# Mutual Fund Flows and Performance in (Imperfectly) Rational Markets?*

Nikolai Roussanov†, Hongxun Ruan‡ and Yanhao Wei §

October 26, 2020

## Abstract

Does the observed relationship between mutual fund flows and recent performance represent irrational "return chasing" or rational learning about unobserved fund manager skill in the presence of decreasing returns to scale? We estimate a structural model of investor beliefs implicit in the fund flows and compare it with the rational Bayesian benchmark that estimated from past performance data. Our estimates imply that investors are more optimistic about fund manager's average skill level than warranted by the historical data. They over-weight recent performance in a manner consistent with models based on the representativeness heuristic, yet respond slowly to changes in these beliefs, consistent with limited attention and/or informational frictions. Flows to retail funds imply more strongly biased beliefs than those to institutional funds.

*Keywords*: mutual fund performance, flow-performance relationship, Bayesian learning, representativeness

# 1 Introduction

Investor flows in and out of mutual funds respond to past fund performance. Conventional wisdom has long held that this flow-performance relationship represents irrational "return chasing" by investors, since past performance does not appear to reliably predict future performance (e.g. Carhart, 1997). Berk and Green (2004) (hereafter BG) upended this consensus by pointing out that rational Bayesian learning about unobserved fund manager skill is consistent with a positive flow-performance relationship, while decreasing returns to scale imply a lack of persistence in observed performance. Since both views are *qualitatively* consistent with the existing empirical evidence, distinguishing between the two is a *quantitative* challenge (e.g., as argued by Cochrane, 2013). It is this challenge that we undertake in this paper.

We proceed in two steps. First, building on Roussanov et al. (2018), we estimate the optimally filtered dynamics of beliefs about fund manager skill for the entire cross-section of U.S. equity mutual funds, together with the shape of returns to scale and the prior beliefs about each fund manager's skill, based on their historical *performance* data.[1] In this step, our estimation does not make any assumption on investor behavior driving mutual fund flows. As a result, the estimates serve as a rational benchmark that reflects the econometrician's best estimate of manager's skill at a given point in time. We uncover substantial variation in latent manager skill, which is persistent over time but is subject to fairly steeply decreasing returns to scale.

Second, we show that fund flows respond very little to variation in the estimated rational beliefs. As a result, a measure of misallocation - the difference between actual fund size and the "efficient" fund size implied by the BG model (given the estimated beliefs) - strongly predicts subsequent fund performance. Funds that are "too small" relative to their "efficient" fund size subsequently outperform, and vice versa (funds that are "too big" tend to underperform). We proceed to estimate the parameters of the belief process that best fits the observed relationship between past performance and fund flows (assuming standard Bayesian updating of those beliefs). Our estimation reveals that the "average" mutual fund investor is substantially more optimistic about the underlying fund manager skill than is warranted by the historical performance data. Consistent with the "return chasing" view, investor flows appear to respond too strongly to recent performance relative to more distant historical performance. This response of flows to recent performance suggests either mistaken beliefs about mean reversion in manager skill and the relative role of skill and luck in generating performance, or a type of "recency" bias in beliefs, consistent with models based on the representativeness heuristic of Kahnemann and Tversky.[2] At the same time, flows adjust very sluggishly towards these implied beliefs,

---

[1]This generalizes the estimation approach we follow in Roussanov et al. (2018).

[2]The representativeness heuristic was introduced by Kahneman and Tversky (1972a), Tversky and Kahneman (1974), Tversky and Kahneman (1983) and can be used as an organizing principle for explaining several related biases in probabilistic decision making that imply over-inference from small samples and excessive reliance on the most recent or salient observations at the expense of prior beliefs or "base rates," as detailed in a survey

suggesting: limited attention by investors, or information/search frictions, or both.

Specifically, the first step of our exercise starts by considering that the BG model relies on fund-level decreasing returns to scale (hereafter DRS) to maintain a non-degenerate cross-sectional distribution of funds.[3] However, existing empirical evidence on this key component of the model is mixed.[4] Early studies quantify the magnitude of DRS by regressing the fund performance on assets under management (hereafter AUM).[5] A potential concern in these studies is that fund size is not randomly assigned. An omitted factor, such as manager's skill, can affect both fund size and fund performance, leading to biases in DRS estimates. Pástor et al. (2015) use fund fixed effects to remove this bias under the assumption of constant skill and find statistically insignificant DRS. Building on BG's framework with a flexible parameterization of DRS, we use Kalman filter to express the conditional posterior of manager's skill at a given point in time as a function of historical fund performance and AUM. This allows us to obtain a consistent estimate of the DRS using maximum-likelihood method (MLE), by matching the model-predicted performance based on these posteriors to the performance data. The estimation does not rely on the cross-sectional correlation between size and performance, which is the source of bias in regression-based estimates. A key difference between our method and Pástor et al. (2015) is that the MLE is efficient, giving us enough power to reject the null that DRS is zero. Our estimate indicates that the annual fund performance would decrease by 101 bps for a median fund that increases its assets under management by $100 million, or about 40%, which is smaller than but not far from the estimated DRS in Chen et al. (2004) and Ferreira et al. (2013).

The estimated model for fund performance allows us to compute the posterior belief for each fund manager's skill in each period. In BG's theory, fund size is adjusted by investor in- or outflows in response to the most recent performance to the extent that the expected performance (based on the posterior belief and the impact of DRS) is equated to the fund's expense ratio. This relationship allows us to compute the fund size predicted by the rational (and frictionless) benchmark. We then construct a misallocation measure which is the difference between observed AUM and the model-predicted fund size. We find that this misallocation is large and persistent.[6] More importantly, this measure of misallocation strongly predicts subsequent net performance, indicating that it is not simply an artifact of our mismeasuring fund skill. A standard deviation increase in misallocation leads to a statistically significant

chapter by Benjamin (2018). Rabin (2002), Rabin and Vayanos (2010), Gennaioli and Shleifer (2010), Bordalo et al. (2018), Bordalo et al. (2019), and Bordalo et al. (Forthcoming) develop theoretical models based on this principle that apply to various aspects of decision-making in finance and macroeconomics.

[3]In the absence of frictions and decreasing returns to scale, the cross-section of fund sizes collapses to a point mass where the most skilled fund captures the entire market.

[4]Chen et al. (2004), Ferreira et al. (2013), Yan (2008) use OLS method and find a significant DRS at the fund-level. Pástor et al. (2015) use a recursive demeaning method and fail to find DRS at the fund-level.

[5]Including Chen et al. (2004); Ferreira et al. (2013); Yan (2008).

[6]Roussanov et al. (2018) document large *cross-sectional* misallocation but do not explore its variation over time. Pástor et al. (2020) measure misallocation by looking at properties of fund portfolio returns.

58 bps decrease in net performance. In addition, fund flows also respond to misallocation but much slower than predicted by the frictionless model. This evidence indicates that investors are either adjusting their fund investments slowly in response to new information, or systematically deviating from the rational beliefs, or both.

Motivated by the above findings, we estimate a model of fund flows that allows for all the aforementioned biases to be presented. Specifically, the model allows investors to use a belief updating process that deviates from the rational benchmark as measured by the performance data, at least in terms of the parameters underpinning the Kalman filter. In addition, we allow for a flexible elasticity of how net flows respond to the (implied) beliefs about manager skill. In contrast to our first step, these parameters are estimated by matching the model to the data on fund *flows*, rather than fund performance.

First, we find that, comparing with the rational benchmark, investors are more optimistic about fund managers' average skill. Investors' mean prior about fund skill is 6.3% compared to 1.4% estimated from the performance data. Intuitively, the stark difference is due to the fact that, despite the relatively poor record of average performance of mutual funds, they still enjoy significant inflows, particularly in the first years after inception. This finding is also consistent with direct survey evidence - e.g., Choi and Robertson (2020) find that a majority of investors in active mutual funds cite the belief that these funds tend to deliver higher average returns than passive index funds as an important factor in their investment decision.

Second, we find that investors tend to over-weight recent realized performance (a noisy signal of skill) relative to the prior. In other words, investors seem to think that a fund's recent performance reflects the manager's skill more than what it actually does, consistent with the "over-inference" behavior implied by the representativeness heuristic of Kahnemann and Tversky (see Rabin, 2002 and Benjamin, 2018 for discussion). Moreover, we find that investors discount the more distant information more heavily than a rational Bayesian would. For example, 5-year-lagged performance receives 0.26 times less weight than last-year's performance in the rational benchmark, whereas it is weighted 0.90 times less by investors, as suggested by our estimates. This pattern of belief updating is also consistent with the representativeness heuristic, as argued by Gennaioli and Shleifer (2010) , Bordalo et al. (2018), and Bordalo et al. (2019). The estimated persistence of managers' skill is also substantially lower, at 0.76 in the flow model versus 0.95 in the performance model. This suggests that investors might "chase returns" or over-extrapolate from recent performance because of an implicit belief that a fund manager's ability to generate outperformance varies over time a lot more than it actually does.

Finally, we find that investors adjust their fund investments slowly in response to shifting posterior beliefs. Over time, fund sizes do tend towards the frictionless model-predicted allocations but at a speed much slower than predicted by the theory. In particular, it is expected to take 6 years for a typical fund to reach halfway of the adjustment towards its efficient allocation. This fact might be explained by information/search costs as in Hortaçsu and Syverson (2004)

and Roussanov et al. (2018), or by models of limited attention and/or adjustment costs (see Gabaix (2017) for a survey).

To the extent that investors are heterogeneous in their beliefs about mutual fund performance, we also consider flows into retail and institutional share classes separately. We estimate beliefs implied by institutional fund flows to be closer to the rational beliefs than those of retail investors, even though they are still quite optimistic about fund managers' skill, perceive skill as much less persistent than it appears to be in the historical data, and are also quite slow in adjusting towards efficient allocations.

## 1.1 Literature

To the best of our knowledge, our paper is the first to formally test the model of Berk and Green (2004) or to estimate the beliefs of mutual fund investors that are implied by the data. The closest to our paper is work by Baks et al. (2001), who develop a Bayesian method of performance evaluation and find that even extremely skeptical prior beliefs still lead to sizable allocation to active mutual funds. Pástor and Stambaugh (2012) show that uncertainty about decreasing returns to scale and the priors on average manager skill can lead to slow learning. In a related but distinct branch of literature focusing on asset prices, Pástor and Stambaugh (2002) develop an econometric framework that combines investor priors about asset pricing models and managerial skill with return data.

Our paper is related to a large literature on learning in finance. Learning has been applied in different areas in finance, such as volatility and predictability of asset returns, stock price bubbles, portfolio choice, IPO, trading volume, firm profitability and etc.[7]Among those, a closely related paper is Huang et al. (2012) which use mutual fund flows to infer how rational mutual fund investors are. They find empirical evidence consistent with rational learning. Based on our structural estimation approach, we find investors significantly deviating from the rational benchmark.

Our paper is also related to the vast literature on flow-performance relationship. Prior works show that fund flows respond to fund performance (Ippolito, 1992; Chevalier and Ellison, 1997; Sirri and Tufano, 1998). A closely related paper is Song (2019) which documents that active mutual funds with positive factor-related past returns (besides market factor, e.g., size, momentum, etc.) gain inflows to the point that they significantly underperform. Our results on performance and misallocation share the same conclusion, i.e., funds that are "too big" relative to their efficient fund size subsequently underperform, and vice versa. However, due to our structural approach, we are able to construct a BG-model-based misallocation measure and assess different "structural" sources of investor misallocation. In general, our contribution to this literature is to empirically estimate the canonical BG model and explore its quantitative

---

[7]For an extensive survey, please refer to Pástor and Veronesi (2009).

implications.

