without suffering catastrophic interference between the competitor items (i.e the items that follow the overlapping items in both lists). Such interference prevents the model from recovering order information, or even accurately recalling a series with repeated or high similar elements (Lashley, 1951; Henson et al., 1996; Kahana and Jacobs, 2000). Despite the failure of many of its predictions, chaining theory retains a powerful appeal based in part on the everyday experience of sequential cuing of memories.

The idea that memories preserve a record of their spatial information, and the complementary observation that spatial cuing offers an aid to learning and retention, formed a centerpiece of the medieval "arts of memory" (Yates, 1966). Soon after the advent of list memory studies in the late 1800s, researchers recognized that subjects often visualized a series of items as occurring within a virtual, mental, space, much as medieval scholars used the "palaces of memory" to commit vast written works to

memory (Ladd and Woodworth, 1911). This idea of positional coding as a means of representing ordinal information offered an alternative to the classic idea of chained associations described above. Positional coding models assumed that in learning a series of items, subjects formed associations between items and positions (location on an array, or position within an ordinal series). Later, at the time of test, subjects used positional information, assumed to be a mental primitive, to cue retrieval of items. Modern positional coding theories offer some of the most successful accounts of memory of short, ordered lists: phone numbers and postal codes and the like (Burgess and Hitch, 2006; Brown et al., 2000).

Although one can imagine items in a series as occupying locations in space, they actually occupy locations in time. For early researchers, this raised the question of whether memories retain information about their time of occurrence. Although the notion of time tagging appeared in some of the earliest psychological literature (James, 1890) memory scholars did not take stock of its significance until the emergence of experimental procedures that required subjects to explicitly judge the temporal order of studied items (Yntema and Trask, 1963; Hinrichs and Buschke, 1968) and models developed to explain these data (Bower, 1972). Consistent with the aforementioned notions of spatial and temporal coding of memories, neural recordings in both human and non-human animals have identified individual neurons in the brain that encode spatial information (O'Keefe and Dostrovsky, 1971; Ekstrom et al., 2003), and temporal information (MacDonald et al., 2011; Pastalkova et al., 2008; Umbach et al., 2020). Further, the activity of neural ensembles coding time and place during memory storage predict aspects of subsequent recall even when the recall task requires neither memory for the time or place of occurrence (Miller et al., 2013; Umbach et al., 2020).

An exemplar theory which posits that agents store a complete record of the attributes describing every memory would support optimal choice. But the marked failures of memory in our daily lives and the presumed finite capacity of the memory system pose a serious challenge for such a model. If memory failures primarily reflect inaccurate storage, there must be a massive degree of "lossy" data compression on

54

the front end. A commonly adopted view is that initial processing occurs via a limited capacity short-term storage system; only a tiny fraction of information survives to make it into long-term storage. Contrary to this limited storage view, Gallistel and King (2009) summarize extensive data, and computational arguments, in favor of the brain's ability to store massive amounts of information. Failures of memory, they argue, reflect failures of retrieval rather than storage (see, also, Tulving and Madigan (1970)). Using sensitive indices of retrieval one can show that a once experienced laboratory event can leave a near-permanent record in memory (Kolers and Magee, 1978; Standing, 1973; Brady et al., 2011).

Thus, to build a choice model upon an exemplar theory with complete memory requires a fully-specified model for retrieval.[55] The principles of memory used to motivate our modeling approach–similarity, recency and contiguity–would need to emerge from the hypothesized retrieval process. Most retrieval choice rules will naturally give rise to similarity effects (Kahana, 2012, Chapter 4); see also Bordalo et al. (2020b). Recency effects arise based on the similarity of temporal codes between study and test (recency is another form of similarity). The challenge, however, facing these models is to explain the contiguity effect and its persistence across time scales. The ability to retrieve temporal codes associated with past memories, and use those codes to retrieve subsequent memories, requires substantial machinery that is not part of exemplar models. Although a chaining model, as described above, can produce associations among nearest neighbors, it does not easily account for associations that span multiple items as seen in the data. These findings require something like the contextual retrieval process used in this paper.

To summarize, a half-century of scholarship has shown that memory depends critically on the cues present at the time of retrieval, and that retrieval operates during both learning and recall. Retrieval during learning determines the information present on subsequent learning and recall trials. Scholars no longer see encoding and retrieval

---

[55]This statement is also true for exemplar models in which time is an explicit characteristic, e.g. SIMPLE (Brown et al., 2007).

as distinct phases of memory; both processes play an important role during the acquisition of new knowledge and the during the recall or recognition of past experiences. The retrieved context framework builds upon classic notions of contextual variability and cue-dependent recall to offer a unified account of the principles of recency, contiguity and similarity, and their persistence across time scales.

The three applications in Section 4 illustrate the force of dynamic contextual retrieval in economic settings. Consider the first illustration: one of the role of early-life experience. A canonical exemplar model with a lossy front-end might be considered that of Sims (2003). Such models and successors (Gabaix, 2019) can explain under-reaction because items take a long time to get into memory. However, they cannot explain the curious persistence of memory. This requires selective retrieval after a long gap, suggesting that the problem did not lie with the initial storage.

Exemplar models extended to incorporate static context, as in Bordalo et al. (2020b) are closer. In these models, re-appearance of a static physical context could trigger long-distant memories. However, while static models appear to offer a simpler explanation for these phenomena, to get them to work requires a greater burden of assumptions. In all three applications, time is a crucial variable.[56] The first and third applications require a latent inner state; pessimism retrieved by a stock market loss in the first and fear retrieved by a scene in the movie in the third. The second application is literally a jump-back-in-time, drawing on the notion of temporal contiguity. To explain this using a static model would require not only collapsing time dynamics (the financial crisis and the Great Depression would have to co-occur), but arbitrarily assigning to the financial crisis the role of context. As this section shows, many of these challenges and questions have had parallels in the literature on human memory.

---

[56]This may seem least clear in the third application, but note that even here the experiment is performed after the scene is viewed, not in the middle of the scene as would be expected if evaluating physical context.

# 6 Conclusion

What makes us know what we know, perceive what we perceive, think what we think? What makes us the same person when we get up in the morning as we were the day before? What makes life a connection of meaningful events, and not just a random set of stimuli? It is our memories – the experience of our lives that is unique to each individual.

The standard model in economics would have it otherwise. Under this framework, individuals maximize expected utility, seeing the form of the utility function and beliefs about the future as fundamentally stable. These assumption arise from pure reason – under a specific notion of rationality. Although empirical and experimental studies have challenged this framework, its parsimony and the appeal of assuming agents smarter than ourselves have led to its continued use.

In this paper we propose an alternative that is also parsimonious, potentially rational, and based on a centuries-old program in experimental psychology, namely, the study of memory. We show that principles emerging from this program offer a very different set of implications than the standard ones in economics.

First, decisions can be affected by seemingly irrelevant information, making revealed preferences far less stable than they might be otherwise. Second, individuals have trouble processing information, not because of the amorphous idea of lack of attention, but because of a stable internal state. Because we carry with us an internal state, we cannot "take in" our surroundings all at once. This internal state, however, allows us to tag memories in time, with the powerful consequence of temporal contiguity, a property of memory noted since ancient times. Temporal contiguity strings together events, pulling up an entire universe with one memory, and forms a basis for conjectures of causal behavior. Finally, retrieval and encoding of memories implies that our beliefs may not converge, regardless of how many data points we observe. The reason is that the same data comes, for everyone, with its own associations – its own context. This context then triggers an actual perception of different data.

This paper represents a start in connecting memory with decision-making. Many questions remain unanswered. We assumed a single decision-maker solving a static problem (though memory itself is dynamic). A key question pertains to the decision-maker's view not only of future self, but of other selves in the economy. Like beliefs about the physical world, these may also be contextually-dependent. A second question pertains to the boundaries between the type of recall of probabilities that we consider here, and model-driven decisions. To some extent we all do use models; at what level does the model come into play? Finally, economics concerns itself with maximization under constraints. Perhaps memory evolved to solve some maximization problem, but if so, which one? These are some of the questions we hope will be answered in future work.

# A  Proofs of results in Section 2.1

The following Lemma clarifies that the elements of $M_t$ can be thought of as proportional to probabilities, with a prior distribution given by $M_0$. The scaling in (3) implies that the absolute magnitude of the elements in $M$ is irrelevant. The time path of retrieved context, will the same if one employs a modified updating rule for $M$ that scales the sum of its elements to equal one. The following Lemma makes this statement precise.

**Lemma A.1.** *Let $\{x_t\}_{t=1}^T$ be the time path of context up to time $T$ given $\{f_t\}_{t=1}^T$, $M_0$, $x_0$, and updating rule (4). Let $t_0 = \iota^\top M_0 \iota$, namely, the sum of the elements in $M_0$. Let $\{\tilde{x}_t\}_{t=1}^T$ be a time path of context given the same features, initial condition $\tilde{M}_0 = \frac{1}{t_0} M_0$, $x_0 = \tilde{x}_0$, and updating rule*

$$\tilde{M}_t = \frac{t_0 + t - 1}{t_0 + t} \tilde{M}_{t-1} + \frac{1}{t_0 + t} \tilde{x}_t f_t^\top. \tag{A.1}$$

*Then $\tilde{x}_t = x_t$ and $\tilde{M}_t = \frac{1}{t_0+t} M_t$.*

**Proof.** Assume by induction that $\tilde{x}_{t-1} = x_{t-1}$ and $\tilde{M}_{t-1} = (t_0 + t - 1)^{-1} M_{t-1}$. It follows from (3) that $x_t^{\text{in}} = \tilde{x}_t^{\text{in}}$, implying that $x_t = \tilde{x}_t$. It remains to show that $\tilde{M}_t = (t_0 + t)^{-1} M_t$. By (A.1),

$$\tilde{M}_t = \frac{t_0 + t - 1}{t_0 + t} \tilde{M}_{t-1} + \frac{1}{t_0 + t} x_t f_t^\top \tag{A.2}$$

Recall $\iota^\top M_0 \iota = t_0$. Further recall that $\iota^\top f_t = \iota^\top x_t = 1$. It follows that the elements of the outer product matrix $x_t f_t^\top$ sum to one. Using the induction step and substituting into (A.2) implies

$$\tilde{M}_t = \frac{1}{t_0 + t} M_t,$$

as required. $\qquad\square$

It might seem that the updating rule (A.1) contains more information than (4) in that (A.1) includes both the current sample size $t$ and a prior sample size $t_0$. Lemma A.1

says that this intuition is not correct and that they contain the same information. The reason is that the extra information in (4) is contained in the size of $M$ itself. The sum of the elements in $M$ equals the size of the sample, whereas the sum of the elements of $\tilde{M}$ equals one. The updating rule (A.1) takes information that was previously embedded in $M$ and puts it into the updating rule.

**Proof of Theorem 1.** We show that, given data on latent states and an observed features $f_t$, $M_t(i, j)$ is proportional to the posterior probability of the co-occurence of latent state $i$ and observed state $j$.

Under the assumptions of the theorem, context $x_t$ is such that $x_t(i) = 1$ if $Z_t = z_i$ and 0 otherwise. Consider a Bayesian agent with a Dirichlet prior, with parameters given by $M_0$. Specifically, let $K = mn$, and let $P$ be the $m \times n$ matrix of prior probabilities $p_{ij}$ over state $(z_i, y_j)$:

$$\text{vec}(P) \sim \text{Dir}(K, \text{vec}(M_0)).$$

Note that the prior mean of $p_{ij}$ is $M_0(i, j)/t_0$, with $t_0 = \iota^\top M_0 \iota$. Suppose the agent observes $\{Y_t, Z_t\}_{t=1}^t$. Let $\hat{Y}_s = (Y_s, Z_s)$. Assume the agent computes a quasi-likelihood function under the assumption that observations are iid:[57]

$$\mathcal{L}(\hat{Y}_1, \ldots, \hat{Y}_t | P) = \prod_{t=1}^t l(\hat{Y}_s | P)$$

Each term $l(\hat{Y}_s | P)$ is multinomial. The posterior distribution $p(P | \hat{Y}_1, \ldots, \hat{Y}_t, M_0)$ is therefore Dirichlet (Gelman et al., 2004):

$$\text{vec}(P) \,|\, \{Y_t, Z_t\}_{t=1}^t \sim \text{Dir}(K, \text{vec}(M_t)).$$

The mean of the posterior distribution for $p_{ij}$ is then $M_t(i, j)/(t_0 + t)$ as required. $\quad\square$

---

[57]Alternatively, the agent could use the true likelihood that would allow for autocorrelation. This would necessitate a prior over $p_{ij}^Z$ for all pairs $(i, j)$. The quasi-likelihood function avoids this complication.