This paper also contributes to the large literature on how agents utilize historical information. Malmendier and Nagel (2011) show that investors who have experienced low stock-market returns throughout their lives so far report lower willingness to take financial risks. Greenwood and Nagel (2009) show that younger mutual fund managers invested more heavily in technology stocks than older mutual fund managers around the peak of the tech bubble. Our contribution to this literature is to show that mutual fund investors over-weight recent performance in a manner consistent with models based on the representativeness heuristic such as those in Bordalo et al. (2018) and Bordalo et al. (2019), which feature decision-makers putting excessive weight on the most recent observations (relative to the optimal Kalman filter). Giglio et al. (2019) find that the average sensitivity of an investor's equity share to that investor's subjective expectations of stock returns is small. This evidence is consistent with our finding that investors adjust their investments slowly in response to shifting posterior beliefs of manager skill.

The rest of the paper is organized as follows. In Section 2, we develop the model. In Section 3, we discuss the estimation methods. In Section 4, we describe the data used for the estimation. In Section 5, we present our estimation results. In Section 6, we explore the implications of the estimates. In Section 7, we extend our model to account for investor heterogeneity. Section 8 summarizes our conclusions.

## 2  Model

The model of rational fund flows articulated by BG has two key components: (i) fund performance dynamics driven by unobservable skill and decreasing returns to scale and (ii) fund flows that are driven by learning about fund skill. To bring these components to data, we generalize them along a few dimensions. With the generalizations, our model still nests the rational benchmark as given by BG. In Section 3, the model is structurally estimated to reveal whether and to what extent our generalizations hold according to the data.

### 2.1  Fund performance

The realized gross alpha, denoted by $r_{j,t}$, for an active fund $j \in \{1, ..., N_t\}$ over a time period $t$ is determined by three factors: (i) the fund manager's skill to generate expected returns in excess of those provided by a passive benchmark in that period, denoted by $a_{j,t}$, (ii) the impact of decreasing returns to scale, given by $D(Q_{j,t})$ where $Q_{j,t}$ is the fund's asset under management (AUM), and (iii) an idiosyncratic shock $\varepsilon_{j,t} \sim \mathcal{N}\left(0, \delta^2\right)$. Accordingly,

$$r_{j,t} = a_{j,t} - D\left(Q_{j,t}\right) + \varepsilon_{j,t}. \tag{1}$$

The above equation is the same as equation (1) in BG. Next, we generalize slightly from BG

5

by allowing the manager's skill to be time-varying. We assume manager's skill follows an AR(1) process:

$$a_{j,t} = (1 - \rho)\mu + \rho a_{j,t-1} + \sqrt{1 - \rho^2} \cdot v_{j,t}, \tag{2}$$

where $v_{j,t} \sim \mathcal{N}(0, \kappa^2)$. We specify that a fund's first-period skill is drawn from $\mathcal{N}(\mu, \kappa^2)$, the stationary distribution of the above process. Parameter $\rho$ captures the persistence of the skill level. In the limiting case $\rho \to 1$, skill is fixed over time. The possibility of $\rho < 1$ captures the fact that fund managers and/or their strategies may change over time. More importantly, this possibility later allows us to examine whether investors' reaction to performance history (how they weight recent vs. earlier performance) aligns with the exact underlying degree of persistence of skill.

Following BG, we assume that the manager's skill is not observable to either the investors or the fund manager herself. Let $\widehat{a}_{j,t}$ be a rational investor's belief about $a_{j,t}$ in period $t$. More specifically, $\widehat{a}_{j,t}$ is the posterior mean of $a_{j,t}$ given all the historical information up to $t-1$ (not including $r_{j,t}$ or $Q_{j,t}$). One can apply Kalman filter to derive an expression for $\widehat{a}_{j,t}$. Intuitively, equation (2) can be thought of as describing how the "hidden state", $a_{j,t}$, evolves over time, and equation (1) implies that a signal of this hidden state is $r_{j,t} + D(Q_{j,t})$. Applying the Kalman filter gives us the following recursive formula for the posterior mean:

$$\widehat{a}_{j,t+1} = \rho \left[ \widehat{a}_{j,t} + \frac{\widehat{\sigma}_{j,t}^2}{\widehat{\sigma}_{j,t}^2 + \delta^2} \left( r_{j,t} + D(Q_{j,t}) - \widehat{a}_{j,t} \right) \right] + (1 - \rho)\mu, \tag{3}$$

where the posterior variance satisfies:

$$\widehat{\sigma}_{j,t+1}^2 = \rho^2 \frac{\delta^2 \widehat{\sigma}_{j,t}^2}{\widehat{\sigma}_{j,t}^2 + \delta^2} + (1 - \rho^2)\kappa^2. \tag{4}$$

For the initial period $t$ when fund $j$ was born, we simply have $\widehat{a}_{j,t} = \mu$ and $\widehat{\sigma}_{j,t}^2 = \kappa^2$. In the special case of $\rho = 1$, these expressions coincide with BG's Proposition 1. For $\rho < 1$, the filtered posterior beliefs assign more weights to the realized performance in the more recent periods compared to earlier periods. Later, we will estimate $\delta$, $\mu$, $\kappa$, and $\rho$ all from the performance data.

As to the functional form of $D(\cdot)$, that is, decreasing returns to scale, we assume the following parameterization:

$$D(Q_{j,t}; \eta, \gamma) = \eta \cdot \frac{Q_{j,t}^\gamma - 1}{\gamma}, \quad \gamma \in [0, 1]. \tag{5}$$

This power function specification is fairly flexible. When $\gamma = 1$, it is linear in $Q_{j,t}$; when $\gamma \to 0$, it converges to $\log(Q_{j,t})$. For the intermediate values $\gamma \in (0, 1)$, the function is somewhere in between. One of the reasons to use this flexible parameterization is that the literature is

6

inconclusive on the appropriate functional form for DRS.[8] Later, we estimate the exact value of $\eta$ and $\gamma$ from the performance data.

Our timing convention follows Berk and Green (2004): fund $j$ enters period $t$ with a fund size of $Q_{j,t}$ and an estimate of managerial skill of $\widehat{a}_{j,t}$. Then gross alpha $r_{j,t}$ is realized and both fund manager and investors update their estimate of the managerial skill of fund $j$ to $\widehat{a}_{j,t+1}$.

## 2.2 Fund flows

In the model of BG, the impact of decreasing returns to scale (DRS) under the "efficient" fund size, $\widehat{Q}_{j,t}^{BG}$, exactly offsets the difference between the posterior belief of fund skill ($\widehat{a}_{j,t}$) and the expense ratio ($p_{j,t}$). In other words, we have

$$D(\widehat{Q}_{j,t}^{BG};\ \eta, \gamma) = \widehat{a}_{j,t} - p_{j,t}.$$

In Section 5, we measure capital "misallocation" as any mismatch between observed fund size and the above $\widehat{Q}_{j,t}^{BG}$, estimated under the optimally filtered beliefs on skills (equation 3). In a world with efficient capital allocation, misallocation should be non-existent or possibly a noise independent across periods, and thus on average should not be able to predict fund net performance. However, this is not what we find. The measured misallocation significantly predicts the net performance — funds that are too small relative to the efficient benchmark outperform due to DRS, and funds that are too big underperform. This result suggests possible deviations from the rational model for fund flows. In addition, we find that misallocations tend to be persistent, suggesting investors are rather slow to correct misallocations.

Accordingly, we generalize BG's model slightly by allowing investors to hold a belief that is different from $\widehat{a}_{j,t}$ (which represents the correct or "rational" belief). Let $\widetilde{a}_{j,t}$ denotes the investor's belief. This belief might or might not coincide with $\widehat{a}_{j,t}$. To keep our model empirically tractable, we still impose a structure on $\widetilde{a}_{j,t}$, which we will make clear in just a bit. In addition, we also allow investors to use a set of DRS parameters that might be different from the underlying correct parameters, $\eta$ and $\gamma$. Denote the investors' parameters as $\widetilde{\eta}$ and $\widetilde{\gamma}$. The BG-implied fund size under investor's beliefs, $\widetilde{Q}_{j,t}^{BG}$, is given by:

$$D(\widetilde{Q}_{j,t}^{BG};\ \widetilde{\eta}, \widetilde{\gamma}) = \widetilde{a}_{j,t} - p_{j,t}. \tag{6}$$

For the ease of notation, in what follows, we use the lower-case $q$ to denote the log transformation of any value represented by $Q$. Thus, $\widetilde{q}_{j,t}^{BG} \equiv \log(\widetilde{Q}_{j,t}^{BG})$. Using equation (6) and equation (5), we have

$$\widetilde{q}_{j,t}^{BG} = \frac{1}{\widetilde{\gamma}} \log\left[1 + \frac{\widetilde{\gamma}}{\widetilde{\eta}}(\widetilde{a}_{j,t} - p_{j,t})\right]. \tag{7}$$

---

[8]Linear specifications were used in Pástor et al. (2015); logarithm specifications were used in Chen et al. (2004), Yan (2008), Elton et al. (2012), Ferreira et al. (2013), Reuter and Zitzewitz (2015).

Note that for $\widetilde{\gamma} \to 0$, the above equation reduces to

$$\widetilde{q}_{j,t}^{BG} = \frac{\widetilde{a}_{j,t} - p_{j,t}}{\widetilde{\eta}}.$$

Intuitively, $\widetilde{q}_{j,t}^{BG}$ is the log size for fund $j$ that would be achieved if the capital allocations are fully adjusted to reflect investor beliefs within the period. For empirical application, we maintain the assumption that fund size adjusts towards $\widetilde{q}_{j,t}^{BG}$, albeit allowing a possibly slower rate:

$$q_{j,t} - q_{j,t-1} = \phi(\widetilde{q}_{j,t}^{BG} - q_{j,t-1}) + \xi_{j,t} \tag{8}$$

In the above, $q_{j,t} = \log Q_{j,t}$ is the log of fund size observed in the data, $\phi$ governs the rate of the convergence towards the efficient fund size, and $\xi_{j,t}$ is a shock term. If $\phi = 1$ then fund flows adjust immediately in response to performance information, as is the case of the BG model. Solving for $q_{j,t}$ in equation (8) gives

$$q_{j,t} = \phi\widetilde{q}_{j,t}^{BG} + (1 - \phi)q_{j,t-1} + \xi_{j,t}.$$

We may assume that the error term $\xi_{j,t}$ is an independent innovation at time $t$. However, it is likely that $\xi$ is serially correlated – a fund may carry on growing from one year to the next. So we allow serial correlation through an AR(1) process:

$$\xi_{j,t} = \beta\xi_{j,t-1} + \sqrt{1 - \beta^2} \cdot \zeta_{j,t},$$

and $\zeta_{j,t} \sim \mathcal{N}(0, \omega^2)$ is an innovation at time $t$.

To complete the model, we need to specify $\widetilde{a}_{j,t}$. We assume that it follows the same structure as $\widehat{a}_{j,t}$, but under a potentially different set of performance-model parameters: $\widetilde{\mu}$, $\widetilde{\kappa}$, $\widetilde{\delta}$, and $\widetilde{\rho}$, as well as the DRS parameters: $\widetilde{\gamma}$ and $\widetilde{\eta}$. In other words, we assume that the skill-performance process as perceived by investors follows the same distributional family as we specified in Section 2.1. As a result, $\widetilde{a}_{j,t+1}$ and $\widetilde{\sigma}_{j,t+1}^2$ follow the same recursive formula as in equation (3) and (4) but with tilded parameters:

$$\widetilde{a}_{j,t+1} = \widetilde{\rho}\left[\widetilde{a}_{j,t} + \frac{\widetilde{\sigma}_{j,t}^2}{\widetilde{\sigma}_{j,t}^2 + \widetilde{\delta}^2}\left(r_{j,t} + D(Q_{j,t}; \widetilde{\eta}, \widetilde{\gamma}) - \widetilde{a}_{j,t}\right)\right] + (1 - \widetilde{\rho})\widetilde{\mu}, \tag{9}$$

$$\widetilde{\sigma}_{j,t+1}^2 = \widetilde{\rho}^2\frac{\widetilde{\delta}^2\widetilde{\sigma}_{j,t}^2}{\widetilde{\sigma}_{j,t}^2 + \widetilde{\delta}^2} + (1 - \widetilde{\rho}^2)\widetilde{\kappa}^2. \tag{10}$$

For the initial period $t$ when fund $j$ was born, we have $\widetilde{a}_{j,t} = \widetilde{\mu}$ and $\widetilde{\sigma}_{j,t}^2 = \widetilde{\kappa}^2$. Imposing a structure on $\widetilde{a}_{j,t}$ (instead of estimating it non-parametrically) keeps our model empirically tractable. Moreover, by keeping the rational belief and investor belief in the same distributional family, it facilitates a direct comparison between the two.