**Lemma A.2.** *Given Assumption 2, the following two invariance conditions are equivalent:*

1. *For any positive integer $s, t$, if $f_s = f_t$ then $x_s^{\text{in}} = x_t^{\text{in}}$.*

2. *For any basis vector $\bar{f}$, and any non-negative integers $s, t$, $M_s \bar{f} \propto M_t \bar{f}$.*

**Proof.** Assume condition 1 above. Assume $f_s = f_t \equiv \bar{f}$. By (3) and condition 1

$$M_{s-1} \bar{f} = M_{s-1} f_s \propto x_s^{\text{in}} = x_t^{\text{in}} \propto M_{t-1} f_t = M_{t-1} \bar{f},$$

proving condition 2.

Now assume condition 2. Let $f_s = f_t = \bar{f}$. Then by (3) and condition 2,

$$x_s^{\text{in}} \propto M_{s-1} f_s \propto M_{t-1} \bar{f} \propto x_t^{\text{in}}.$$

Equality follows because elements of retrieved context must sum to 1. This proves condition 1. $\square$

Note that the proof does not, strictly speaking, require Assumption 2. It holds for any subset of possible basis vectors (assuming we restrict $x^{\text{in}}$ accordingly). In what follows, we will mainly be concerned with the case in which the features are basis vectors.

**Proof of Theorem 3.** Given Lemma A.2 and $\zeta = 1$, it suffices to show that, for any (basis) features vector $\bar{f}$, and any non-negative integer $s, t$, $M_s \bar{f} \propto M_t \bar{f}$.

For convenience, normalize one of the times to 0. We prove, by induction, that for non-negative integers $t$:
$$M_t \bar{f} \propto M_0 \bar{f}.$$

Clearly the statement holds for $t = 0$. Assume

$$M_{t-1} \bar{f} \propto M_0 \bar{f}. \tag{A.3}$$

It follows from (4) that

$$M_t = M_{t-1} + x_t f_t^\top. \tag{A.4}$$

It follows from (2), (3), and $\zeta = 1$ that

$$x_t = (\|M_{t-1} f_t\|)^{-1} M_{t-1} f_t. \tag{A.5}$$

Substituting (A.5) into (A.4) implies

$$M_t = M_{t-1} + (\|M_{t-1} f_t\|)^{-1} M_{t-1} f_t f_t^\top. \tag{A.6}$$

Consider some basis vector $\bar{f}$. It follows from (A.6) that

$$M_t \bar{f} = M_{t-1} \bar{f} + (\|M_{t-1} f_t\|)^{-1} (M_{t-1} f_t)(f_t^\top \bar{f}). \tag{A.7}$$

First assume $\bar{f} \neq f_t$. Then $f_t^\top \bar{f} = 0$, the second term on the right-hand side of (A.7) equals zero and

$$M_t \bar{f} = M_{t-1} \bar{f}.$$

Now assume $\bar{f} = f_t$. Then $f_t^\top \bar{f} = 1$ and

$$M_t \bar{f} = M_{t-1} \bar{f} + (\|M_{t-1} f_t\|)^{-1} M_{t-1} \bar{f} \propto M_{t-1} \bar{f}.$$

Thus for any basis vector $\bar{f}$,

$$M_t \bar{f} \propto M_{t-1} \bar{f} \propto M_0 \bar{f},$$

where the second statement of proportionality follows from the induction step (A.3). $\square$

**Lemma A.3.** *Assume thus far an agent has experienced an event at times* $\{t_1, \ldots, t_\ell\}$.

*Then Assumption 2 implies*

$$x_{t_\ell}^{\text{in}} = \frac{1}{\ell} \left( x_{t_1}^{\text{in}} + \sum_{k=1}^{\ell-1} x_{t_k} \right). \tag{A.8}$$

**Proof.** Without loss of generality, assume the event is represented by basis vector $e_1$. A direct application of (3) implies that, should this occur at time $t$:

$$x_t^{\text{in}} \propto M_{t-1} e_1 \tag{A.9}$$

Substituting in from (5), we find:

$$x_t^{\text{in}} \propto M_0 e_1 + \sum_{s=1}^{t-1} x_s (f_s^\top e_1) \tag{A.10}$$

If the agent experienced the event at time $s < t$, $f_s^\top e_1 = 1$; otherwise it is equal to zero (note that the Lemma assumes all features are basis vectors). Therefore:

$$x_t^{\text{in}} \propto M_0 e_1 + \sum_{s \in \{t_1, \dots, t_\ell\}} x_s$$

Note that the vector on the right hand side has elements summing to $\ell$. The result follows from the fact that elements of $x_t^{\text{in}}$ must sum to 1. $\square$

**Proof of Theorem 4.** Consider retrieved context for the $\ell$th event. A slight rewriting of (A.8) implies

$$
\begin{aligned}
x_{t_\ell}^{\text{in}} &= \frac{1}{\ell} \left( x_{t_1}^{\text{in}} + \sum_{k=1}^{\ell-2} x_{t_k} + x_{t_{\ell-1}} \right) \\
&= \frac{1}{\ell} \left( \left( x_{t_1}^{\text{in}} + \sum_{s=1}^{\ell-2} x_{t_k} \right) + (1-\zeta) x_{t_{\ell-1}-1} + \zeta x_{t_{\ell-1}}^{\text{in}} \right) \tag{A.11}
\end{aligned}
$$

The second line follows from (2):

$$x_{t_{\ell-1}} = (1 - \zeta)x_{t_{\ell-1}-1} + \zeta x^{\text{in}}_{t_{\ell-1}}.$$

The term in the inner parentheses in (A.11) equals $(\ell - 1)x^{\text{in}}_{t_{\ell-1}}$. This follows from (A.8), applied to $x^{\text{in}}_{t_{\ell-1}}$. Therefore,

$$x^{\text{in}}_{t_\ell} = \frac{1}{\ell}\left((\ell - 1)x^{\text{in}}_{t_{\ell-1}} + (1 - \zeta)x_{t_{\ell-1}-1} + \zeta x^{\text{in}}_{t_{\ell-1}}\right).$$

Collecting terms in $x^{\text{in}}_{t_\ell-1}$ establishes the result. □

**Proof of Theorem 2.** Assume $\{Z_t\}$ is iid. Define a matrix $P$ such that $P(i, j) = p(z_i, y_j)$, the joint probability of $z_i$ and $y_j$. Then $P(i, j) \propto M_{t-1}$, with the constant of proportionality equal to the sum of the elements of $M_{t-1}$. Suppose $Y_t = y_j$. Then $f_t$ is the $j$th basis vector and

$$
\begin{aligned}
x^{\text{in}}_t &= \frac{M_{t-1}f_t}{||M_{t-1}f_t||} \\
&= \frac{Pe_j}{||Pe_j||} \\
&= \left(\sum_{i=1}^{m} p(z_i, y_j)\right)^{-1}\begin{bmatrix} p(z_1, y_j) \\ \vdots \\ p(z_m, y_j) \end{bmatrix}.
\end{aligned}
$$

Note that $\sum_i p(z_i, y_j)$ is simply the unconditional probability of $y_j$. Thus

$$x^{\text{in}}_t(i) = p(z_i, y_j)\left(\sum_{i=1}^{m} p(z_i, y_j)\right)^{-1} = p(z_i \,|\, y_j),$$

the conditional probability of $Z_t = z_i$, given $Y_t = y_j$. □

**Proof of Theorem 5.** It suffices to show that, as a function of features elements, $x^{\text{in}}_t$ is uniformly continuous, where we define continuity by the $L^1$-norm. We first show

64

that unscaled $x_t^{\text{in}}$ is uniformly continuous. Define $t_0 = \iota^\top M_0 \iota$. We show that, for any $\epsilon > 0$, there exists a $\delta > 0$ such that

$$\left\| \frac{1}{t_0 + t}(M_0 + \sum_{s=1}^{t} x_s f_s^\top)(f_t - \hat{f}_t), \right\| < \epsilon \tag{A.12}$$

provided that $\|f_t - \hat{f}_t\| < \delta$.

Standard triangle inequality argument imply that it suffices to show that

$$\left\| \frac{1}{t} \sum_{s=1}^{t} x_s f_s^\top (f_t - \hat{f}_t) \right\| < \epsilon/2 \tag{A.13}$$

$$\left\| \frac{1}{t_0} M_0 (f_t - \hat{f}_t) \right\| < \epsilon/2 \tag{A.14}$$

For (A.13), note that $f_s^\top (f_t - \hat{f}_t)$ is a scalar, so that

$$\left\| \frac{1}{t} \sum_{s=1}^{t} x_s f_s^\top (f_t - \hat{f}_t) \right\| = \frac{1}{t} \sum_{s=1}^{t} \|x_s\| \|f_s^\top (f_t - \hat{f}_t)\|$$

$$\leq \frac{1}{t} \sum_{s=1}^{t} \|x_s\| \|f_s\| \|f_t - \hat{f}_t\|$$

$$= \frac{1}{t} \sum_{s=1}^{t} \|f_t - \hat{f}_t\| = \|f_t - \hat{f}_t\|. \tag{A.15}$$

For (A.14), note that

$$\left\| \frac{1}{t_0} M_0 (f_t - \hat{f}_t) \right\| = \left\| \frac{1}{t_0} \sum_{j=1}^{m} M_0(i,j)(f_t(j) - \hat{f}_t(j)) \right\|$$

$$\leq \left\| \frac{1}{t_0} \sum_{j=1}^{m} M_0(i,j) \right\| \max_j \{|f_t(j) - \hat{f}_t(j)|\}. \tag{A.16}$$

It suffices then to choose $\delta$ so that the right-hand side of (A.15) and (A.16) are less than $\epsilon/2$.

We now extend this argument to show that $x_t^{\text{in}}$ is uniformly continuous.[58] Define

---

[58] We treat the previous history of contexts and features, as well as $M_0$, as fixed. Then $x_t^{\text{in}}$ is a

65

$x_{*t}^{\text{in}}$ to be unscaled $x_t^{\text{in}}$:

$$
\begin{aligned}
x_{*t}^{\text{in}} &= \frac{1}{t_0 + t}(M_0 + \sum_{s=1}^{t} x_s f_s^{\top}) f_t \\
\hat{x}_{*t}^{\text{in}} &= \frac{1}{t_0 + t}(M_0 + \sum_{s=1}^{t} x_s f_s^{\top}) \hat{f}_t.
\end{aligned}
$$

We have shown $\|\hat{x}_{*t}^{\text{in}} - x_{*t}^{\text{in}}\| < \epsilon$. Our aim is to show $\|\hat{x}_t^{\text{in}} - x_t^{\text{in}}\| < \epsilon$. Because we are interested in the limit for a fixed $f_t$ (given $t$), and therefore a fixed $x_{*t}^{\text{in}}$, it suffices to show, with a suitable adjustment to $\epsilon$, that

$$
\|\hat{x}_t^{\text{in}} - x_t^{\text{in}}\| \|x_{*t}^{\text{in}}\| < \epsilon.
$$

Finally note that

$$
\begin{aligned}
\|\hat{x}_t^{\text{in}} - x_t^{\text{in}}\| \|x_{*t}^{\text{in}}\| &= \left\| x_t^{\text{in}} \|x_{*t}^{\text{in}}\| - \hat{x}_t^{\text{in}} \|\hat{x}_{*t}^{\text{in}}\| + \hat{x}_t^{\text{in}} \|\hat{x}_{*t}^{\text{in}}\| - \hat{x}_t^{\text{in}} \|x_{*t}^{\text{in}}\| \right\| \\
&= \left\| (x_{*t}^{\text{in}} - \hat{x}_{*t}^{\text{in}}) - \hat{x}_t^{\text{in}} \left( \|\hat{x}_{*t}^{\text{in}}\| - \|x_{*t}^{\text{in}}\| \right) \right\| \\
&\leq \|x_{*t}^{\text{in}} - \hat{x}_{*t}^{\text{in}}\| + \|\hat{x}_t^{\text{in}}\| \left| \|\hat{x}_{*t}^{\text{in}}\| - \|x_{*t}^{\text{in}}\| \right| \\
&= \|x_{*t}^{\text{in}} - \hat{x}_{*t}^{\text{in}}\| + \left| \|\hat{x}_{*t}^{\text{in}}\| - \|x_{*t}^{\text{in}}\| \right| < \epsilon,
\end{aligned}
$$

provided that $\|x_{*t}^{\text{in}} - \hat{x}_{*t}^{\text{in}}\| < \epsilon/2$. $\qquad\square$

# B   Proofs of results in Section 2.2

**Definition** (Associated features). *Features vector $\bar{f}$ and state $Z_t = z_i$ are associated at time $t$ if either one of the two conditions hold:*

*1. $\hat{e}_i^{\top} M_0 \bar{f} \neq 0$*

---

function of time and of $f_t$. We let $\hat{f}_t \to f_t$ and show that the convergence of $x_t^{\text{in}}$ does not depend on $t$. It will, however, depend on the choice of $f_t$ because it depends on the scale of $\frac{1}{t_0+t}(M_0 + \sum_{s=1}^{t} x_s f_s^{\top}) f_t$ as the subsequent argument makes clear.

*2. There exists an $s \leq t$ such that for $f_s^\top \bar{f} \neq 0$, $x_s(i) \neq 0$.*

If $\bar{f}$ is associated with state $i$, then the agent has either experienced features $\bar{f}$ in a context that places weight on state $i$, or initial memory associates $\bar{f}$ with state $i$.

**Definition** (Uniquely associated features). *Features vector $\bar{f}$ and state $Z_t = z_i$ are uniquely associated at time $t$ if $\bar{f}$ is only associated with state $i$ at $t$.*

If $\bar{f}$ is uniquely associated with state $i$, it can only retrieve state $i$.

**Notation.** *Let $\Omega_{i,t} \subset \mathcal{B}^n$ denote the set of features uniquely associated with state $i$ at time $t$. Let $\Omega_{i,t}^\perp \subset \mathcal{B}^n$ denote the set of features not associated with state $i$ at time $t$.*

**Lemma B.1.**    *1. Features retrieve a context placing weight on $i$ if and only if these features are associated with state $i$.*

*2. Features uniquely associated with state $i$ retrieve only state $i$ (if $f_{t+1} \in \Omega_{i,t}$, then $x_{t+1}^{in}(i) = 1$).*

**Proof.** Let $\bar{f}$ be features at time $t+1$. Consider retrieved context (3):

$$
\begin{aligned}
x_{t+1}^{in} &\propto M_t \bar{f} \\
&\propto M_0 \bar{f} + \sum_{s=0}^{t} x_s (f_s^\top \bar{f}),
\end{aligned}
$$

it follows that

$$
x_{t+1}^{in}(i) \propto \hat{e}_i^\top M_0 \bar{f} + \sum_{s=0}^{t} x_s(i)(f_s^\top \bar{f})
$$

The right hand side is nonzero, if and only if $\bar{f}$ is associated with state $i$.

Now assume $\bar{f}$ is uniquely associated with state $i$. Then, for $k \neq i$, $x_{t+1}^{in}(k) = 0$. Because the elements of the context vector sum to 1, $x_{t+1}^{in}(i) = 1$.    □

**Lemma B.2** (Context reset). *For a given integer $\tau > 0$, assume the agent experiences a sequence of features $f_{t+1}, \ldots, f_{t+\tau} \in \Omega_{i,t}^\perp$. Also assume the features are orthogonal to one another. Then*

$$
x_{t+\tau}(i) = (1 - \varsigma)^\tau x_t(i). \tag{B.1}
$$

*Thus as $\tau \to \infty$, $x_{t+\tau}(i) \to 0$.*

**Proof.** We prove (B.1) by induction on $\tau$. It holds trivially for $\tau = 0$. Assume (B.1) holds for $\tau - 1$. Given features $f_{t+1}, \ldots, f_{t+\tau}$, it follows from (4) that

$$M_{t+\tau-1} = M_t + x_{t+1} f_{t+1}^\top + \cdots + x_{t+1} f_{t+t+\tau-1}^\top$$

Assume features $f_{t+\tau}$ not associated with $i$ at $t$ and orthogonal to $f_{t+1}, \ldots, f_{t+\tau-1}$. It follows from (3) that

$$
\begin{aligned}
x_{t+\tau}^{\text{in}} &\propto M_{t+\tau-1} f_{t+\tau} \\
&\propto M_t f_{t+\tau} + x_{t+1} f_{t+1}^\top f_{t+\tau} + \cdots + x_{t+1} f_{t+\tau-1}^\top f_{t+\tau} \\
&\propto M_t f_{t+\tau}
\end{aligned}
$$

Thus, lack of association at time $t$, $x_{t+\tau}^{\text{in}}(i) = 0$. Recall that we have assumed by induction that $x_{t+\tau-1}(i) = (1 - \zeta)^{\tau-1} x_t(i)$. Then (B.1) follows from (2). $\square$

Exponential decay of context, combined with a sufficiently large number of orthogonal features, implies the possibility of context "reset." Context reset is not a mere mathematical construction: novel features are used in the memory laboratory to reset context. These novel features, presumably orthogonal to the features that the agent has recently experienced, are introduced through "distractor tasks" that often involve solving arithmetic problems under time constraints (Howard and Kahana, 1999).

**Proof of Theorem 6.** See Lemma B.2. $\square$

The sets $\Omega_{i,t}$ and $\Omega_{i,t}^\perp$ define subjective associations. It is useful to have notation for the analogous concepts in the physical world.

**Notation.** *Let $\mathcal{Y}_i \subset \mathcal{Y}$ denote the set of outcomes that can* only *occur in state $i$. That is:*

$$\mathcal{Y}_i \equiv \{y_j \in \mathcal{Y} : p(y_j \,|\, z_i) > 0 \ \& \ p(y_j \,|\, z_k) = 0, \forall k \neq i\}.$$

Let $\mathcal{Y}_i^\perp$ denote the set of outcomes that cannot occur in state $i$:

$$\mathcal{Y}_i^\perp \equiv \{y_j \in \mathcal{Y} : p(y_j|z_i) = 0\}.$$

**Definition** (Correct associations). *The agent has correct associations with state $i$ if the agent's associations reflect reality. That is: $\Omega_i = \mathcal{Y}_i$, and $\Omega_i^\perp = \mathcal{Y}_i^\perp$.*

If features unambiguously signal a state, then context shifts in the direction of that state.

**Lemma B.3.** *Consider a nonempty state $i$ such that $\mathcal{Y}_i^\perp = \mathcal{Y} \setminus \mathcal{Y}_i$. Assume at time $t-1$ that the agent has correct associations with state $i$. Then*

1. *For $x_{t-1}(i) \in [0,1)$, $x_t(i) > x_{t-1}(i)$ if and only if $Z_t = z_i$.*

2. *For $x_{t-1}(i) = 1$, $x_t(i) = x_{t-1}(i)$ if and only if $Z_t = z_i$.*

**Proof.** First rewrite (2) as:

$$\Delta x_t(i) = \zeta(x_t^{\mathrm{in}}(i) - x_{t-1}(i)). \tag{B.2}$$

Under the stated assumptions, $f_t \in \Omega_{i,t-1}$ implies $Y_t \in \mathcal{Y}_i$ and $f_t \in \Omega_{i,t-1}^\perp$ implies $Y_t \in \mathcal{Y}_i^\perp$. Moreover, by Lemma B.1, $f_t \in \Omega_{i,t-1}$ implies $x_t^{\mathrm{in}}(i) = 1$ and $x_t^{\mathrm{in}}(i) = 0$ otherwise. It follows that if $x_t^{\mathrm{in}}(i) = 1$, we must have $Y_t \in \mathcal{Y}_i$. If $x_t^{\mathrm{in}}(i) = 0$, we must have $Y_t \in \mathcal{Y}_i^\perp$. Therefore, $x_t^{\mathrm{in}}(i) = 1$ if and only if $Z_t = z_i$. It then follows from (B.2) that

$$\Delta x_t(i) = \begin{cases} \zeta(1 - x_{t-1}(i)) & \text{if } Z_t = z_i \\ -\zeta x_{t-1}(i) & \text{otherwise} \end{cases} \tag{B.3}$$

Suppose first that $x_{t-1}(i) \in [0,1)$. If $Z_t = z_i$, $x_t(i) - x_{t-1}(i) = \zeta(1 - x_{t-1}(i)) > 0$. If $Z_t \neq z_i$, $x_t(i) - x_{t-1}(i) = -\zeta x_{t-1}(i) \leq 0$, establishing the first statement.

Now suppose $x_{t-1}(i) = 1$. If $Z_t = z_i$, then $x_t(i) = x_{t-1}(i) = 1$. If $Z_t \neq z_i$. $x_t(i) < x_{t-1}(i) = 1$, establishing the second statement. $\square$

We use the notation $\mathbf{E}$ to denote the expectation taken under the econometrician's measure, whereas $\mathbb{E}$ is the expectation taken under the agent's subjective probability.

**Proof of Theorem 7.** We calculate $\mathbf{E}[\Delta x_{t+1}(i) \,|\, \Delta x_t(i) > 0]$, where $\mathbf{E}$ denotes the expectation that the econometrician calculates. We assume throughout that $x_{t-1}(i) < 1$. Then lemma B.3 and the assumptions of the theorem together imply that $\Delta x_t(i) > 0$ if and only if $Z_t = z_i$. We therefore need only calculate $\mathbf{E}[\Delta x_{t+1}(i) \,|\, Z_t = z_i]$.

Given that $Z_t = z_i$, we consider what happens at $t+1$. If $Z_{t+1} \neq z_i$, $f_{t+1} \neq f_t$ by the assumptions of the theorem. We have $y_{t+1} \in \mathcal{Y}_i^\perp$, and $f_{t+1} \in \Omega_{i,t}^\perp$. Thus, $x_{t+1}^{\text{in}}(i) = 0$ by Lemma B.1. If, on the other hand, $Z_{t+1} = z_i$, $y_{t+1} \in \mathcal{Y}_i$, and $f_{t+1} \in \Omega_{i,t-1}$. Recall

$$x_{t+1}^{\text{in}} \propto (M_{t-1} + x_t f_t^\top) f_{t+1}. \tag{B.4}$$

If the features are again novel $(f_{t+1} \neq f_i)$, then $x_{t+1}^{\text{in}}(i) = 1$ by Lemma B.1. However, if $f_{t+1} = f_t$, (B.4) becomes

$$x_{t+1}^{\text{in}} \propto M_{t-1} f_{t+1} + x_t.$$

The first term is proportional to $\hat{e}_i$ by assumption. Its magnitude will depend on the magnitude of the elements in $M_{t-1}$. Let $t_0$ be the length of the prior sample.[59] Lemma A.1, implies

$$x_{t+1}^{\text{in}} \propto (t + t_0 - 1)\hat{e}_i + x_t$$

and therefore that

$$x_{t+1}^{\text{in}}(i) = (\|(t + t_0 - 1)\hat{e}_i + x_t\|)^{-1}(t + t_0 - 1 + x_t(i)). \tag{B.5}$$

To summarize, conditional on $Z_{t+1} = z_i$, we have

$$x_{t+1}^{\text{in}}(i) = \begin{cases} 1 & \text{if } f_{t+1} \neq f_t \\ (\|(t + t_0 - 1)\hat{e}_i + x_t\|)^{-1}(t + t_0 - 1 + x_t(i)) & f_{t+1} = f_t \end{cases}$$

[59]The length of the prior sample is the sum of the elements in $M_0$. See Lemma A.1 for further discussion.