It is important to note that $q_{j,t}$ as specified in our model depends on the ratio $\widetilde{\kappa}/\widetilde{\delta}$ but not the absolute sizes of $\widetilde{\kappa}$ or $\widetilde{\delta}$. To see this, note that $q_{j,t}$ depends on the investor's posterior $\widetilde{a}_{j,t}$; however, in Bayesian updating, posterior means are not affected by an increase in the prior variance $\widetilde{\kappa}^2$ *and* a proportional increase in the signal variance $\widetilde{\delta}^2$. Hence, $\widetilde{\kappa}$ and $\widetilde{\delta}$ cannot be separately identified using only data on flows. Instead we can estimate the relative precision

8

parameter $\lambda \equiv \widetilde{\kappa}/\widetilde{\delta}$, which captures how much of variation in the observed outperformance is perceived to be driven by skill (because $\widetilde{\kappa}$ controls the dispersion in the implied distribution of skill) relative to luck (because $\widetilde{\delta}$ is the implied volatility of random shocks to outperformance).

To summarize, our model allows the investor's behavior to deviate from the rational benchmark in several aspects. The first aspect is the way they form beliefs on the managers' skills. We assume investor beliefs to follow the same distributional family as the rational benchmark, but we allow their beliefs to follow a potentially different set of parameters from the "correct" parameters consistent with the performance data. Second and related, we allow investor behaviors to reflect a different degree of DRS from that consistent with the performance data. Third, we allow investors to adjust capital allocations slowly rather than immediately to the efficient benchmark.

At this stage, it is important to point out that the "learning" in our model happens on the same subject as in BG: the investors learn about the managers' skills. One can think of a further level of learning where investors also learn about the parameters (such as $\eta$ and $\gamma$). However, this would lead to a much more complex model beyond the scope of this paper. We partially extend our analysis to address this issue in Section 6.4.

Finally, we do not tempt to specify an explicit model for fund pricing, $p_{j,t}$. A benefit of this approach is to avoid biasing our estimation with a possibly misspecified pricing model. However, consistent with BG, we do assume that $p_{j,t}$ does not reveal about the underlying skill $a_{j,t}$ beyond $\widehat{a}_{j,t}$.

## 3 Estimation

There are three sets of parameters to be estimated: (i) $\eta$, $\mu$, $\kappa$, $\delta$, $\rho$, $\gamma$, which together govern the evolution of fund performance, (ii) $\widetilde{\eta}$, $\widetilde{\mu}$, $\lambda$, $\widetilde{\rho}$, $\widetilde{\gamma}$, which together govern investor's beliefs, and (iii) $\phi$, $\beta$, $\omega$, which affect investor's choices. Below, we describe the estimation strategy for these parameters.

For ease of notation, let $Y_{j,t} = \{r_{j,t}, q_{j,t}, p_{j,t}\}$ denote the data about fund $j$ from period $t$. By equation (1), we can write down the conditional likelihood of observing $r_{j,t}$ as

$$\Pr(r_{j,t} \mid q_{j,t}, p_{j,t}, Y_{j,t-1}, Y_{j,t-2}, ...) \sim$$
$$\mathcal{N}\left[\widehat{a}_{j,t} - D(Q_{j,t}; \eta, \gamma), \ \widehat{\sigma}_{j,t}^2 + \delta^2\right]. \tag{11}$$

In the above, $\widehat{a}_{j,t}$ is the rational posterior on skill conditional on $\{Y_{j,t-1}, Y_{j,t-2}, ...\}$. Its recursive expression is derived in equation (3). Particularly, note that $\widehat{a}_{j,t}$ does not change upon observing the current-period fund size $q_{j,t}$ or price $p_{j,t}$. This is because: (i) price $p_{j,t}$ is assumed not to reveal about the underlying skill $a_{j,t}$ beyond $\widehat{a}_{j,t}$, and (ii) $q_{j,t}$ is a function of $\{Y_{j,t-1}, Y_{j,t-2}, ...\}$ and $p_{j,t}$, plus innovation $\zeta_{j,t}$ that does not hold information about the underlying skill.

From equation (8), we can write down the conditional likelihood of observing $q_{j,t}$:

$$\Pr(q_{j,t} \mid p_{j,t}, Y_{j,t-1}, Y_{j,t-2}, ...) \sim$$
$$\mathcal{N}\left(\phi \widetilde{q}_{j,t}^{BG} + (1 - \phi)q_{j,t-1} + \beta \xi_{j,t-1},\ (1 - \beta^2)\omega^2\right), \tag{12}$$

where $\xi_{j,t-1}$ can be backed out using equation (8) from the observed previous-period fund sizes as follows,[9]

$$\xi_{j,t-1} = q_{j,t-1} - \left[\phi \widetilde{q}_{j,t-1}^{BG} + (1 - \phi)q_{j,t-2}\right].$$

Combining equation (11) and (12), we can write the partial likelihood function (Wooldridge, 2010) as

$$\prod_{j,t} \Pr(r_{j,t}, q_{j,t} \mid p_{j,t}, Y_{j,t-1}, ...) = \prod_{j,t} \underbrace{\Pr(r_{j,t} \mid q_{j,t}, p_{j,t}, Y_{j,t-1}, ...)}_{\text{Performance}} \cdot \underbrace{\Pr(q_{j,t} \mid p_{j,t}, Y_{j,t-1}, ...)}_{\text{Flows}}. \tag{13}$$

In the above, the first part of the likelihood (labeled "performance") tries to fit the observed returns, particularly how returns correlate across periods. Note this part only relies on the performance model and does *not* make any assumptions on how fund sizes are determined (in other words, how investors choose funds). The observed fund sizes do enter this part of the likelihood, but only as conditional variables to account for the DRS. In contrast, the second part of the likelihood (labeled "flows") tries to fit the observed fund sizes.

Maximizing the likelihood in equation (13) estimates the performance model and flow model jointly. It is important to point out that, instead of a joint estimation, we can also conduct our estimation in a separate manner. We can either: (i) maximize the performance part in (13) alone to estimate $\eta, \mu, \kappa, \delta, \rho, \gamma$ (which together govern the evolution of fund performance), or (ii) maximize the flow part alone to estimate $\widetilde{\eta}, \widetilde{\mu}, \lambda, \widetilde{\rho}, \widetilde{\gamma}$ (which together govern investor's beliefs) and $\phi, \beta, \omega$ (which affect investor's choices).

In general, a joint estimation has advantages and disadvantages. It offers more efficiency, however, mis-specification in any part of the likelihood may "contaminate" the estimates in the other part, if the two parts depend on some common parameters. Fortunately, in our context, the performance model and flow model use two completely different sets of parameters. As a result, a joint estimation of equation (13) is equivalent to separate estimations of the performance and flow models.

## 3.1 Identification of decreasing returns to scale

BG model relies on the existence of fund-level DRS to maintain a non-degenerate cross-sectional distribution of funds.[10] However, previous findings on the fund-level DRS are mixed in the literature. Early studies employ the ordinary least square (OLS) method to quantify the magnitude

---

[9]For the likelihood at period $t$, our estimation needs to use the information of $q_{t-1}$ and $q_{t-2}$, so that we restrict the sample to periods of $t \geq 3$.

[10]In the absence of frictions and decreasing returns to scale, the cross-section of fund sizes collapses to a point mass where the most skilled fund captures the entire market.

of DRS by directly regressing fund returns on lagged fund sizes.[11] This method generates biased estimates due to the fact that fund sizes are not randomly assigned to mutual funds. There could exist omitted factors (such as fund skill) that affect both fund size and fund returns.

Recognizing this identification challenge, Reuter and Zitzewitz (2015) use the impact of Morningstar star change on inflow as an exogenous shock to fund size to gauge the causal impact of fund size on performance. Pástor et al. (2015) use a recursive demeaning method to remove the impact of skill on fund size and fund performance. Both studies failed to find statistically significant fund-level DRS.[12]

In this paper, we employ a different method which is to structurally estimate the BG model so that we can explicitly account for manager skills. While operationally this can be implemented through a partial MLE as described above, it is still important to understand, in theory, whether it is possible at all to identify DRS in the BG model from just fund performance and size data. To this end, we offer an argument on the identification of DRS in the BG model. With the generalizations of BG that we made in Section 2, the identification should be much more complex and we do not attempt it here. Nevertheless, our Monte Carlo experiments in the Appendix show that all the parameters in our generalized model can be recovered.

As in the BG model, let $D(Q_{j,t}) = \eta q_{j,t}$, where $q_{j,t}$ is the log size of the fund $j$ in period $t$ (in other words, $\gamma \to 0$). Also, let a fund's skill be persistent over time, that is, $a_{j,t} = a_j$ for all $t$ (in other words, $\rho \to 1$). For the ease of notation, here we will assume that every fund is born at $t = 1$. The performance of fund $j$ in period $t$ is given by

$$r_{j,t} = a_j - \eta q_{j,t} + \varepsilon_{j,t}.$$

So $r_{j,t} + \eta q_{j,t}$ can be regarded as a normally distributed signal centered around the underlying skill $a_j$. At period $T$, the posterior about the skill as seen by econometrician is given by

$$\mathbb{E}(a_j \mid Y_{j,T}, ..., Y_{j,1}) = \frac{\kappa^{-2}\mu + \delta^{-2}\sum_{t=1}^{T}(r_{j,t} + \eta q_{j,t})}{\kappa^{-2} + \delta^{-2}T},$$

where $Y_{j,t} = \{r_{j,t}, q_{j,t}, p_{j,t}\}$ again denotes the data about fund $j$ from period $t$.[13] We assume that $q_{j,T+1}$ does not provide more information about $a_j$ beyond $\{Y_{j,T}, ..., Y_{j,1}\}$, which holds in both BG and our model. With this assumption, we have:

$$\mathsf{E}(r_{j,T+1} \mid q_{j,T+1}, Y_{j,T}, ..., Y_{j,1}) = \mathsf{E}(a_j|Y_{j,T}, ..., Y_{j,1}) - \eta q_{j,T+1}$$

$$= \frac{\kappa^{-2}\mu}{\kappa^{-2} + \delta^{-2}T} + \frac{\delta^{-2}\sum_{t=1}^{T} r_{j,t}}{\kappa^{-2} + \delta^{-2}T} + \frac{\delta^{-2}\eta\sum_{t=1}^{T} q_{j,t}}{\kappa^{-2} + \delta^{-2}T} - \eta q_{j,T+1}.$$

As an identification argument, consider $T \to +\infty$, in which case the role of prior diminishes:

$$\mathsf{E}(r_{j,T+1} \mid q_{j,T+1}, Y_{j,T}, ..., Y_{j,1}) \to \bar{r}_{j,T} + \eta\left(\bar{q}_{j,T} - q_{j,T+1}\right),$$

---

[11]Including Chen et al. (2004); Ferreira et al. (2013); Yan (2008).

[12]Pástor et al. (2015) find decreasing returns to scale at the industry level.

[13]See Proposition 1 in BG for more details about the derivation.

where

$$\overline{r}_{j,T} \equiv \frac{1}{T} \sum_{t=1}^{T} r_{j,t},$$

$$\overline{q}_{j,T} \equiv \frac{1}{T} \sum_{t=1}^{T} q_{j,t}.$$

The identification of $\eta$ should be clear from the above expression. The expression also offers an intuitive way to relate DRS to the observed features of the data.. For large $T$, the DRS $\eta$ essentially manifests itself as the elasticity at which the deviation of return from historical average (i.e., $r_{j,T+1} - \overline{r}_{j,T}$) responds to a deviation of fund size from historical average (i.e., $\overline{q}_{j,T} - q_{j,T+1}$).

# 4    Data

We collect data from CRSP and Morningstar. Our sample contains 2,885 well-diversified actively managed domestic equity mutual funds from the United States between 1965 and 2014. Our sample has 31,089 fund-year observations. We closely follow data-cleaning procedures in Berk and van Binsbergen (2015) and Pástor et al. (2015).

There are three main data variables to be used in estimation: annual gross realized alpha (i.e. fund performance), fund size, and expense ratio. To compute the annual realized alpha $r_{j,t}$, we start with monthly return data. We first augment each fund's monthly net return with its monthly expense ratio (1/12th of the annual expense ratio) to get the monthly gross return. Then, we regress the excess gross return (over the 1-month U.S. T-bill rate) on risk factors throughout the life of the fund to get the betas for each fund.[14] We multiply betas with the factor returns to get the benchmark returns for each fund at each point in time. We subtract the benchmark return from the excess gross return to get the monthly gross alpha. Last, we aggregate the monthly gross alpha to the annual realized alpha $r_{j,t}$. We use the Fama-French six-factor model of Fama and French (2018) as the benchmark. There is an alternative way of risk adjustment which is to use Morningstar benchmark index data. But this data is not freely available. However, Pástor et al. (2015) show that in terms of studying the impact of decreasing returns to scale and performance, the two methods yield similar results. Therefore, we take the route of factor adjustment.

Fund size for each year is the fund's AUM at the end of the previous year. To make fund size comparable across time, we inflate all the fund sizes to December 2011 dollars by following Pástor et al. (2015)'s method which is to use the ratio of the total market value of all CRSP

---

[14]An alternative way to estimate the betas, adopted by several studies in the literature, is to carry out the regressions here with expanding windows. Please see the appendix where we estimate our model with the alphas resulted from expanding-window betas. None of our results change qualitatively.

stocks in December 2011 to its value at the end of the previous year.[15] In our dataset, there is a huge skewness in funds' AUM. From the summary statistics, we can see that the mean of funds' AUM is much larger than the median. The funds at the 99 percentile are over 2,600 times larger than the funds at the 1 percentile. And the fund size at the third quartile is over 11 times larger than the fund size at the first quartile. In the literature, to study the impact of decreasing returns of scale, dollar amount of the funds' AUM were used in Pástor et al. (2015), and the logarithm of the funds' AUM were used in Chen et al. (2004), Yan (2008), Elton et al. (2012), Ferreira et al. (2013), and Reuter and Zitzewitz (2015). Our flexible functional form of the DRS allows the data to inform us about the shape of the function. To lessen the effects of "incubation bias"[16], following Fama and French (2010), we limit the tests to funds that reach 15 million 2011 dollars in AUM. Once a fund passes the threshold, it is included for all subsequent periods, so this requirement does not create selection bias.

In the mutual fund industry, a single mutual fund may provide several share classes to investors that differ in their fees structures. Following much of the literature (with some exceptions, e.g., Bergstresser et al., 2009), we conduct our analysis at the fund level instead of the share class level. We compute a fund's AUM by summing AUM across the fund's share classes, and compute the fund's realized alphas, expense ratios by using AUM weighted average across share classes.

[Table 1 about here.]

# 5   Parameter Estimates

The parameter estimates of the performance and flow models are presented in Table 2.

[Table 2 about here.]

## 5.1   Performance model parameters

Our estimate of the decreasing returns to scale (DRS), $\eta$, is 22 basis points with a $t$-statistic of 2.2. To see the economic magnitude of this estimate, consider an increase in fund size by $100 million, which is about a 40% increase in the size for the median size fund. The estimate of $\eta$ indicates that such an increase in size is associated with a decrease in expected annual fund performance of 101 bps ($0.0022 \times \log(100) = 0.0101$). To compare, Chen et al. (2004) find

---

[15]Alternatively, we can use consumer price index (see https://fred.stlouisfed.org/series/CPIAUCSL) to inflate all the fund sizes to December 2011 dollars. We did so as a robustness check and all of our results remain qualitatively the same.

[16]For details, please refer to Evans (2010).

that the same increase in the fund size leads to around 110 bps decrease in fund performance.[17] Ferreira et al. (2013) find that the same increase in the fund size leads to around 124 bps decrease in fund performance.[18] Note that both papers relied on OLS regressions. Pástor et al. (2015) overcomes the potential bias in OLS by a method of recursive demeaning. Their estimate of DRS is statistically insignificant, potentially due to a lack of statistical power.

To put the DRS estimate further in perspective, we note that the volatility of annual performance in our data is around 7.3 percent. Hence, a decrease of 101 bps of a fund's performance is approximately 14% of the annual volatility in mutual fund's performance.

Parameter $\gamma$, which measures the curvature of DRS, is estimated to be very close to zero. The estimate implies that a log specification of decreasing returns to scale is most in line with the data.

The mean of the prior distribution of managerial skill is 1.4% (per annum). This number is positive and statistically significant, which means that an average active mutual fund manager is skilled. This result is consistent with previous literature, for example, Berk and van Binsbergen (2015). Based on our estimates, for a fund with the average skill level (1.4%) and expense ratio (1.18%), its size should be no larger than $3 million[19], otherwise, DRS would make it produce negative expected performance. However, in the data, the average fund size is $1.45 billion. This observation indicates that the mutual fund industry on average might be too large. Meanwhile, we do find a significant variation in skill across funds. For example, for a fund with the skill level one standard deviation (1.4%, i.e., the estimated value of the standard deviation of skill prior) higher than the average, charging avearge expense ratio, its size can grow to around $1.6 billion[20] before generating negative expected performance.

Parameter $\rho$, which measures the persistence in fund manager's skill is estimated to be 0.95. The estimate indicates a fund's skill changes slowly over time and consequently, distant past performance are likely still relevant in predicting a manager's skill. This result of skill persistence is consistent with Berk and van Binsbergen (2015), who find that cross-sectional differences in value added persist for as long as 10 years.

## 5.2 Flow model parameters

First, we see that parameter $\widetilde{\mu}$, the mean prior belief on manager skill, is estimated to be 6.3%. This estimate represents the (implicit) beliefs of mutual fund investors as it is estimated from the flows data (instead of the performance data). This estimate is close to the value of the prior mean (6.5%) shown in BG's Table 1. In their paper, the value of the prior mean is picked

---

[17]We use their coefficient of DRS for the monthly 4-factor gross alpha (in their Table 3), 0.020. The result 110 bps equals 0.020 times the logarithm of 100 million dollars times 12.

[18]We use their coefficient of DRS for the quarterly 4-factor alpha for US funds (in their Table 5), 0.0675. The result 124 bps equals 0.0675 times the logarithm of 100 million dollars times 4.

[19]We compute the BG implied size as $exp\left((0.014 - 0.0118)/0.0022\right) = 2.7$.

[20]We compute the BG implied size as $exp\left((0.014 + 0.014 - 0.0118)/0.0022\right) = 1,578$.

using a calibration procedure to match the empirical relation between the flow of funds and performance, which is conceptually similar to our estimation based on flow data.

These magnitude is striking, suggesting extreme over-optimism of investors about fund manager skill, as the difference between the mean of investor prior of skills ($\widetilde{\mu}$) and the rational benchmark ($\mu$) is about 4.9% (the former is almost five times as large as the latter). However, the expected outperformance depends not just on managers' skill, but also on its erosion by fund size. Our estimate of the decreasing returns to scale (DRS) from the flow model, $\widetilde{\eta}$, is 70 basis points which is highly statistically significant. To gauge the economic magnitude of this estimate, considering an similar increase in fund size of $100 million, the estimate of $\widetilde{\eta}$ indicates that such an increase in size is associated with a decrease in expected annual fund performance of 346 bps. (We compute this decrease in future performance based on Eq. (5) with $\widetilde{\eta} = 0.0070$ , $\widetilde{\gamma} = 0.03$, and $Q = 100$.) In other words, investors flows react in a way that anticipating more rapid erosion in perfromance due to inflow of capital than is evidenced by actual fund returns. Interestingly, this is somewhat in contrast with the survey evidence in Choi and Robertson (2020), who show that only a small minority (about one fifth) of (individual) mutual fund investors believe in decreasing returns to scale, although that fraction is higher among wealthier investors, who are more likely to influence total fund flows.

Despite this more conservative estimate of the DRS, the difference in the prior beliefs about manager skill between the rational benchmark and the flows-based model indicates that for an average fund ($1.45 billion), investors expect that the fund generates "extra" outperformance of 81bps per annum compared to what the performance data tell us.[21] Intuitively, the stark difference is due to the fact that despite the relatively mediocre performance record of mutual funds, on average, they still enjoy significant inflows, especially in the early years after a fund's inception. In Figure 1, we plot the average annual net inflow of funds as a function of fund age. We can see that before the age of 5, the average inflows are statistically greater than zero (whereas older funds tend to experience net outflows, which are statistically significant starting at 10 years).

[Figure 1 about here.]

Parameter $\phi$, which measures the adjustment rate of flow towards the efficient fund size, is estimated to be 0.07. Recall that for $\phi = 1$, we have $q_{j,t} = \widetilde{q}_{j,t}^{BG}$, that is, the market adjusts capital allocations completely in each period so as to fully reflect investor beliefs. The other polar case, $\phi = 0$, means the fund flows do not adjust towards the efficient allocations at all. Our point estimate of $\phi$ indicates that fund flows do respond to past performance, but slowly. Based on this estimate, it takes about 6 years for a typical fund to reach halfway of $\widetilde{Q}_{j,t}^{BG}$,

---

[21]The average optimism is computed as follows: $\left(\widetilde{\mu} - \widetilde{\eta} \times \left(1453^{\widetilde{\gamma}} - 1\right)/\widetilde{\gamma}\right) - \left(\mu - \eta \times \log(1453)\right) = \left(0.063 - 0.007 \times \left(1453^{0.03} - 1\right)/0.03\right) - (0.014 - 0.0022 \times \log(1453)) = 0.0081.$

the efficient fund size under investor belief.[22] This slow adjustment might be explained by information/search costs as in Hortaçsu and Syverson (2004) and Roussanov et al. (2018), or by models of limited attention and adjustment costs (see Gabaix, 2017 and the references therein).

Parameter $\lambda \equiv \widetilde{\kappa}/\widetilde{\delta}$ is estimated to be 0.694, which is substantially larger than $\kappa/\delta = 0.203$. Recall that $\delta^2$ is the variance of $\varepsilon_{j,t}$, the noise component of fund performance, and $\kappa^2$ is the variance of the prior, which represents the dispersion of skill across managers, and hence the variation in performance due to skill. Our estimates indicate that the average investor preceives realized performance to be a more precise signal of underlying skill than it is in reality. Put differently, investors seem to under-estimate the role of luck (relative to skill) in driving variation in funds' performance. Consequently, investors tend to over-react to recent fund performance. We discuss this evidence of "recency" in greater detail in the following Section 6.2. Overreaction to recent performance needs not be in conflict with the evidence of slow adjustment mentioned above. Intuitively, the parameter $\phi$ captures the fraction of active investors (or investor-dollars) that adjust allocations in a particular period to reflect their beliefs, and parameter $\rho$ captures the degree of perceived mean reversion in the posterior beliefs of this fraction of active investors. From an identification perspective, intuitively, $\rho$ takes a small value if the sign of flow responds immediately to very recent performance, even when (possibly) the more distant historical performance points to a different sign of flow response. On the other hand, $\phi$ takes a small value when the magnitude of flow is small responding to past performance overall.