Let
$$\bar{x}_t^{\text{in}}(i; x_t) = \mathbf{E}_t[x_{t+1}^{\text{in}}(i)|Z_t = Z_{t+1} = z_i] \tag{B.6}$$

Note that $x_t(i) < \bar{x}_t(i; x_t) < 1$. A large prior sample or many elements of $\mathcal{Y}_i$ give us $\bar{x}_t(i; x_t)$ close to 1.

Substituting (B.6) into (B.2) implies

$$\mathbf{E}_t[\Delta x_{t+1}(i)|Z_{t+1} = Z_t = z_i] = \begin{cases} \zeta(\bar{x}^{\text{in}}(i; x_t) - x_t(i)) & Z_{t+1} = z_i \\ -\zeta x_t(i) & \text{otherwise} \end{cases}$$

Taking the expectation over the possible outcomes of $Z_{t+1}$, we find:

$$\begin{aligned} \mathbf{E}_t[\Delta x_{t+1}(i)|\Delta x_t(i) > 0] &= \mathbf{E}_t[\Delta x_{t+1}(i)|Z_t = z_i] \\ &= \zeta(p_{ii}^Z \bar{x}^{\text{in}}(i; x_t) - x_t(i)). \end{aligned}$$

The theorem follows. $\qquad\qquad\qquad\qquad\qquad\qquad\qquad\qquad\qquad\qquad\qquad\square$

**Proof of Theorem 8.** Because $Z_t$ is iid, $x_t^{\text{in}}(i) = p(z_i \,|\, Y_t)$ (Theorem 2). Applying (B.2):
$$\Delta x_{t+1}(i) = \zeta(p(z_i|Y_{t+1}) - x_t(i)). \tag{B.7}$$

The assumption of iid $Z_t$ implies that $Y_t$ is also iid. Thus,

$$\mathbf{E}_t p(z_i|Y_{t+1}) = \mathbf{E} p(z_i|Y_{t+1}) = p(z_i). \tag{B.8}$$

We now calculate $\mathbf{E}_t x_t(i)$, conditional on an upward revision in beliefs from $t-1$ to $t$. First note $p(z_i) = \sum_{j=1}^n p(z_i|y_j)p(y_j)$. For $y_j \in \mathcal{Y}_i^\perp$, $p(z_i|y_j) = 0$. Because $p(z_i)$ averages these zero terms with terms $j$ such that $p(z_i|y_j) > 0$,

$$\mathbf{E}\left[p(z_i|Y_t) \,|\, Y_t \notin \mathcal{Y}_i^\perp\right] > p(z_i).$$

If we further add the provision that $Y_t$ is such that $p(z_i|Y_t) > \bar{x} \geq 0$, we weakly increase

the expectation on the left hand side above. That is,

$$
\begin{aligned}
\mathbf{E}\left[p(z_i|Y_t)\,|\,p(z_i|Y_t) > \bar{x}\right] &= \mathbf{E}\left[p(z_i|Y_t)\,|\,Y_t \notin \mathcal{Y}_i^\perp \ \& \ p(z_i|Y_t) > \bar{x}\right] \\
&\geq \mathbf{E}\left[p(z_i|Y_t)\,|\,Y_t \notin \mathcal{Y}_i^\perp\right] \\
&> p(z_i), \quad\quad\quad\quad\quad\quad\quad\quad\quad\quad\text{(B.9)}
\end{aligned}
$$

because $p(z_i|y_j) > \bar{x}$ implies $p(z_i|y_j) > 0$ which implies $y_j \notin \mathcal{Y}_i^\perp$.

By definition, an upward revision in expectations about state $i$ at time $t$ occurs if and only if $x_t(i) > x_{t-1}(i)$, which in turn occurs if and only if $x_t^{\text{in}}(i) > x_{t-1}(i)$. Finally, recall $p(z_i\,|\,Y_t) = x_t^{\text{in}}(i)$. Putting these pieces together:

$$
\begin{aligned}
\mathbf{E}\left[x_t^{\text{in}}(i)|\Delta x_t(i) > 0\right] &= \mathbf{E}\left[p(z_i|Y_t)\,|\,\Delta x_t(i) > 0\right] \\
&= \mathbf{E}\left[p(z_i|Y_t)\,|\,p(z_i|Y_t) > x_{t-1}(i)\right] \\
&> p(z_i), \quad\quad\quad\quad\quad\quad\quad\quad\quad\quad\text{(B.10)}
\end{aligned}
$$

by (B.9). It then follows from (B.7) that

$$
\begin{aligned}
\mathbf{E}\left[\Delta x_{t+1}(i)|\Delta x_t(i) > 0\right] &= \zeta\left(p(z_i) - \mathbf{E}[(1-\zeta)x_{t-1}(i) + \zeta x_t^{\text{in}}(i)|\Delta x_t(i) > 0]\right) \\
&= \zeta\left(p(z_i) - \zeta\mathbf{E}\left[x_t^{\text{in}}(i)|\Delta x_t(i) > 0\right]\right) - (1-\zeta)\mathbf{E}\left[x_{t-1}(i)|\Delta x_t(i) > 0\right]
\end{aligned}
$$

For $\zeta$ sufficiently large, the first term $p(z_i) - \zeta\mathbf{E}\left[x_t^{\text{in}}(i)|\Delta x_t(i) > 0\right]$ is negative by (B.10), whereas the second term is small. $\qquad\square$

## C  Bayesian updating from rare events

Consider a Bayesian agent learning about the probability of a rare event from observations of the event. In the terminology of Section 2, we assume for the purpose of this section that $Z_t$ is iid, and that there are two outcomes of $\{Z_t\}$, and that the set $\mathcal{Y}$ partitions into outcomes possible in one of these states and those possible in the other.

This is of course the same as saying that $Z_t$ is observable.

Let $p$ denote the probability of the rare event. The agent has prior

$$p \sim \text{Beta}(p^*\tau + 1, (1 - p^*)\tau + 1), \tag{C.1}$$

for $p^*, \tau \geq 0$. This prior corresponds to beliefs if the agent had begun with a prior that is uniform on $[0, 1]$ and observed a sample of length $\tau$, of which there were $p^*\tau$ occurrences of the rare event. The density function corresponding to (C.1) is given by

$$f(p) \propto p^{p^*\tau}(1 - p)^{(1-p^*)\tau},$$

where the constant of proportionality does not depend on $p$.

Assume $T$ years of data. For concreteness, we will call the rare event a crisis. Conditional on the probability $p$, the likelihood of exactly $N$ occurrences of the event equals

$$\mathcal{L}(N \text{ crises} \,|\, p) = \binom{T}{N} p^N (1 - p)^{T-N}. \tag{C.2}$$

Therefore the posterior distribution equals

$$\begin{aligned} f(p \,|\, N \text{ crises}) &\propto \mathcal{L}(N \text{ crises}|p) f(p) \\ &\propto p^{N+p^*\tau}(1 - p)^{T+\tau-(N+p^*\tau)} \end{aligned}$$

where once again we have ignored terms that do not depend on $p$. This is proportional to the Beta density, so

$$p \,|\, N \text{ crises} \sim \text{Beta}(N + p^*\tau + 1, T + \tau - (N + p^*\tau) + 1).$$

It follows from properties of the Beta distribution that the posterior mean equals

$$\mathbb{E}[p \,|\, N \text{ crises}] = \frac{N + p^*\tau + 1}{T + \tau + 2}. \tag{C.3}$$

The posterior mean depends on the sample path. Figure 3 shows the average posterior mean, assuming the likelihood (C.2):

$$
\begin{aligned}
\mathbb{E}_{\#\mathrm{crises}}\left[\mathbb{E}[p \mid N\mathrm{crises}]\right] &= \int \frac{N + p^*\tau + 1}{T + 2}\mathcal{L}(N \text{ crises} \mid p)\, dN \\
&= \frac{pT + p^*\tau + 1}{T + \tau + 2}
\end{aligned}
$$

where we have used the fact that, conditional on $p$, $N$ has a binomial distribution, and therefore $\mathbb{E}[N \mid p] = pT$. The figure corresponds to the case of $\tau = 0$, however the results are very similar for $\tau > 0$, provided that the actual sample is large relative to the prior sample.

# D   Proof of the result in Section 2.3

**Proof of Theorem 9.** Following the proof and notation of of Theorem 2, note that for any latent state $i$,

$$
\begin{aligned}
f_{it}^{\mathrm{in}} &= \frac{M_{t-1}^{\top}\hat{e}_i}{\|M_{t-1}^{\top}\hat{e}_i\|} \\
&= \frac{P^{\top}\hat{e}_i}{\|P^{\top}\hat{e}_i\|} \\
&= \left(\sum_{j=1}^{n} p(z_i, y_j)\right)^{-1} \begin{bmatrix} p(z_i, y_1) \\ \vdots \\ p(z_i, y_n) \end{bmatrix}.
\end{aligned}
$$

Note that $\sum_{j=1}^{n} p(z_i, y_j)$ is simply the unconditional probability of $z_i$. Thus, the $j$th entry of $f_{it}^{\mathrm{in}}$ is the conditional probability $p(y_j|z_i)$. the unconditional probability of

$p(y_j)$ is therefore:

$$
\begin{aligned}
p(y_j) &= \sum_{i=1}^{m} f_{it}^{\text{in}} p(z_i) \\
&= \sum_{i=1}^{m} f_{it}^{\text{in}} x_t^{\text{in}}(i)
\end{aligned}
$$

as required by (10).

$\square$

# E  Generalizing results in Section 2.4

Consider the setting of Section 2.4: the agent experiences a crisis at time $t - 1$ (say, 1929), followed by a depression at time $t$. Section 2.4 considers the effect of re-appearance of crisis features at time $t' > t$ under two simplifying assumptions, (1) $x_{1929}$ is a basis vector and (2) $x_{t'}^{\text{in}} = x_{t'}$. In this appendix we relax these assumptions. Given the relative uniqueness of the 1929 stock market crash and the Great Depression we fix ideas by assuming that $f_{\text{crisis}}$ and $f_{\text{depression}}$ were novel events, and that prior associations through $M_0$ were sufficiently weak as to be negligible. We assume that depression features last for $k$ periods. Because our primary interest is the jump-back-in time and not details of associations per se, we assume that $f_{\text{depression}}$ did not reoccurred prior to $t'$, and that the context retrieved by the *original* depression features is not associated with the context in 1929.[60]

The reappearance of crisis features, $f_{t'} = f_{\text{crisis}}$ implies $x_{t'}^{\text{in}} = x_{1929}$ as in Section 2.4. Context at time $t'$ is now:

$$
x_{t'} = (1 - \zeta)x_{t'-1} + \zeta x_{1929}. \tag{E.1}
$$

Assume that $t'$ represents the first appearance of crisis features, so that $x_{t'-1}^{\top} x_{1929} = 0$.

---

[60]This would occur if depression features were novel between $t$ and $t + k - 1$. We disregard the small effect of learning in this initial period.

Let $w_{1929,i}$ be the projection of $x_{1929}$ onto basis vector $i$, so that

$$x_{1929} = \sum_{i=1}^{m} w_{1929,i} \hat{e}_i \qquad w_{1929,i} \geq 0. \tag{E.2}$$

with $\sum_{i=1}^{m} w_{1929,i} = 1$.