# 6 Implications

In this section, we use several exercises to illustrate the implications of our structural estimated model. First, we illustrate the persistence of misallocations by examining whether misallocation predicts performance. Second, we graphically illustrate how investors weight the historical performance of a fund when updating their beliefs about the fund's skill, and compare their weighting scheme to what a rational investor would do. Third, we relate our interpretation of the "recency" bias with the interpretation based on the representativeness heuristic. Lastly, we compare the estimated rational beliefs and investor beliefs in terms of the ability to predict fund performance out-of-sample. We also compare them with simple prediction rules, such as 3-year, 5-year and 10-year average past performance.

## 6.1 Misallocation and performance

The central prediction of BG's model is that, if capital allocation to funds is efficient (so that DRS offsets any positive skill net of fees) then fund (out)performance net of fees should not be

---

[22]For this calculation, we start with a fund whose: (i) skill equals $\mu$, (ii) initial fund size equals the median fund size ($64 million), and (iii) fund price fixes at the median expense ratios (1.14%). We set all shocks in the model to zero and check how long it takes for the fund size to reach half of $\widetilde{Q}_{j,t}^{BG}$.

forecastable. Following this line of thought, we analyze whether misallocation, measured as the difference between fund's actual size and efficient size that we estimate, can predict net fund performance. If misallocation is non-existent or simply due to noise and thus independent across periods, then on average it should not predict fund net performance. However, if fund flows don't adjust properly to reflect rational updating of the fund managers' skill then misallocation will be able to predict net performance. Funds that are "too small" relative to the benchmark would subsequently outperform due to DRS, and funds that are "too big" would subsequently underperform.

The results are provided in Table 3. In column (1), we regress the net performance ($r_{j,t} - p_{j,t}$) onto misallocation, computed as the difference between the actual and efficient fund size, $q_{j,t} - \widehat{q}_{j,t}^{BG}$. Here, the efficient fund size is defined by (compare to equation 6)

$$D(\widehat{Q}_{j,t}^{BG};\ \eta, \gamma) = \widehat{a}_{j,t} - p_{j,t}.$$

We find a statistically significant coefficient in front of the misallocation measure. In terms of economic magnitude, a standard deviation increase in misallocation leads to 58 bps decrease in the expected performance. To check whether our misallocation measure is sensible, we break it up into $\widehat{q}_{j,t}^{BG}$ and $q_{j,t}$, separately. The results are provided in column (2). As intended, the coefficient in front of $\widehat{q}_{j,t}^{BG}$ is positive, while the coefficient in front of $q_{j,t}$ is negative, and the magnitudes of the two coefficients are virtually the same.

As a robustness check, we repeat the above analysis but replace $\widehat{q}_{j,t}^{BG}$ with $\widetilde{q}_{j,t}^{BG}$, which is the efficient fund size computed using investor's belief (see equation 6). The results are provided in columns (3) and (4), which are not qualitatively different from columns (1) and (2). Similarly, if we re-run all the above regressions with the dependent variable as the gross performance (instead of the net performance), all of the results remain qualitatively the same: misallocation strongly predicts future fund (out-)performance.

The results in Table 3 show that there is a significant negative relationship between net performance and misallocation. We further quantify the extent to which funds that are "too small" relative to their efficient scale outperform and, in turn, funds that are "too big" underperform. Table 4 presents evidence for ten decile portfolios of funds sorted on the misallocation measure. For each portfolio, in Panel A, we report the equal weighted average net performance and in Panel B, we report the value weighted average net performance. We find that, the funds that are "too big" underperform subsequently with underperformance ranging from -1.2% to -3.0% per annum. The underperformance is generally statistically significant relative to the Fama-French 5 factor model as well as momentum for the funds in top half of the distribution of misallocation. We also find that the funds that are "too small" outperform subsequently with outperformance ranging from 0.9% to 1.8% per annum, although such outperformance is significant only for the bottom decile of equal weighted portfolios of funds (based on misallocation).

[Table 3 about here.]

[Table 4 about here.]

Lastly, we test the prediction of the BG model that underpins the (lack of) persistence in performance: investor flows respond to misallocation of capital, whereby inflows inrease the size of funds that are "too small" given their updated skill, while outflows shrink the funds that are "too big." Thus, we regress flow $(q_{j,t+1} - q_{j,t})$ onto the misallocation measure and other controls. The results are provided in Table 5. If investors are quick to correct misallocation and push funds' AUMs towards their efficient levels, the magnitude of the coefficient in front of the misallocation should be close to 1. Meanwhile, we find that the magnitude of the coefficient is significantly *smaller* than 1. This result is consistent with our estimate of $\phi$ in the model, meaning that fund flows respond to misallocation but the response is much weaker than predicted by the frictionless model. As a robustness check, we repeat the above analysis but replace $\widehat{q}_{j,t}^{BG}$ with $\widetilde{q}_{j,t}^{BG}$, which is the efficient fund size computed using implied investor beliefs (see equation 6). The results, provided in columns (3) and (4), are not qualitatively different from columns (1) and (2), albeit the sensitivity of flows to misallocation is slightly stronger (and even more significant statistically). In the appendix, we show that our results are also robust to the inclusion of several additional control variables such as fund family level advertising spending, the average performance of other funds' in the same Morningstar category, and Morningstar ratings.

[Table 5 about here.]

## 6.2 Recency in flows-implied beliefs

A common behavioral interpretation of mutual funds return-chasing behavior is a "recency" bias in beliefs. Our model of posterior belief about fund skill can be interpreted as a weighting scheme that flexibly incorporates information about the past performance of a fund in a way that is controlled by several key parameters. There are two important aspects in the investor's weighting scheme: (i) the extend to which more distant information is discounted, as measured by parameter $\widetilde{\rho}$, and (ii) how informative the realized performance is about the underlying skill, in comparison to the prior, which is measured by $\lambda = \widetilde{\kappa}/\widetilde{\delta}$.

[Figure 2 about here.]

In Figure 2, we plot the weighting schemes for four cases. Each weighting scheme shows how to weight the historical DRS-adjusted performance, $\{r_{j,t} + D(Q_{j,t}), \ t < T\}$, when forming the posterior about $a_{j,T}$.[23]

---

[23]The weighting scheme can be easily computed by exploiting the fact that the posterior should be a linear function of the prior and signals. Technically, we first simulate the model for a number of funds and $T$ periods, then regress $\widehat{a}_{j,T}$ or $\widetilde{a}_{j,T}$ on $\{r_{j,t} + D(Q_{j,t}), \ t < T\}$ and the prior. If implemented correctly, the regression will have a $R^2 = 1$ and the coefficients will sum up to 1.

The blue curve with stars plots the weights implied by the estimated performance model (i.e., how the rational posterior $\hat{a}_{j,t}$ weights historical information). The yellow curve with circles plots the weighting scheme under the estimated investor beliefs (i.e., how the $\tilde{a}_{j,t}$ weights historical information). Comparing these two curves, we can clearly see that relative to the rational benchmark, investors over-weight recent performance (lag period 1 to 4) in a manner consistent with models based on the "representativeness" heuristic, and under-weight distant performance information (lag period 5 onward).

To obtain more intuition for the role of different parameters in the weighting scheme, we plot two additional curves. The red curve with squares plots the same investor's weighting scheme as the yellow curve except that we impose $\tilde{\rho} = 0.95$ that is more consistent with the Bayesian estimate based on past performance (the estimated flows-implied $\tilde{\rho} = 0.76$). Comparing the red and yellow curves, we see that a larger $\tilde{\rho}$ (red curve) implies more weight on the distant signals, which is intuitive because a larger $\tilde{\rho}$ means manager skill is more persistent over time. Another useful result here is that the weight on the prior for the red curve is 0.25; for the yellow curve is 0.54. Hence, when $\tilde{\rho}$ is larger (red curve) the posterior puts less weight on the prior. Intuitively, this is because a smaller $\tilde{\rho}$ means that the manager skill reverts to the stationary distribution faster, so that the posterior should rely more on the prior. In the opposite extreme case where $\tilde{\rho} \to 1$, the weight on the prior will go to zero as $T \to \infty$.

The gray curve with asterisks plots the same investor's weighting scheme as the yellow curve except that we impose $\lambda = 0.45$ (the estimated $\lambda = 0.69$). The weight on prior for gray curve is 0.68; for the yellow curve is 0.54. We can see that as $\lambda$ gets smaller (gray curve), the weight on the prior increases. Intuitively, this is because $\lambda = \tilde{\kappa}/\tilde{\delta}$ measures how one perceives the precision of the signal relative to the precision of the prior.

Finally, it is important to note that the weighting curve changes differently when we vary $\tilde{\rho}$ and when we vary $\lambda$. In particular, $\tilde{\rho}$ mainly calibrates the relative weights of recent vs. distant signals, while $\lambda$ mainly calibrates the relative weights of signals vs. prior. Conceptually, this difference explains how the two parameters can be separately identified from the data.

## 6.3 Representativeness heuristic

As shown above, investor flows appear to respond overly strongly to recent performance relative to more distant historical performance, as well as to prior beliefs. This finding can be rationalized with either: (i) mistaken beliefs about mean reversion in manager skill ($\rho$) and the relative role of skill and luck in generating performance ($\kappa/\delta$), or (ii) a bias in beliefs with overweighting the recent observations that is consistent with models of beliefs based on the representativeness heuristic of Kahneman and Tversky (1972a) (e.g., Gennaioli and Shleifer, 2010; Rabin and Vayanos, 2010; Bordalo et al., 2018; Bordalo et al., 2019; see survey by Benjamin, 2018 for further discussion).

To gain some intuition about how representativeness applies to our context, consider a simplified model where an investor learns about the underlying skills of a fund manager $(a_t)$ using noisy signals, the realized returns $(r_t)$. For the simplicity of exposition, let's assume the hidden skill is constant over time (i.e., $a_t = a$). For a rational investor, conditional on observing a positive realized return $(r_t > 0)$, the probability that the manager is skilled (i.e., $a > 0$) is given by the Bayes rule:

$$P(a > 0 | r_t > 0) = \frac{P(r_t > 0 | a > 0)P(a > 0)}{P(r_t > 0)},$$

where

$$P(r_t > 0) = P(r_t > 0 | a > 0)P(a > 0) + P(r_t > 0 | a \leq 0)P(a \leq 0).$$

Representativeness heuristic is often described as neglect of "base rates" - in the present context that would mean ignoring the fact that the probability of generating a positive return simply due to luck is quite high, and so the denominator $P(r_t > 0)$ is quite large (e.g., well in access of 0.5), as well as potentially overestimating the probability that a manager is skilled, $P(a > 0)$. An investor who is applying representativeness heuristic would thus overweight the "representative" characteristic: $P(r_t > 0 | a > 0)$, that is, the signal $r_t > 0$ is "representative" of $a > 0$. This overweighting would inflate the posterior $P(a > 0 | r_t > 0)$ (and consequently deflate $P(a \leq 0 | r_t > 0)$). The same intuition applies to the opposite posterior $P(a \leq 0 | r_t \leq 0)$, which will also be too inflated.

Hence, an investor who is relying on the representativeness heuristic would put too much probability on high (or low, respectively) manager skills after observing good (or bad, respectively) returns, down-weighting the role of luck in generating these returns. Consequently, the representativeness heuristic creates return-chasing behavior in a way that our model interprets as putting a disproportionate weight on skill relative to luck $(\kappa/\delta)$.

Variations on the representativeness heuristic collectively known as "local representativeness," whereby population properties are mistakenly ascribed to small samples ("law of small numbers" and "gamblers' fallacy"), can also lead to a form of "recency" bias. In particular, Rabin (2002) and Rabin and Vayanos (2010) show how versions of this bias can lead to excessive extrapolation from recent observations in a way that is consistent with a downward-biased persistence parameter $\rho$. Bordalo et al. (2018) and Bordalo et al. (2019) introduce "diagnostic expectations" as a way of modeling the dynamic evolution of beliefs that are subject to the representativeness heuristic. We explore incorporating diagnostic expectations into our baseline model in the Appendix and show that it is subsumed by our flexible Bayesian approach that treats belief parameters as free rather than constrained by historical data.

### 6.4 Out-of-sample prediction

In this section, we explore the out-of-sample prediction of fund performance using estimated rational beliefs $(\hat{a}_{j,t})$, flows-implied investor beliefs $(\tilde{a}_{j,t})$, and naive predictors such as moving averages of past outperformance.

Specifically, we estimate our model parameters for both the performance-based and flows-based models using data from 1965 up to 2009. Then, we use the estimated parameters to generate the rational posteriors and investor posteriors of fund skills from 2010 to 2014. These posteriors are computed using equation (3) and (9), respectively. Importantly, only the performance data after 2009 are used in computing these posteriors — it is the parameter estimates that are obtained without using post-2009 data. Next, we subtract the impact of DRS from the posterior to construct the predictor for fund performance $r_{j,t}$. Aside from filtered posteriors, we also consider simple moving averages, such as 5-year average past performance, $\sum_{k=1}^{5} r_{j,t-k}/5$, as predictors for $r_{j,t}$. Finally, to make sure this exercise is a truly out-of-sample prediction, we replace our measure of gross performance of funds with expanding-window alphas. This way, the betas used to estimate fund performance before and in 2009 do not rely on any data beyond 2010.[24]

We compute the mean squared error (MSE) of the various predictors for $r_{j,t}$ from 2010 to 2014. The results are provided in Table 6. We find that the rational posterior has the smallest MSE, which means that it performs the best in the out of sample prediction. Both rational posterior and the investor posterior outperform the naive predictors. The result suggests that it offers some advantage to use a properly specified econometric model to extract useful information from historical data.

[Table 6 about here.]

## 7 Extension: institutional vs. retail investors

Up to now, we have assumed a representative investor. While this is clearly a dramatic simplification, we cannot identify the heterogeneous beliefs of individual investors without account-level data, which is not available to us. We can, however, consider different groups of investors aggregated into (admittedly coarse) segments. In particular, it is natural to ask whether institutional investors are more sophisticated (and, by extension, more "rational") than households (see, e.g., Glode et al., 2017).

In order to explore this question, we exploit the fact that there are usually multiple share classes of the same fund, some are marketed to retail investors and others are only available to

---

[24]More specifically, we fix the starting point of the window at the birth time of the fund and the endpoint of the window progresses along time. The initial size of the window is 24 months. Then, from the 25th month onwards, we expand the window monthly to compute the betas. Last, we aggregate monthly alpha into annual alpha.

institutions. We extend our model of fund flows to allow for two types of investors, who hold different beliefs and invest in the two different classes of the same fund. We assume that each type of investors are dogmatic in their beliefs (ignoring the fact that their beliefs might differ from those of the other type of investors). While this is admittedly a simplification, we make this assumption for the sake of tractability.

## 7.1 Model and estimation

Let there be two different classes for each fund: (1) institutional and (2) retail. Then the sizes of the two share classes add up to the total net assets of the fund:

$$Q_{j,t} = Q_{j,t}^{(1)} + Q_{j,t}^{(2)}.$$

In the above, we use superscripts 1 and 2 to denote the two share classes. For the investors in class $k$, where $k \in \{1, 2\}$, we denote their posterior beliefs as $\widetilde{a}_{j,t}^{(k)}$ and $\widetilde{\sigma}_{j,t}^{(k)}$, which follow the same structure as in equation (9) and (10), but under a different set of parameters $\widetilde{\mu}^{(k)}$, $\widetilde{\kappa}^{(k)}$, $\widetilde{\delta}^{(k)}$, $\widetilde{\rho}^{(k)}$, $\widetilde{\eta}^{(k)}$, and $\widetilde{\gamma}^{(k)}$. So their perceived efficient fund size is given by (compare to equation 7)

$$\widetilde{q}_{j,t}^{BG(k)} = \frac{1}{\widetilde{\gamma}^{(k)}} \log \left[ 1 + \frac{\widetilde{\gamma}^{(k)}}{\widetilde{\eta}^{(k)}} (\widetilde{a}_{j,t}^{(k)} - p_{j,t}^{(k)}) \right]. \tag{14}$$

Note that the above is the *fund* size, not the share class size, as perceived by class-$k$ investors. Importantly, this is because DRS happens at the fund level, not the share-class level. Because the investment strategies are the same at the two share classes within the same fund, as the size of one share class increases, it should cause decreasing return on the other share class as well.

The flow for the share class $k$ in fund $j$ is specified as follows (compare to equation 8)

$$q_{j,t}^{(k)} - q_{j,t-1}^{(k)} = \phi^{(k)} \cdot \left( \log \left( \psi^{(k)} \right) + \widetilde{q}_{j,t}^{BG(k)} - q_{j,t-1}^{(k)} \right) + \xi_{j,t}^{(k)}. \tag{15}$$

The new parameter $\psi^{(k)}$ presents the proportion of class $k$ in the fund's efficient size.[25] We will estimate the parameter from the data. Given that the sizes of the two share classes should add up to the fund size, one may impose that $\psi^{(1)} + \psi^{(2)} = 1$ in the estimation. However, from the standpoint of estimation, this restriction is not necessary.[26]

Again, we allow serial correlation in the residual $\xi_{j,t}^{(k)}$ as in Section 2.2, but with class-specific parameters: $\beta^{(k)}$ and $\omega^{(k)}$, $k \in \{1, 2\}$.

The estimation of the extended model follows Section 3, with a few differences. First, the data is expanded to included class-specific records:

$$Y_{j,t} \equiv \left\{ r_{j,t}, p_{j,t}^{(1)}, p_{j,t}^{(2)}, q_{j,t}^{(1)}, q_{j,t}^{(2)} \right\}.$$

---

[25]Notice that $\log \left( \psi^{(k)} \cdot \widetilde{Q}_{j,t}^{BG(k)} \right) = \log \left( \psi^{(k)} \right) + \widetilde{q}_{j,t}^{BG(k)}$.

[26]Without the restriction, the interpretation of $\psi^{(k)}$ is more akin to what class-$k$ investors *perceive* as the proportion of fund size that they should be responsible for.

Second, the likelihood function fits the sizes of share classes, instead of the fund size. The partial likelihood on the flow model is

$$\prod_{j,t} \left[ \prod_{k \in K(j,t)} \Pr \left( q_{j,t}^{(k)} \mid p_{j,t}^{(1)}, p_{j,t}^{(2)}, Y_{j,t-1}, Y_{j,t-2}, \dots \right) \right].$$

In the above $K(j, t) \subseteq \{1, 2\}$ denotes the share classes that fund $j$ offers in period $t$ (some funds did not have an institutional class in earlier years). Ideally, we want to estimate the model with the sample that has both share classes available at the same time so that the comparison between the estimated parameters for different share classes would be meaningful. To achieve that purpose, we focus on the sample period of 2000-2014, because before 2000, there were very few institutional share classes in the data. Finally, note that the partial likelihood on the performance model stays exactly the same as in Section 3, because we have made no changes to the performance model.

## 7.2   Results

The estimates of the two-class model are reported in Table 7. It reveals that there is indeed some heterogeneity in beliefs, at least between the (broadly defined) institutional and retail investors. Overall, compared to the retail classes, the beliefs implied by the flows at institutional-share classes are closer to the rational benchmark as estimated from the performance data (Table 2).

Specifically, first, institutional investors hold a prior about the skill of fund managers that is lower than retail investors and closer to the benchmark ($\widetilde{\mu}^{(1)} = 3.3\%$ vs. $\widetilde{\mu}^{(2)} = 5.4\%$, and $\mu = 1.4\%$). Second, compared to retail investors, institutional investors believe that realized performance are driven more by luck rather than skill, which is more aligned with the performance data ($\lambda^{(1)} = 0.429$ vs. $\lambda^{(2)} = 0.705$, and $\kappa/\delta = 0.203$). Third, institutional investors exhibits a less degree of "over-extrapolation" compared to retail investors, in the sense of believing that skill is more persistent; however, the degree of "over-extrapolation" is still very strong compared to the benchmark ($\widetilde{\rho}^{(1)} = 0.737$ vs. $\widetilde{\rho}^{(2)} = 0.700$, and $\rho = 0.954$). Fourth, institutional investors adjust capital allocations somewhat faster than retail investors ($\phi^{(1)} = 0.077$ vs. $\phi^{(2)} = 0.064$). Lastly, institutional investors perceive a degree of DRS that is smaller than retail investors and closer to the benchmark.

[Table 7 about here.]

## 8   Conclusion

We estimate a structural model of investor beliefs implicit in the mutual fund flows. We compare this estimated model with the rational Bayesian benchmark that is based on past performance. Our estimates imply that investors are more optimistic about fund manager's average skill than warranted by the historical data. They over-weight recent performance in a manner consistent

23

with models based on the "representativeness" heuristic, yet respond slowly to changes in these beliefs, consistent with limited attention and/or informational frictions. Beliefs implied by flows to institutional funds display less bias than those estimated from retail funds. These results offer new perspective on mutual fund investor's behaviors beyond the flow-performance relationship and pave roads for fruitful future research on household finance.

# References

**Baks, Klaas P, Andrew Metrick, and Jessica Wachter**, "Should investors avoid all actively managed mutual funds? A study in Bayesian performance evaluation," *The Journal of Finance*, 2001, *56* (1), 45–85.

**Ben-David, Itzhak, Jiacui Li, Andrea Rossi, and Yang Song**, "What do mutual fund investors really care about?," *Fisher College of Business Working Paper*, 2019, (2019-03), 005.

**Benjamin, Daniel J**, "Errors in probabilistic reasoning and judgment biases," Technical Report, National Bureau of Economic Research 2018.

**Bergstresser, Daniel, John MR Chalmers, and Peter Tufano**, "Assessing the costs and benefits of brokers in the mutual fund industry," *Review of financial studies*, 2009, *22* (10), 4129–4156.

**Berk, Jonathan B and Jules H van Binsbergen**, "Measuring skill in the mutual fund industry," *Journal of Financial Economics*, 2015, *118* (1), 1–20.

_ **and Richard C Green**, "Mutual fund flows and performance in rational markets," *Journal of political economy*, 2004, *112* (6), 1269–1295.

**Bordalo, Pedro, Nicola Gennaioli, and Andrei Shleifer**, "Diagnostic expectations and credit cycles," *The Journal of Finance*, 2018, *73* (1), 199–227.

_ , _ , **Rafael La Porta, and Andrei Shleifer**, "Diagnostic expectations and stock returns," *The Journal of Finance*, 2019, *74* (6), 2839–2874.

_ , _ , **Yueran Ma, and Andrei Shleifer**, "Over-reaction in macroeconomic expectations," *American Economic Review*, Forthcoming.

**Brown, Keith C, W Van Harlow, and Laura T Starks**, "Of tournaments and temptations: An analysis of managerial incentives in the mutual fund industry," *The Journal of Finance*, 1996, *51* (1), 85–110.