From (9), it follows that features retrieved by (basis) context vector $\hat{e}_i$ at time $t'$ equal

$$f_{i,t'}^{\text{in}} = \alpha_{i,t'} M_{t'-1}^{\top} \hat{e}_i, \qquad i = 1, \ldots, m, \tag{E.3}$$

where $\alpha_{i,t'} = ||M_{t'-1}^{\top} \hat{e}_i||^{-1}$ ensures that the elements of $f_{i,t'}^{\text{in}}$ sum to 1. According to (10), we calculate the features retrieved by context at time $t'$ by taking the weighted sum of (E.3). Because context at time $t'$ is itself a weighted sum of prior context and retrieved context, we can consider each of these terms separately.

Define the notation

$$f_{1929,t'}^{\text{in}} \equiv \sum_{i=1}^{m} w_{1929,i} f_{i,t'}^{\text{in}}.$$

By (10), $f_{1929,t'}^{\text{in}}$ are the features retrieved by $x_{t'}^{\text{in}} = x_{1929}$. The subjective probability of a depression implied by $f_{1929,t'}^{\text{in}}$ equals the inner product $(f_{\text{depression}})^{\top} f_{1929,t'}^{\text{in}}$. Note that $f_{\text{depression}}$ is the basis vector representing physical depression features, so the inner product is simply the entry of $f_{1929,t'}^{\text{in}}$ that corresponds to depressions among all the elements of $\mathcal{Y}$. The inner product equals

$$(f_{\text{depression}})^{\top} f_{1929,t'}^{\text{in}} = \sum_{i=1}^{m} w_{1929,i} (f_{\text{depression}})^{\top} f_{i,t'}^{\text{in}}.$$

Using (E.3):

$$
\begin{aligned}
(f_{\text{depression}})^\top f_{i,t'}^{\text{in}} &= \alpha_i (f_{\text{depression}})^\top M_{t'-1}^\top \hat{e}_i && \text{(E.4)} \\
&= \alpha_i (f_{\text{depression}})^\top \left( M_0^\top \hat{e}_i + \sum_{s=1}^{t'-1} f_s x_s^\top \right) \hat{e}_i && \text{(E.5)} \\
&= \alpha_i \sum_{s=1}^{t'-1} (f_{\text{depression}})^\top f_s)(x_s^\top \hat{e}_i) && \text{(E.6)} \\
&= \alpha_i (x_t + \cdots + x_{t+k-1})^\top \hat{e}_i && \text{(E.7)} \\
&= \alpha_i (\sum_{l=1}^{k} (1-\zeta)^l x_{1929} + x_{1929}^\perp)^\top \hat{e}_i, && \text{(E.8)}
\end{aligned}
$$

where $x_{1929}^\perp$ is a vector orthogonal to $x_{1929}$. Equation E.6 follows from the lack of prior associations for depression features through $M_0$. Equation E.7 follows from the assumption of no depression features after time $t + k - 1$, and (E.8) follows from the orthogonality of retrieved context and $x_{1929}$ between $t$ and $t + k - 1$. Because $x_{1929}^\top \hat{e}_i = w_{1929,i}$,

$$
(f_{\text{depression}})^\top f_{i,t'}^{\text{in}} = \left( \frac{1}{\zeta} - 1 \right) \left( 1 - (1-\zeta)^k \right) \sum_{i=1}^{m} \alpha_i w_{1929,i}^2 \tag{E.9}
$$

In the case of basis $x_{1929}$, $w_{1929,i}$ equals 1 for exactly one $i$ and is otherwise 0. The weight $\alpha_I$ is determined by how common the context $x_{1929}$ is. If uncommon, then $\alpha_i$ simply equals $\left( \frac{1}{\zeta} - 1 \right) \left( 1 - (1-\zeta)^k \right)$ and $x_{1929}$ retrieves a probability of 1.

Note however, that total context is not $x_{t'}^{\text{in}} = x_{1929}$ but rather (E.1). Under the assumption of orthogonality, the probability of a depression goes from zero (features retrieved by $x_{t'-1}$ to $\zeta$ multiplied by (E.9).

# F  Proofs for Section 4.1

This Appendix contains proofs generalizing the results in Section 4.1 to $\zeta < 1$. The following is a simple extension of (5) to retrieved features.

**Lemma F.1.** *Features retrieved by context $\hat{e}_i$ as of time $t$ satisfy:*

$$f_{t,i}^{\text{in}} = \left( \iota^\top M_0^\top \hat{e}_i + \sum_{s=1}^{t-1} (x_s^\top \hat{e}_i) \right)^{-1} \left( M_0^\top \hat{e}_i + \sum_{s=1}^{t-1} f_s (x_s^\top \hat{e}_i) \right) \tag{F.1}$$

*where $f_s$ can be physical or retrieved features, and where $\iota$ is the $n \times 1$ vector of ones.*

**Proof.** Combining (9) with (4) implies

$$f_{t,i}^{\text{in}} \quad \propto \quad M_0^\top \hat{e}_i + \sum_{s=1}^{t} (f_s x_s^\top) \hat{e}_i$$

$$\propto \quad M_0^\top \hat{e}_i + \sum_{s=1}^{t} f_s (x_s^\top \hat{e}_i)$$

which is analogous to (5) for retrieved features. To construct the normalizing constant, pre-multiply (F.2) with $\iota$, and recall that $\iota^\top f_s = 1$ for all $s$. $\qquad\qquad\square$

The following is an extension of Theorem 4 to features retrieval.

**Theorem F.2.** *Let $k_i = \|M_0^\top \hat{e}_i\|$. Assume retrieved features are encoded with context. Then retrieved features obey the following recursion*

$$f_{t,i}^{\text{in}} = \left( 1 - \frac{(x_{t-1}^\top \hat{e}_i)(1 - x_{t-1}^\top \hat{e}_i)}{k_i + \sum_{s=1}^{t-1} x_s^\top \hat{e}_i} \right) f_{t-1,i}^{\text{in}} + \frac{(x_{t-1}^\top \hat{e}_i)(1 - x_{t-1}^\top \hat{e}_i)}{k_i + \sum_{s=1}^{t-1} x_s^\top \hat{e}_i} f_{t-1,i}^{\text{in},\perp}, \tag{F.2}$$

*where*

$$f_{t-1,i}^{\text{in},\perp} = (1 - x_{t-1}^\top \hat{e}_i)^{-1} \sum_{j \neq i} x_{t-1}(j) f_{t-1,j}^{\text{in}}, \tag{F.3}$$

*namely $f_{t-1,i}^{\text{in},\perp}$ represents features retrieved at $t-1$ by context elements other than $i$. The initial condition is $f_{0,i} \propto M_0^\top \hat{e}_i$.*

**Proof.** We apply (F.1), setting $k_i = \iota^\top M_0^\top \hat{e}_i$ and $f_s = f_s^{\text{in}}$:

$$f_{t,i}^{\text{in}} = \left( k_i + \sum_{s=1}^{t-1} (x_s^\top \hat{e}_i) \right)^{-1} \left( M_0^\top \hat{e}_i + \sum_{s=1}^{t-1} f_s^{\text{in}} (x_s^\top \hat{e}_i) \right) \tag{F.4}$$

78

We rewrite (F.4), using the same recursive reasoning as in the proof of Theorem 4:

$$
f_{t,i}^{\text{in}} = \left(k_i + \sum_{s=1}^{t-1}(x_s^\top \hat{e}_i)\right)^{-1} \underbrace{\left(M_0^\top \hat{e}_i + \sum_{s=1}^{t-2} f_s^{\text{in}}(x_s^\top \hat{e}_i)\right)}_{(1+\sum_{s=1}^{t-2} x_s^\top \hat{e}_i)f_{t-1,i}^{\text{in}}}
$$

$$
+ \left(k_i + \sum_{s=1}^{t-1} x_s^\top \hat{e}_i\right)^{-1} f_{t-1}^{\text{in}}(x_{t-1}^\top \hat{e}_i). \quad \text{(F.5)}
$$

Note that we apply (F.1) at $t-1$ to conclude $M_0^\top \hat{e}_i + \sum_{s=1}^{t-2} f_s^{\text{in}}(x_s^\top \hat{e}_i) = (1+\sum_{s=1}^{t-2} x_s^\top \hat{e}_i)f_{t-1,i}^{\text{in}}$.

Retrieved features at time $t-1$ are a weighted average of those retrieved by $\hat{e}_i$, and those retrieved by the other elements of context.

$$
f_{t-1}^{\text{in}} = f_{t-1,i}^{\text{in}}(x_{t-1}^\top \hat{e}_i) + f_{t-1,i}^{\text{in},\perp}(1 - x_{t-1}^\top \hat{e}_i), \quad \text{(F.6)}
$$

where $f_{t-1,i}^{\text{in},\perp}$ is as defined in (F.3). Combining (F.5) and (F.6) shows that $f_{t,i}^{\text{in}}$ is a weighted average of $f_{t-1,i}^{\text{in}}$ and $f_{t-1,i}^{\text{in},\perp}$. Moreover, the coefficient multiplying $f_{t-1,i}^{\text{in},\perp}$ must equal $\left(k_i + \sum_{s=1}^{t-1} x_s^\top \hat{e}_i\right)^{-1}(1 - x_{t-1}^\top \hat{e}_i)(x_{t-1}^\top \hat{e}_i)$. Because the elements of $f_{t,i}^{\text{in}}$ must sum to 1, it follows that the coefficient on $f_{t-1,i}^{\text{in}}$ equals one minus this quantity, as shown in (F.2). $\qquad \square$

The following Lemma generalizes Lemma A.3 to encoding of retrieved features. We assume that only the physical event triggers encoding of features that are non-orthogonal to the event. For example, if the event is a stock market loss, we disregard events that were not losses, but nonetheless reminded the agent of losses, as second-order.

**Lemma F.3.** *Assume the agent experiences an event at $\{t_1, \ldots, t_\ell\}$, and that features vectors are otherwise orthogonal to the event. Then retrieved context in response to $f_{t_\ell} = e_i$ is proportional to*

$$
x_{t_\ell}^{\text{in}} \propto M_0 e_i + \sum_{k=2}^{\ell-1} x_{t_k}(f_{t_k}^\top e_i) \quad \text{(F.7)}
$$

Note that when features are basis vectors, (F.7) reduces to (A.8).

**Proof of Lemma F.3.** By (5) and the fact that $t_1$ is the first occurence of the event:

$$x_{t_1}^{\text{in}} \propto M_0 e_i.$$

Moreover,

$$M_{t_1-1} e_i = M_0 e_i$$

Assume by induction that (F.7) holds for the $(\ell-1)$st occurrence of the event:

$$x_{t_{\ell-1}}^{\text{in}} \propto M_0 e_i + \sum_{k=2}^{t_{\ell-2}} x_{t_k} (f_{t_k}^\top e_i). \tag{F.8}$$

and that

$$M_{t_{\ell-1}-1} e_i = M_0 e_i + \sum_{k=2}^{t_{\ell-2}} x_{t_k} (f_{t_k}^\top e_i). \tag{F.9}$$

By (2), context equals

$$x_{t_{\ell-1}} = (1-\zeta) x_{t_{\ell-1}-1} + \zeta x_{t_{\ell-1}}^{\text{in}}$$

and memory is updated as:

$$M_{t_{\ell-1}} = M_{t_{\ell-1}-1} + x_{t_{\ell-1}} f_{t_{\ell-1}}^\top. \tag{F.10}$$

where it is not necessary to take a stance on whether $f_{t_{\ell-1}} = e_i$ or features retrieved by $x_{t_{\ell-1}}$.