**Carhart, Mark M**, "On persistence in mutual fund performance," *The Journal of finance*, 1997, *52* (1), 57–82.

**Chen, Joseph, Harrison Hong, Ming Huang, and Jeffrey D Kubik**, "Does fund size erode mutual fund performance? The role of liquidity and organization," *The American Economic Review*, 2004, *94* (5), 1276–1302.

**Chevalier, Judith and Glenn Ellison**, "Risk taking by mutual funds as a response to incentives," *Journal of Political Economy*, 1997, *105* (6), 1167–1200.

**Choi, James J and Adriana Z Robertson**, "What Matters to Individual Investors? Evidence from the Horse's Mouth," *The Journal of Finance*, 2020, *75* (4), 1965–2020.

**Christoffersen, Susan EK, David K Musto, and Russ Wermers**, "Investor flows to asset managers: Causes and consequences," *Annu. Rev. Financ. Econ.*, 2014, *6* (1), 289–310.

**Cochrane, John H.**, "Finance: Function Matters, Not Size," *Journal of Economic Perspectives*, May 2013, *27* (2), 29–50.

**Elton, Edwin J, Martin J Gruber, and Christopher R Blake**, "Does mutual fund size matter? The relationship between size and performance," *The Review of Asset Pricing Studies*, 2012, *2* (1), 31–55.

**Evans, Richard B**, "Mutual fund incubation," *The Journal of Finance*, 2010, *65* (4), 1581–1611.

**Fama, Eugene F and Kenneth R French**, "Luck versus skill in the cross-section of mutual fund returns," *The journal of finance*, 2010, *65* (5), 1915–1947.

_ **and** _ , "Choosing factors," *Journal of Financial Economics*, 2018, *128* (2), 234–252.

**Ferreira, Miguel A, Aneel Keswani, António F Miguel, and Sofia B Ramos**, "The determinants of mutual fund performance: A cross-country study," *Review of Finance*, 2013, *17* (2), 483–525.

**Gabaix, Xavier**, "Behavioral inattention," Technical Report, National Bureau of Economic Research 2017.

**Gallaher, Steven, Ron Kaniel, and Laura T Starks**, "Madison Avenue meets Wall Street: Mutual fund families, competition and advertising," *Working paper*, 2006.

**Gennaioli, Nicola and Andrei Shleifer**, "What comes to mind," *The Quarterly journal of economics*, 2010, *125* (4), 1399–1433.

**Giglio, Stefano, Matteo Maggiori, Johannes Stroebel, and Stephen Utkus**, "Five facts about beliefs and portfolios," Technical Report, National Bureau of Economic Research 2019.

**Glode, Vincent, Burton Hollifield, Marcin Kacperczyk, and Shimon Kogan**, "Is investor rationality time varying? Evidence from the mutual fund industry," in "Behavioral Finance: WHERE DO INVESTORS'BIASES COME FROM?," World Scientific, 2017, pp. 67–113.

**Greenwood, Robin and Stefan Nagel**, "Inexperienced investors and bubbles," *Journal of Financial Economics*, 2009, *93* (2), 239–258.

**Hortaçsu, Ali and Chad Syverson**, "Product differentiation, search costs, and competition in the mutual fund industry: A case study of S&P 500 index funds," *The Quarterly Journal of Economics*, 2004, *119* (2), 403–456.

**Huang, Jennifer C, Kelsey D Wei, and Hong Yan**, "Investor learning and mutual fund flows," *Working paper*, 2012.

**Ippolito, Richard A**, "Consumer reaction to measures of poor quality: Evidence from the mutual fund industry," *The Journal of Law and Economics*, 1992, *35* (1), 45–70.

**Kahneman, Daniel and Amos Tversky**, "Subjective probability: A judgment of representativeness," *Cognitive psychology*, 1972a, *3* (3), 430–454.

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

**Pástor, Luboš and Pietro Veronesi**, "Learning in financial markets," *Annual Review of Financial Economics*, 2009, *1* (1), 361–381.

__ **and Robert F Stambaugh**, "Investing in equity mutual funds," *Journal of Financial Economics*, 2002, *63* (3), 351–380.

__ **and __** , "On the size of the active management industry," *Journal of Political Economy*, 2012, *120* (4), 740–781.

__ **, __ , and Lucian A Taylor**, "Scale and skill in active management," *Journal of Financial Economics*, 2015, *116* (1), 23–45.

**Pástor, L'uboš, Robert F Stambaugh, and Lucian A Taylor**, "Fund tradeoffs," *Journal of Financial Economics*, 2020.

**Rabin, Matthew**, "Inference by believers in the law of small numbers," *The Quarterly Journal of Economics*, 2002, *117* (3), 775–816.

__ **and Dimitri Vayanos**, "The gambler's and hot-hand fallacies: Theory and applications," *The Review of Economic Studies*, 2010, *77* (2), 730–778.

**Reuter, Jonathan and Eric Zitzewitz**, "How much does size erode mutual fund performance? A regression discontinuity approach," Technical Report, National Bureau of Economic Research 2015.

**Roussanov, Nikolai, Hongxun Ruan, and Yanhao Wei**, "Marketing mutual funds," *Working paper*, 2018.

**Sirri, Erik R. and Peter Tufano**, "Costly Search and Mutual Fund Flows," *The Journal of Finance*, 1998, *53* (5), 1589–1622.

**Song, Yang**, "The Mismatch Between Mutual Fund Scale and Skill," *The Journal of Finance,forthcoming*, 2019.

**Tversky, Amos and Daniel Kahneman**, "Judgment under uncertainty: Heuristics and biases," *science*, 1974, *185* (4157), 1124–1131.

_ **and** _ , "Extensional versus intuitive reasoning: The conjunction fallacy in probability judgment.," *Psychological review*, 1983, *90* (4), 293.

**Wooldridge, Jeffrey M**, *Econometric analysis of cross section and panel data*, MIT press, 2010.

**Yan, Xuemin**, "Liquidity, investment style, and the relation between fund size and fund performance," *Journal of Financial and Quantitative Analysis*, 2008, pp. 741–767.

Figure 1: Average annual net flow as a function of age



This figure plots the average annual net flow against fund age. The dotted lines indicate the 95% confidence intervals.

Figure 2: Belief weights on historical information

A posterior belief about a fund's skill is a weighted sum of the historical DRS-adjusted performances of the fund. The weight decays as the lag of historical performance grows. This plot displays how the weights decay in different sets of posterior beliefs. The blue curve with stars plots the weighting scheme under the estimated performance-based model (i.e., how the rational posterior $\hat{a}_{j,t}$ weights historical information). The yellow curve with circles plots the weighting scheme under the estimated investor beliefs based on fund flows (i.e., how the $\tilde{a}_{j,t}$ weights historical information). The red curve with squares plots the same investor's weighting scheme as the yellow curve except that we impose $\tilde{\rho} = 0.95$ that is close to that under the rational belief (the estimated $\tilde{\rho} = 0.76$). The gray curve with asterisks plots the same investor's weighting scheme as the yellow curve except that we impose $\lambda = 0.45$, which is in the middle between that under the rational belief (0.203) and the estimated $\lambda = 0.69$ from the flow model (see Table 2).

Table 1: Summary statistics

|  | Mean | SD | P1 | P25 | P50 | P75 | P99 |
|---|---|---|---|---|---|---|---|
| Annual alpha (%) | 0.33 | 7.33 | -17.01 | -3.27 | -0.04 | 3.19 | 23.93 |
| Annual expense ratio (%) | 1.18 | 0.50 | 0.14 | 0.90 | 1.14 | 1.43 | 2.55 |
| Fund size ($million) | 1,453 | 5,603 | 9 | 73 | 238 | 840 | 23,466 |

This table presents summary statistics for our sample of U.S. equity mutual funds. The sample period is from 1965 to 2014. Each observation is a fund-year combination. Annual alpha is computed using Fama-French six-factor model as in Fama and French (2018). Fund size is the fund's total AUM aggregated across share classes. Both annual expense ratio and annual alpha are computed as AUM-weighted averages across share classes.

Table 2: Parameter estimates

| | Value | SE | Description |
|---|---|---|---|
| **Performance Model** | | | |
| $\eta$ | 2.2e-3 | (0.001) | Size of DRS |
| $\gamma$ | 2.9e-4 | (8.7e-2) | Shape of DRS |
| $\mu$ | 0.014 | (0.003) | Mean of skill prior |
| $\kappa$ | 0.014 | (5.2e-4) | Stdv. of skill prior |
| $\delta$ | 0.069 | (2.0e-4) | Stdv. of return noise |
| $\rho$ | 0.954 | (0.014) | Skill persistence |
| **Flow Model** | | | |
| $\widetilde{\eta}$ | 0.007 | (0.001) | Size of DRS |
| $\widetilde{\gamma}$ | 0.030 | (0.003) | Shape of DRS |
| $\widetilde{\mu}$ | 0.063 | (0.005) | Mean of skill prior |
| $\lambda := \widetilde{\kappa}/\widetilde{\delta}$ | 0.694 | (0.035) | Ratio of prior and noise stdv. |
| $\widetilde{\rho}$ | 0.762 | (0.015) | Skill persistence |
| $\phi$ | 0.070 | (0.005) | Flow adjustment rate |
| $\beta$ | 0.347 | (0.005) | Serial corr. in flow residual |
| $\omega$ | 0.340 | (0.001) | Stdv. of flow residual |

This table reports the maximum likelihood estimates for our model. The upper part reports parameters in the model of fund performance; the lower part reports parameters in the model of fund flows. For more details about the definitions of the parameters, please refer to Section 2. The standard errors are in parentheses.

Table 3: Sensitivity of net performance to misallocation

| Net performance ($r_{j,t} - p_{j,t}$) | | | | |
|---|---|---|---|---|
| | (1) | (2) | (3) | (4) |
| $q_{j,t} - \widehat{q}_{j,t}^{BG}$ | -0.18*** | | | |
| | (3.53) | | | |
| $\widehat{q}_{j,t}^{BG}$ | | 0.18*** | | |
| | | (3.51) | | |
| $q_{j,t}$ | | -0.17** | | -0.15* |
| | | (2.55) | | (1.71) |
| $q_{j,t} - \widetilde{q}_{j,t}^{BG}$ | | | -0.22* | |
| | | | (1.96) | |
| $\widetilde{q}_{j,t}^{BG}$ | | | | 0.24* |
| | | | | (1.87) |
| Constant | -0.17 | -0.23 | -0.96* | -1.49** |
| | (0.27) | (0.38) | (1.79) | (2.47) |
| N | 25,530 | 25,530 | 25,530 | 25,530 |
| Adj $R^2$ | 0.007 | 0.007 | 0.005 | 0.006 |

This table reports the regressions of net performance on misallocation. We regress net performance $r_{j,t} - p_{j,t}$ (in percentage) onto misallocation in period $t$. For the definitions of misallocations, please see Section 6.1. We double-cluster standard errors by Morningstar category and year. The $t$-statistics are in parentheses. Significance at the 1%, 5%, and 10% levels are indicated by ***, **, and *, respectively.