By (3),

$$x_{t_\ell}^{\text{in}} \propto M_{t_\ell-1} e_i.$$

By definition, $t_{\ell-1}$ is the occurrence of the event just before the occurence at $t_\ell$. Because of this and additional assumptions of the theorem, all features between $t_{\ell-1}$ and $t_\ell$ are

80

orthogonal to $e_i$. Therefore we can ignore terms in $M_{t_\ell - 1}$ that occur after $t_{\ell-1}$, and

$$x_{t_\ell}^{\text{in}} \propto M_{t_{\ell-1}} e_i$$

Substituting in from (F.10):

$$x_{t_\ell}^{\text{in}} \propto (M_{t_{\ell-1}-1} + x_{t_{\ell-1}} f_{t_{\ell-1}}^\top) e_i.$$

Substituting in from (F.9):

$$x_{t_\ell}^{\text{in}} \propto M_0 e_i + \sum_{k=2}^{t_{\ell-2}} x_{t_k} (f_{t_k}^\top e_i) + x_{t_{\ell-1}} (f_{t_{\ell-1}}^\top e_i).$$

Because we can ignore terms in $M_{t_\ell - 1}$ that occur after $t_{\ell-1}$:

$$M_{t_\ell - 1} e_i = M_{t_{\ell-1}} e_i$$

It follows from (F.9) and (4) that

$$M_{t_{\ell-1}} e_i = M_0 e_i + \sum_{k=2}^{t_{\ell-2}} x_{t_k} (f_{t_k}^\top e_i) + x_{t_{\ell-1}} (f_{t_{\ell-1}}^\top e_i),$$

completing the proof. □

# References

Abbott, L. F. and Blum, K. I. (1996). Functional significance of long-term potentiation for sequence learning and prediction. *Cerebral Cortex*, 6(3):406–416.

Allen, F. and Gale, D. (2009). *Understanding Financial Crises*. Clarendon lectures in finance. Oxford University Press.

Angeletos, G.-M., Huo, Z., and Sastry, K. A. (2020). *Imperfect Macroeconomic Expectations: Evidence and Theory*. NBER Macroeconomic Annual. University of Chicago Press.

Arrow, K. J. (1971). *Essays in the Theory of Risk-Bearing*. Markham Publishing Co, Chicago.

Azeredo da Silveira, R. and Woodford, M. (2019). Noisy memory and over-reaction to news. *AEA Papers and Proceedings*, 109:557–61.

Barberis, N., Greenwood, R., Jin, L., and Shleifer, A. (2015). X-CAPM: An extrapolative capital asset pricing model. *Journal of Financial Economics*, 115(1):1–24.

Barberis, N., Shleifer, A., and Vishny, R. (1998). A model of investor sentiment. *Journal of Financial Economics*, page 37.

Barberis, N. C. (2013). Thirty years of prospect theory in economics: A review and assessment. *Journal of Economic Perspectives*, 27(1):173–196.

Bassi, A., Colacito, R., and Fulghieri, P. (2013). 'o sole mio: An experimental analysis of weather and risk attitudes in financial decisions. *The Review of Financial Studies*, 26(7):1824–1852.

Bordalo, P., Coffman, K., Gennaioli, N., Schwerter, F., and Shleifer, A. (2021). Memory and representativeness. *Psychological Review*, 128(1):71–85.

Bordalo, P., Gennaioli, N., Porta, R. L., and Shleifer, A. (2020a). Expectations of fundamentals and stock market puzzles. Working Paper 27283, National Bureau of Economic Research.

Bordalo, P., Gennaioli, N., and Shleifer, A. (2018). Diagnostic expectations and credit cycles. *The Journal of Finance*, 73(1):199–227.

Bordalo, P., Gennaioli, N., and Shleifer, A. (2020b). Memory, attention, and choice. *The Quarterly Journal of Economics*, 135(3):1399–1442.

Bower, G. H. (1972). Stimulus-sampling theory of encoding variability. In Melton, A. W. and Martin, E., editors, *Coding Processes in Human Memory*, chapter 5, pages 85–121. John Wiley and Sons, New York.

Brady, T. F., Konkle, T., and Alvarez, G. A. (2011). A review of visual memory capacity: Beyond individual items and toward structured representations. *Journal of vision*, 11(5):4–4.

Brown, G. D. A., Neath, I., and Chater, N. (2007). A temporal ratio model of memory. *Psychological Review*, 114(3):539–576.

Brown, G. D. A., Preece, T., and Hulme, C. (2000). Oscillator-based memory for serial order. *Psychological Review*, 107(1):127–181.

Brunnermeier, M. K. and Sannikov, Y. (2014). A macroeconomic model with a financial sector. *American Economic Review*, 104(2):379–421.

Burgess, N. and Hitch, G. J. (2006). A revised model of short-term memory and long-term learning of verbal sequences. *Journal of Memory and Language*, 55:627–652.

Caballero, R. J. and Simsek, A. (2020). A risk-centric model of demand recessions and speculation*. *The Quarterly Journal of Economics*, 135(3):1493–1566.

Campbell, J. Y. (2016). Restoring rational choice: The challenge of consumer financial regulation. *American Economic Review*, 106(5):1–30.

Carr, H. A. (1931). The laws of association. *Psychological Review*, 38:212–228.

Cohen, R. T. and Kahana, M. J. (2020). Retrieved-context theory of memory in emotional disorders. Working paper, University of Pennsylvania.

Cohn, A., Engelmann, J., Fehr, E., and Marechal, M. A. (2015). Evidence for countercyclical risk aversion: An experiment with financial professionals. *The American Economic Review*, 105(2):860–885.

Conlin, M., O'Donoghue, T., and Vogelsang, T. J. (2007). Projection bias in catalog orders. *The American Economic Review*, 97(4):1217–1249.

Crowder, R. G. (1976). *Principles of learning and memory*. Lawrence Erlbaum and Associates, Hillsdale, NJ.

Cuculiza, C., Antoniou, C., Kumar, A., and Maligkris, A. (2020). Terrorist attacks, analyst sentiment, and earnings forecasts. forthcoming, *Management Science*.

Daniel, K., Hirshleifer, D., and Subrahmanyam, A. (1998). Investor psychology and security market under- and overreactions. *The Journal of Finance*, 53(6):1839–1885.

Ebbinghaus, H. (1885/1913). *Memory: a contribution to experimental psychology.* Teachers College, Columbia University, New York.

Eich, E. (1995). Searching for mood dependent memory. *Psychological Science*, 6(2):67–75.

Ekstrom, A. D., Kahana, M. J., Caplan, J. B., Fields, T. A., Isham, E. A., Newman, E. L., and Fried, I. (2003). Cellular networks underlying human spatial navigation. *Nature*, 425:184–187.

Enke, B., Schwerter, F., and Zimmermann, F. (2020). Associative memory and belief formation. working paper, Harvard University.

Estes, W. K. (1955). Statistical theory of spontaneous recovery and regression. *Psychological Review*, 62:145–154.

Estes, W. K. (1959). *Component and pattern models with Markovian interpretations. Studies in Mathematical Learning Theory*. Stanford University Press, Stanford, CA.

Estes, W. K. (1986). Array models for category learning. *Cognitive Psychology*, 18(4):500–549.

Folkerts, S., Rutishauser, U., and Howard, M. (2018). Human episodic memory retrieval is accompanied by a neural contiguity effect. *Journal of Neuroscience*, 38(17):4200–4211.

French, K. R., Baily, M. N., Campbell, J. Y., Cochrane, J. H., Diamond, D. W., Duffie, D., Kashyap, A. K., Mishkin, F. S., Rajan, R. G., Scharfstein, D. S., Shiller, R. J., Shin, H. S., Slaughter, M. J., Stein, J. C., and Stulz, R. M. (2010). *The Squam Lake Report: Fixing the Financial System*. Princeton University Press.

Fuster, A., Laibson, D., and Mendel, B. (2010). Natural expectations and macroeconomic fluctuations. *Journal of Economic Perspectives*, 24(4):67–84.

Gabaix, X. (2019). Behavioral inattention. In Bernheim, D., DellaVigna, S., and Laibson, D., editors, *Handbook of Behavioral Economics*, volume 2. Elsevier.

Gallistel, C. R. and King, A. P. (2009). *Memory and the Computational Brain: Why Cognitive Science Will Transform Neuroscience*. Wiley-Blackwell.

Gelman, A., Carlin, J. B., Stern, H. S., and Rubin, D. B. (2004). *Bayesian Data Analysis*. Chapman & Hall/CRC, Boca Raton, FL.

Gennaioli, N. and Shleifer, A. (2018). *A Crisis of Beliefs: Investor Psychology and Financial Fragility*. Princeton University Press, Princeton, NJ.

Gilboa, I. and Schmeidler, D. (1995). Case-based decision theory. *The Quarterly Journal of Economics*, 110(3):605–639.

Godden, D. R. and Baddeley, A. D. (1975). Context-dependent memory in two natural environments: On land and under water. *British Journal of Psychology*, 66:325–331.

Gödker, K., Jiao, P., and Smeets, P. (2019). Investor Memory. Working paper, Maastrich University.

Goetzmann, W. N., Kim, D., and Shiller, R. J. (2017). Crash beliefs from investor surveys. NBER Working Paper No. 22143.

Gomes, J. F., Grotteria, M., and Wachter, J. A. (2019). Cyclical dispersion in expected defaults. *Review of Financial Studies*, 32(4):1275–1308.

Gorton, G. and Metrick, A. (2012). Securitized banking and the run on repo. *Journal of Financial Economics*, 104(3):425–451.

Gourio, F. (2012). Disaster risk and business cycles. *American Economic Review*, 102(6):2734–2766.

Greene, R. L. (1992). *Human memory: Paradigms and paradoxes*. Lawrence Erlbaum and Associates, Hillsdale, New Jersey.

Grossman, S. J. and Stiglitz, J. E. (1980). On the impossibility of informationally efficient markets. *American Economic Review*, 70(3):393–408.

Guiso, L., Sapienza, P., and Zingales, L. (2018). Time varying risk aversion. *Journal of Financial Economics*, 128(3):403–421.

Hamilton, J. D. (1994). *Time Series Analysis*. Oxford University Press, Princeton, NJ.

He, Z. and Krishnamurthy, A. (2013). Intermediary asset pricing. *American Economic Review*, 103(2):732–70.

Healey, M. K., Long, N. M., and Kahana, M. J. (2019). Contiguity in episodic memory. *Psychonomic Bulletin & Review*, 26(3):699—720.

Henson, R., Norris, D. G., Page, M. P. A., and Baddeley, A. D. (1996). Unchained memory: Error patterns rule out chaining models of immediate serial recall. *The Quarterly Journal of Experimental Psychology*, 49A:80–115.

Herbart, J. F. (1834). Lehrbuch zur Psychologie. In Kehrbach, K., editor, *Sämtliche Werke*, volume 4, pages 376–377. August Wilhelm Unzer, Königsberg, 2nd edition.

Hinrichs, J. V. and Buschke, H. (1968). Judgment of recency under steady-state conditions. *Journal of experimental psychology*, 78(4p1):574.

Hintzman, D. L. (1988). Judgments of frequency and recognition memory in multiple-trace memory model. *Psychological Review*, 95:528–551.

Hong, H. and Stein, J. C. (1999). A unified theory of underreaction, momentum trading, and overreaction in asset markets. *The Journal of Finance*, 54(6):2143–2184.

Howard, M. W., Jing, B., Rao, V. A., Provyn, J. P., and Datey, A. V. (2009). Bridging the gap: Transitive associations between items presented in similar temporal contexts. *Journal of Experimental Psychology: Learning, Memory, and Cognition*, 35:391–407.

Howard, M. W. and Kahana, M. J. (1999). Contextual variability and serial position effects in free recall. *Journal of Experimental Psychology: Learning, Memory, and Cognition*, 25(4):923–941.

Howard, M. W. and Kahana, M. J. (2002). A distributed representation of temporal context. *Journal of Mathematical Psychology*, 46(3):269–299.

Howard, M. W., Viskontas, I. V., Shankar, K. H., and Fried, I. (2012). Ensembles of human MTL neurons "jump back in time" in response to a repeated stimulus. *Hippocampus*, 22:1833–1847.

Hume, D. (1748). *An Enquiry Concerning Human Understanding*. Retrieved from Project Gutenberg: http://www.gutenberg.org/files/9662/9662-h/9662-h.htm.

James, W. (1890). *The principles of psychology*. Henry Holt and Co, Inc., New York, NY, US.

Jegadeesh, N. and Titman, S. (1993). Returns to buying winners and selling losers: Implications for stock market efficiency. *Journal of Finance*, 48:65–91.

Jones, C. M. (2002). A century of stock market liquidity and trading costs. Working paper, Columbia University.

Jost, A. (1897). Die Assoziationsfestigkeit in ihrer Abhängigkeit von der Verteilung der Wiederholungen. *Zeitschrift für Psychologie und Physiologie der Sinnesorgane*, 16:436–472.

Kahana, M. J. (1996). Associative retrieval processes in free recall. *Memory & Cognition*, 24(1):103–109.

Kahana, M. J. (2012). *Foundations of Human Memory*. Oxford University Press, New York, NY.

Kahana, M. J. and Jacobs, J. (2000). Inter-response times in serial recall: Effects of intraserial repetition. *Journal of Experimental Psychology: Learning, Memory and Cognition*, 26:1188–1197.

Kahle, K. M. and Stulz, R. (2013). Access to capital, investment, and the financial crisis. *Journal of Financial Economics*, 110(2):280–299.

Keim, D. B. (1983). Size-related anomalies and stock return seasonality: Further empirical evidence. *Journal of Financial Economics*, 12(1):13–32.

Keim, D. B. and Stambaugh, R. F. (1984). A further investigation of the weekend effect in stock returns. *The Journal of Finance*, 39(3):819–835.

Kempter, R., Gerstner, W., and van Hemmen, J. L. (1999). Hebbian learning and spiking neurons. *Phys. Rev. E*, 59:4498–4514.

Kolers, P. A. and Magee, L. E. (1978). Specificity of pattern-analyzing skills in reading. *Canadian Journal of Psychology*, 32:43–51.

Kuhn, J. R., Lohnas, L. J., and Kahana, M. J. (2018). A spacing account of negative recency in final free recall. *Journal of Experimental Psychology: Learning, Memory, and Cognition*, 44(8):1180–1185.

Ladd, G. T. and Woodworth, R. S. (1911). *Elements of physiological psychology: A treatise of the activities and nature of the mind from the physical and experimental point of view*. Charles Scribner's Sons, New York, NY.

Lashley, K. (1951). The problem of serial order in behavior. In *Cerebral Mechanisms in Behavior*. Wiley, New York.

Loewenstein, G. (2000). Emotions in economic theory and economic behavior. *The American Economic Review*, 90(2):426–432.

Lohnas, L. J. and Kahana, M. J. (2014a). Compound cuing in free recall. *Journal of Experimental Psychology: Learning, Memory and Cognition*, 40(1):12–24.

Lohnas, L. J. and Kahana, M. J. (2014b). A retrieved context account of spacing and repetition effects in free recall. *Journal of Experimental Psychology: Learning Memory and Cognition*, 40(3):755–764.

Lohnas, L. J., Polyn, S. M., and Kahana, M. (2015). Expanding the scope of memory search: Intralist and interlist effects in free recall. *Psychological Review*, 122(2):337–363.

Long, N. M., Danoff, M. S., and Kahana, M. J. (2015). Recall dynamics reveal the retrieval of emotional context. *Psychonomic Bulletin and Review*, 22(5):1328–1333.

Longstaff, F. A. and Piazzesi, M. (2004). Corporate earnings and the equity premium. *Journal of Financial Economics*, 74:401–421.

Lucas, R. E. (1975). An equilibrium model of the business cycle. *Journal of Political Economy*, 83(6):1113–1144.

Lucas, R. E. (1978). Asset prices in an exchange economy. *Econometrica*, 46:1429–1445.

MacDonald, C., Lepage, K., Eden, U., and Eichenbaum, H. (2011). Hippocampal "time cells" bridge the gap in memory for discontiguous events. *Neuron*, 71(4):737–749.

Makino, H. and Komiyama, T. (2015). Learning enhances the relative impact of top-down processing in the visual cortex. *Nature neuroscience*, 18(8):1116–1122.

Malmendier, U. and Nagel, S. (2011). Depression babies: Do macroeconomic experiences affect risk taking? *The Quarterly Journal of Economics*, 126(1):373–416.

Malmendier, U. and Nagel, S. (2016). Learning from inflation experiences. *The Quarterly Journal of Economics*, 131(1):53–87.

Malmendier, U., Nagel, S., and Yan, Z. (2017). The making of hawks and doves: Inflation experiences on the FOMC. NBER Working Paper No. 23228.

Malmendier, U. and Shen, L. S. (2018). Scarred consumption. NBER Working Paper No. 24696.

Manning, J., Hulbert, J., Williams, J., Piloto, L., Sahakyan, L., and Norman, K. (2016). A neural signature of contextually mediated intentional forgetting. *Psychonomic Bulletin & Review*, 23:1534–1542.

Manning, J. R., Polyn, S. M., Baltuch, G., Litt, B., and Kahana, M. J. (2011). Oscillatory patterns in temporal lobe reveal context reinstatement during memory search. *Proceedings of the National Academy of Sciences, USA*, 108(31):12893–12897.

Markowitz, H. (1952). Portfolio selection. *The Journal of Finance*, 7(1):77–91.

Matt, G. E., Vázquez, C., and Campbell, W. K. (1992). Mood-congruent recall of affectively toned stimuli: A meta-analytic review. *Clinical Psychology Review*, 12(2):227–255.

McGeoch, J. A. (1932). Forgetting and the law of disuse. *Psychological Review*, 39:352–70.

Miller, J. F., Neufang, M., Solway, A., Brandt, A., Trippel, M., Mader, I., Hefft, S., Merkow, M., Polyn, S. M., Jacobs, J., Kahana, M. J., and Schulze-Bonhage, A. (2013). Neural activity in human hippocampal formation reveals the spatial context of retrieved memories. *Science*, 342(6162):1111–1114.

Miller, J. F., Weidemann, C. T., and Kahana, M. J. (2012). Recall termination in free recall. *Memory & Cognition*, 40:540–550.

Moreton, B. J. and Ward, G. (2010). Time scale similarity and long-term memory for autobiographical events. *Psychonomic Bulletin & Review*, 17(4):510–515.

Mullainathan, S. (2002). A memory-based model of bounded rationality. *The Quarterly Journal of Economics*, 117(3):735–774.

Müller, G. E. and Pilzecker, A. (1900). Experimental contributions to memory theory. *Zeitschrift für Psychologie Eganzungsband*, 1:1–300.

Müller, G. E. and Schumann, F. (1894). Experimentelle Beiträge zur Untersuchung des Gedächtnisses. *Zeitschrift für Psychologie und Physiologie der Sinnesorgane*, 6(80-191):257–339.

Murdock, B. B. (1974). *Human memory: Theory and data*. Lawrence Erlbaum and Associates, Potomac, MD.

Murdock, B. B. (1985). An analysis of the strength-latency relationship. *Memory & Cognition*, 13:511–521.

Nagel, S. and Xu, Z. (2018). Asset pricing with fading memory. Working paper, University of Chicago and University of Michigan.

Nosofsky, R. M. (1992). Exemplar-based approach to relating categorization, identification, and recognition. In Ashby, F. G., editor, *Multidimensional models of perception and cognition*, pages 363–394. Lawrence Erlbaum and Associates, Hillsdale, New Jersey.

O'Keefe, J. and Dostrovsky, J. (1971). The hippocampus as a spatial map: Preliminary evidence from unit activity in the freely-moving rat. *Brain Research*, 34:171–175.

Pastalkova, E., Itskov, V., Amarasingham, A., and Buzsáki, G. (2008). Internally generated cell assembly sequences in the rat hippocampus. *Science*, 321:1322 – 1327.

Polyn, S. M., Norman, K. A., and Kahana, M. J. (2009). A context maintenance and retrieval model of organizational processes in free recall. *Psychological Review*, 116(1):129–156.

Ramadorai, T., Anagol, S., and Balasubramaniam, V. (2020). Learning from noise: evidence from india's ipo lotteries. forthcoming, *Journal of Financial Economics*.

Reynolds, J. H. and Chelazzi, L. (2004). Attentional modulation of visual processing. *Annual Review of Neuroscience*, 27(1):611–647.

Rubin, D. C. and Berntsen, D. (2009). The frequency of voluntary and involuntary autobiographical memories across the lifespan. *Memory & Cognition*, 37(5):679–688.

Rubin, D. C., Berntsen, D., and Johansen, M. K. (2008). A memory-based model of posttraumatic stress disorder: Evaluating basic assumptions underlying the PTSD diagnosis. *Psychological Review*, 115(4):985–1011.

Samuelson, P. A. (1969). Lifetime portfolio selection by dynamic stochastic programming. *The Review of Economics and Statistics*, 51(3):239–46.

Savage, L. (1954). *The Foundations of Statistics*. Wiley, New York.

Sederberg, P. B., Gershman, S. J., Polyn, S. M., and Norman, K. A. (2011). Human memory reconsolidation can be explained using the temporal context model. *Psychonomic Bulletin & Review*, 18(3):455–468.

Siddiqui, A. P. and Unsworth, N. (2011). Investigating the role of emotion during the search process in free recall. *Memory & Cognition*, 39(8):1387–1400.

Siegel, L. L. . and Kahana, M. J. (2014). A retrieved context account of spacing and repetition effects in free recall. *Journal of Experimental Psychology: Learning Memory and Cognition*, 40(3):755–764.

Sims, C. A. (2003). Implications of rational inattention. *Journal of Monetary Economics*, 50(3):665–690.

Solway, A., Murdock, B. B., and Kahana, M. J. (2012). Positional and temporal clustering in serial order memory. *Memory & Cognition*, 40(2):177–190.

Standing, L. (1973). Learning 10000 pictures. *The Quarterly journal of experimental psychology*, 25(2):207–222.

Teasdale, J. D. and Fogarty, S. J. (1979). Differential effects of induced mood on retrieval of pleasant and unpleasant events from episodic memory. *Journal of Abnormal Psychology*, 88(3):248–257.

Tulving, E. and Madigan, S. A. (1970). Memory and verbal learning. *Annual Review of Psychology*, 21:437–484.

Uitvlugt, M. G. and Healey, M. K. (2019). Temporal proximity links unrelated news events in memory. *Psychological Science*, 30(1):92–104.

Umbach, G., Kantak, P., Jacobs, J., Kahana, M., Pfeiffer, B. E., Sperling, M., and Lega, B. (2020). Time cells in the human hippocampus and entorhinal cortex support

episodic memory. *Proceedings of the National Academy of Sciences of the United States of America.*, 117(45):28463–28474.

Underwood, B. J. (1948). 'Spontaneous recovery' of verbal associations. *Journal of Experimental Psychology*, 38:429–439.

Wachter, J. A. and Kahana, M. J. (2020). Associative learning and representativeness. Working paper, University of Pennsylvania.

Wachter, J. A. and Zhu, Y. (2019). Learning with rare disasters. Working paper, University of Pennsylvania.

Wixted, J. T. (2007). Dual-process theory and signal-detection theory of recognition memory. *Psychological Review*, 114(1):152–176.

Woodford, M. (2020). Modeling imprecision in perception, valuation, and choice. *Annual Review of Economics*, 12(1):579–601.

Woodworth, R. S. (1938). *Experimental Psychology.* H. Holt and Company, New York.

Yaffe, R. B., Kerr, M. S., Damera, S., Sarma, S. V., Inati, S. K., and Zaghloul, K. A. (2014). Reinstatement of distributed cortical oscillations occurs with precise spatiotemporal dynamics during successful memory retrieval. *Proceedings of the National Academy of Sciences*, 111(52):18727–8732.

Yates, F. A. (1966). *The art of memory.* Routledge and Kegan Paul, London, England.

Yntema, D. B. and Trask, F. P. (1963). Recall as a search process. *Journal of Verbal Learning and Verbal Behavior*, 2:65–74.

Zaromb, F. M., Howard, M. W., Dolan, E. D., Sirotin, Y. B., Tully, M., Wingfield, A., and Kahana, M. J. (2006). Temporal associations and prior-list intrusions in free recall. *Journal of Experimental Psychology: Learning, Memory, and Cognition*, 32(4):792–804.
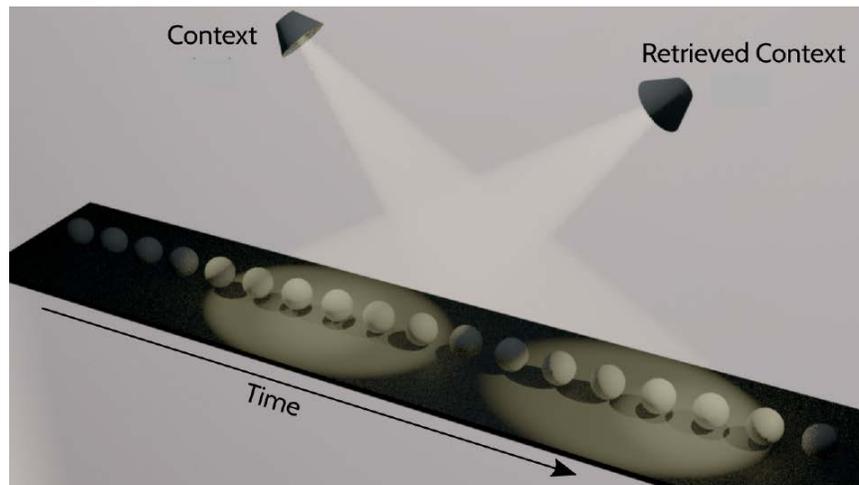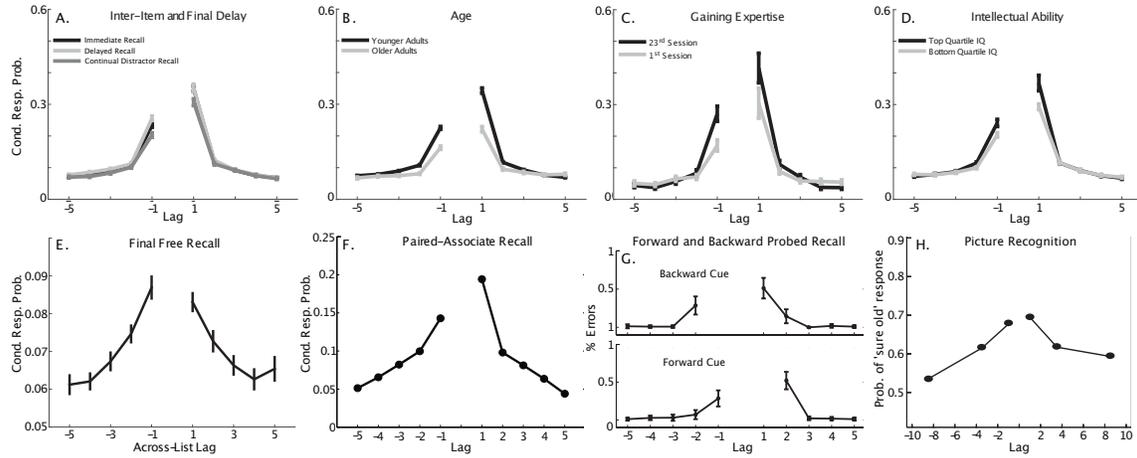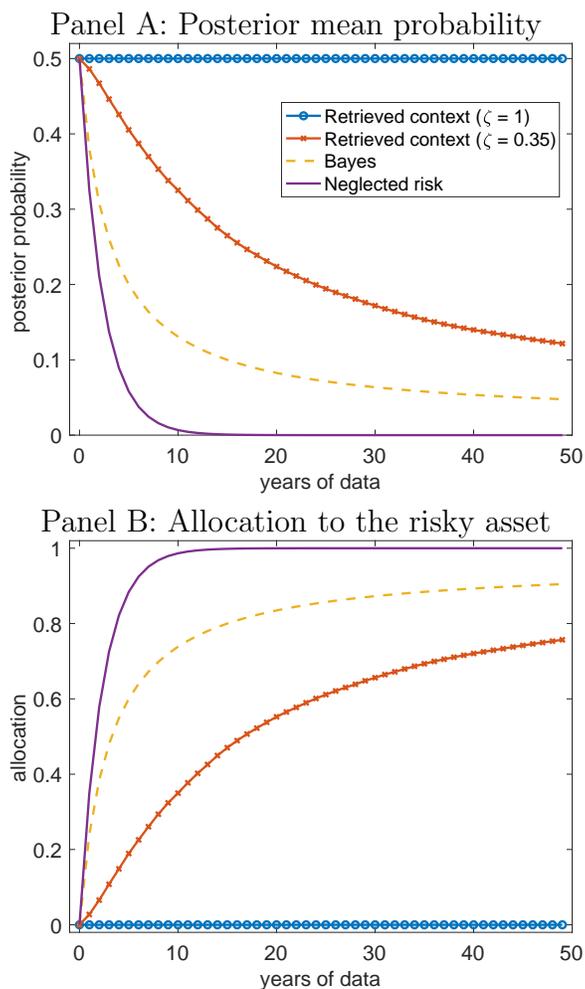
Figure 1: Retrieved Context and the spotlights of memory. In this illustration, memories appear as circles on the stage of life. All experiences that enter memory, as gated by perception and attention, take their place upon the stage. Context serves as a set of spotlights, each shining into memory and illuminating its associated features. The prior state of context illuminates recent memories, whereas the context retrieved by the preceding experience illuminates temporally and semantically contiguous memories. Due to the recursive nature of context and the stochastic nature of retrieval, the lamps can swing over time and illuminate different sets of prior features.

Figure 2: **Universality of Temporal Contiguity. A.** When freely recalling a list of studied items, people tend to successively recall items that appeared in neighboring positions. This temporal contiguity effect (TCE) appears as an in increase in the conditional-response probability as a function of the lag, or distance, between studied items (the lag-CRP). The TCE appears invariant across conditions of immediate recall, delayed recall, and continual-distractor recall, where subjects perform a demanding distractor task between each of the studied items. **B.** Older adults exhibit reduced temporal contiguity, indicating impaired contextual retrieval **C.** Massive practice increases the TCE, as seen in the comparison of 1st and 23rd hour of recall practice. **D.** Higher-IQ subjects exhibit a stronger TCE than individuals with average IQ. **E.** The TCE is not due to inter-item associations as it appears in transitions across different lists, separated by minutes, in a delayed final test given to subjects who studied and recalled many lists. **F.** The TCE appears in conditional error gradients in cued recall, where subjects tend to mistakenly recall items from pairs studied in nearby list positions. **G.** When probed to recall the item that either followed or preceded a cue item, subjects occasionally commit recall errors whose distribution exhibits a TCE both for forward and backward probes. **H.** The TCE also appears when subjects are asked to recognize previously seen travel photos. When successive test items come from nearby positions on the study list, subjects tendency to make high confidence "old" responses exhibits a TCE when the previously tested item was also judged old with high confidence. This effect is not observed for responses made with low confidence. Healey et al. (2019) provide references and descriptions of each experiment.
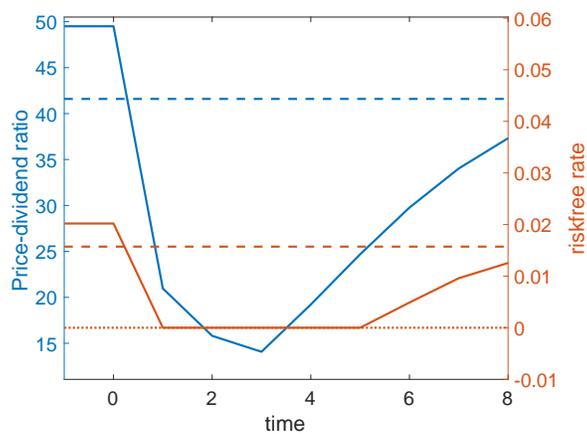
Figure 3: Posterior probability and asset allocation as a function of sample length



Panel A: Posterior mean probability
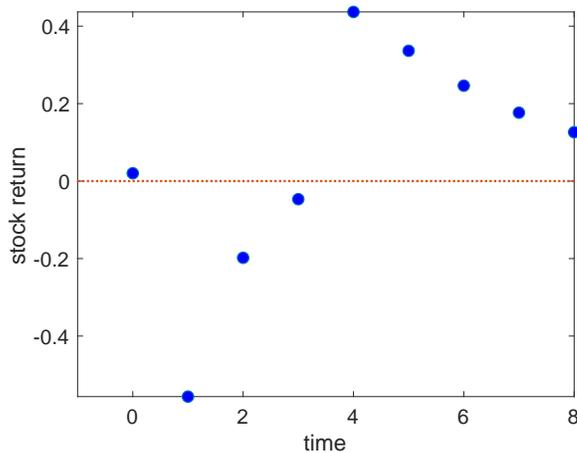
Panel B: Allocation to the risky asset

Notes: The figure shows posterior mean of the probability of a depression (Panel A) and the resulting asset allocation (Panel B) for the model presented in Section 4.1. 'Retrieved context ($\zeta = 1$)' is the retrieved context model when the agent places no weight on prior context. 'Retrieved context ($\zeta = 0.35$)' corresponds to when the agent places a weight of 1-0.35 on the prior context. 'Bayes' corresponds to the mean posterior probability when a Bayesian learns about a rare event from occurrences. 'Neglected risk' corresponds to exponential decay of beliefs after observing a rare event, as described in Section 2.2.

97

Figure 4: Response of prices and returns to a financial crisis

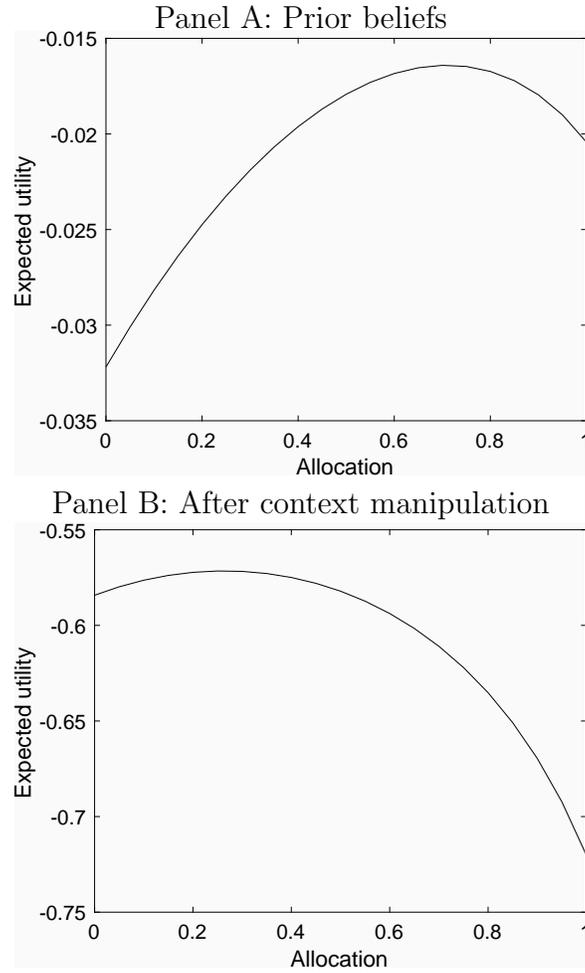Panel A: Price-dividend ratio and riskfree rate



Panel B: Realized equity returns



Notes: Panel A shows the time path of the price-dividend ratio and of the short-term interest rate in response to a period of calm, followed by a financial crisis (in the model of Section 4.2). The dashed lines show full-information values. Panel B shows realized stock returns. The agent observes three periods of crisis features followed by normal features. The figure shows the maximum of the model-implied riskfree rate and zero. Returns are in annual terms.

Figure 5: Expected utility under context manipulation



Panel A: Prior beliefs

Panel B: After context manipulation

Notes: This figure shows expected utility as a function of allocation to the risky asset under the model of Section 4.3. Panel A shows utility prior viewing a scene from a horror movie. Panel B shows utility after context has been manipulated by viewing the scene. In Panel B, curvature of the utility function has increased, so that the optimal allocation to the risky asset falls.