Table 4: Misallocation and net performance

| | Too small | 2 | 3 | 4 | 5 | 6 | 7 | 8 | 9 | Too big |
|---|---|---|---|---|---|---|---|---|---|---|
| **Panel A: Equal weighted portfolio** | | | | | | | | | | |
| **Misallocation decile** | | | | | | | | | | |
| FF5 alpha | 1.84 | 1.00 | 0.33 | -0.14 | -0.17 | -0.34 | -0.69 | -1.07 | -1.22 | -1.21 |
| $t$-stat | 2.12 | 1.71 | 0.66 | -0.38 | -0.37 | -0.84 | -1.65 | -2.55 | -1.94 | -1.83 |
| FF6 alpha | 1.10 | 0.52 | -0.25 | -0.37 | -0.34 | -0.71 | -0.86 | -0.81 | -1.61 | -2.60 |
| $t$-stat | 1.90 | 0.95 | -0.68 | -0.89 | -0.83 | -1.67 | -2.22 | -1.30 | -3.03 | -4.16 |
| **Panel B: Value weighted portfolio** | | | | | | | | | | |
| **Misallocation decile** | | | | | | | | | | |
| FF5 alpha | 1.24 | 0.22 | 0.16 | -0.60 | -0.19 | -0.96 | -0.96 | -1.17 | -1.39 | -1.31 |
| $t$-stat | 1.19 | 0.35 | 0.25 | -1.22 | -0.32 | -2.13 | -1.91 | -2.31 | -1.97 | -2.17 |
| FF6 alpha | 0.91 | -0.27 | -0.44 | -0.54 | -0.78 | -0.64 | -0.88 | -1.50 | -1.97 | -2.98 |
| $t$-stat | 1.00 | -0.52 | -0.99 | -1.07 | -1.67 | -1.27 | -1.92 | -2.48 | -3.68 | -4.60 |

This table reports the net performance $(r_{j,t} - p_{j,t})$ for each of the ten decile portfolios formed on the misallocation measure. For the definitions of misallocations, please see Section 6.1. We calculate the expanding window alpha of each portfolio using the Fama-French five-factor model and the Fama-French six-factor model. Net performance is reported in percentage. For each calendar year, we sort all mutual funds into decile portfolios based on a fund's misallocation measure at the beginning of the year. Then we compute the average portfolio annual net alphas in the year. Panel A reports the results for equal weighted portfolios and Panel B reports the results for value weighted portfolios. In each panel, portfolio "Too small" collects the funds with the lowest decile of misallocation in a given year; portfolio "Too big" collects the funds with the highest decile of misallocation in a given year. The $t$-statistics are based on the heteroskedasticity-consistent standard errors of White (1980).

Table 5: Sensitivity of flows to misallocation

|  | (1) | (2) | (3) | (4) |
|---|---|---|---|---|
| $q_{j,t} - \widehat{q}_{j,t}^{BG}$ | -0.01*** | -0.01*** |  |  |
|  | (6.42) | (4.38) |  |  |
| $q_{j,t} - \widetilde{q}_{j,t}^{BG}$ |  |  | -0.03*** | -0.03*** |
|  |  |  | (7.55) | (6.78) |
| Lag expense ratio |  | 2.54 |  | 1.32 |
|  |  | (1.56) |  | (1.01) |
| Lag load dummy |  | 0.01 |  | 0.01 |
|  |  | (1.61) |  | (1.45) |
| Lag flow |  | 0.14*** |  | 0.13*** |
|  |  | (4.55) |  | (3.99) |
| Lag annual alpha vol |  | -1.04** |  | -1.20* |
|  |  | (2.27) |  | (1.91) |
| Lag log fundsize |  | -0.02*** |  | -0.01* |
|  |  | (4.27) |  | (1.77) |
| Lag age |  | 2e-4 |  | 2e-4 |
|  |  | (0.40) |  | (0.43) |
| Constant | 0.07*** | 0.17*** | 5e-3 | 0.08 |
|  | (5.25) | (3.26) | (0.52) | (1.19) |
| N | 23,018 | 20,716 | 23,018 | 20,716 |
| Adj $R^2$ | 0.01 | 0.05 | 0.05 | 0.07 |

This table reports the regression of flows on misallocation. Flow is defined as $q_{j,t+1} - q_{j,t}$. For the definitions of misallocations, please see Section 6.1. The control variables include: lag expense ratio, lag load fund dummy (takes the value of 1 if the fund has front loads), lag flow, lag annual alpha vol measured as a fund's alpha's standard deviation over the prior year using monthly data, lag log of fund AUM, and lag fund age measured in years. We double-cluster standard errors by Morningstar category and year. The *t*-statistics are in parentheses. Significance at the 1%, 5%, and 10% levels are indicated by ***, **, and *, respectively.

Table 6: Out-of-sample MSE of various predictors for fund future performance

| | (1) | (2) | (3) | (4) | (5) | (6) |
| | Rational posterior | Investor's posterior | Lag 1 year alpha | Lag 3 year alpha | Lag 5 year alpha | Lag 10 year alpha |
|---|---|---|---|---|---|---|
| MSE (%) | 0.22 | 0.25 | 0.56 | 0.32 | 0.31 | 0.29 |

This table presents the mean squared errors of various predictors for realized alpha in 2010-2014. The first column uses the skill posteriors in the performance model, with the model parameters estimated from the performance data in 1965-2009. The second column uses the skill posteriors in the flow model, with the model parameters estimated from the fund size data in 1965-2009. For both columns, DRS is subtracted from the posterior for the current-period skill. The last four columns use moving averages. For the exercises associated with this table, we use expanding-window alphas as our measure of fund performance, to ensure that pre-2010 performance measures do not use any information from 2010-2014.

Table 7: Parameter estimates: different share classes

| | Inst Share | | Retail Share | | | Benchmark | |
|---|---|---|---|---|---|---|---|
| | Value | SE | Value | SE | | Value | SE |
| $\widetilde{\eta}$ | 3.5e-3 | (0.001) | 6.2e-3 | (0.001) | $\eta$ | 2.2e-3 | (0.001) |
| $\widetilde{\gamma}$ | 0.051 | (0.015) | 0.025 | (0.002) | $\gamma$ | 2.9e-4 | (8.7e-4) |
| $\widetilde{\mu}$ | 0.033 | (0.004) | 0.054 | (0.003) | $\mu$ | 0.014 | (0.003) |
| $\lambda = \widetilde{\kappa}/\widetilde{\delta}$ | 0.429 | (0.044) | 0.705 | (0.028) | $\kappa/\delta$ | 0.203 | (0.008) |
| $\widetilde{\rho}$ | 0.737 | (0.025) | 0.700 | (0.014) | $\rho$ | 0.954 | (0.014) |
| $\phi$ | 0.077 | (0.005) | 0.064 | (0.004) | | | |
| $\beta$ | 0.213 | (0.006) | 0.303 | (0.005) | | | |
| $\omega$ | 0.535 | (0.002) | 0.352 | (0.001) | | | |
| $\psi^{Inst.}$ | 0.424 | (0.054) | | | | | |

This table reports parameter estimates for the extended flow model that accounts for institutional share class vs. retail share class. The sample period is from 2000 to 2014. The last two columns reproduce the rational benchmark estimates from the top panel of Table 2. The standard errors are in parentheses.

# Appendix

## A    Monte Carlo study

We conduct Monte Carlo experiments on our partial MLE described in Section 3. The first step is to simulate panel data on funds' return $r_{j,t}$, prices $p_{j,t}$, and size $q_{j,t}$. Our model specifies the data generating process for returns and sizes but not prices. Due to the partial MLE approach, we can generate $p_{j,t}$ as any function of $\{Y_{j,t-1}, Y_{j,t-2}, ...\}$. For the results shown below, we generate $p_{j,t}$ from a simple AR(1) process. In the simulated panel, we retain the same fund identities and years of existence for each fund as in the real data. As a result, the simulated panel data is of the same size as the real data.

The "true" parameters with which we generate the simulated data are set to the values in Table 2. We simulated 100 datasets, and for each dataset, we apply the partial MLE to recover the parameters. Table S1 shows the means and standard errors of the means of the recovered parameter values across the 100 datasets. As we can see, our partial MLE can recover all the parameters.

[Table S1 about here.]

## B    Diagnostic expectations

Bordalo et al. (2019) develop a filtering rule they term "diagnostic expectations," which deviates from the standard Kalman filter in a particular way that embodies the representativeness heuristic of Kahneman and Tversky (1972a). Here, we explore the way that such "diagnostic" beliefs could be inferred from mutual fund flows. Incorporating the diagnostic Kalman filter of Bordalo et al. (2019) into our model of investor beliefs is straightforward. Specifically, we replace the posterior, $\widetilde{a}_{j,t}$, that enters the flow equations (6) - (8) with the diagnostic posterior, $\widetilde{a}_{j,t}^{\theta}$, computed as follows:

$$\widetilde{a}_{j,t+1}^{\theta} = \widetilde{\rho}\left[\widetilde{a}_{j,t} + (1+\theta)\frac{\widetilde{\sigma}_{j,t}^2}{\widetilde{\sigma}_{j,t}^2 + \widetilde{\delta}^2}\left(r_{j,t} + D(Q_{j,t}; \widetilde{\eta}, \widetilde{\gamma}) - \widetilde{a}_{j,t}\right)\right] + (1-\widetilde{\rho})\widetilde{\mu},$$

where $\widetilde{a}_{j,t}$ on the right side still follows Kalman filter (equation 9), and the new parameter $\theta \geq 0$ captures the degree of representativeness. The larger $\theta$ is, the more extra weight that the current-period return, $r_{j,t}$, has on the next-period belief $\widetilde{a}_{j,t+1}^{\theta}$. When $\theta = 0$, the diagnostic filter reduces to the standard Kalman filter.

Table S2 compares the estimation results under the diagnostic expectations and our baseline model. With the diagnostic expectations, there are three parameters that directly affect investors' weighting scheme in our model: $\theta$, $\lambda$, and $\widetilde{\rho}$. To examine their respective roles, we

first estimate three restrictive models. Column 1 reproduces our baseline estimates, where $\theta$ is restricted to be zero. Column 2 displays the estimates while fixing $\widetilde{\rho} = 1$, as in the original BG model. There are two interesting observations. First, the estimate of $\theta$ is significantly positive, implying recency bias by investors. Second, the estimate of $\lambda$ decreases from 0.69 in the baseline model to 0.19 which is much closer to the rational benchmark estimated from the performance model. Intuitively, this result indicates a single parameter $\theta > 0$ is able to capture both: (i) the recency bias as captured by a small $\widetilde{\rho}$, and (ii) prior or "base rate" neglect as captured by a large $\lambda$. To further verify this intuition, in column 3, we estimate a model that fixes $\lambda$ at the rational benchmark. Indeed, we see that the estimate of $\widetilde{\rho}$ moves closer to the rational benchmark while $\theta$ stays significantly positive. In column 4, we estimate all the parameters together; the estimate of $\theta$ is zero, implying that our baseline model is able to capture the representativeness of investors implicit in the observed mutual fund flows. The resulting likelihood is substantially higher than the restrictive models in columns 2 and 3.

To summarize the results, our baseline model is consistent with diagnostic expectations in terms of capturing the return chasing behavior of investors. Diagnostic expectations are more efficient, capturing both recency bias and base rate neglect with a single parameter. Our model offers a somewhat more flexible specification capturing the two elements separately.

[Table S2 about here.]

## C    Robustness

### C.1    Expanding-window alphas

In this section, we show that our main results are robust to using expanding-window alphas as our measure of fund performance. More specifically, we fix the starting point of the window at the birth time of the fund and the endpoint of the window progresses along time. The initial size of the window is 24 months. Then, from the 25th month onwards, we expand the window monthly to compute the betas. Last, we aggregate monthly alpha into annual alpha. Note the same expanding-window alphas are used for our out-of-sample prediction exercise, Section 6.4.

Table S3 displays the parameter estimates of both performance and flow models (to compare with Table 2). The first set of columns reproduces the estimates in the main text. We see that the estimates are similar across the two sets of columns.

Table S4 displays the regressions of net performance on misallocation (to compare with Table 3). We see that the results do not change qualitatively; misallocation still predicts net performance.

Table S5 displays the regressions of flow on misallocation (to compare with Table 5). We see that the results do not change qualitatively; flow still responds to misallocation but sluggishly.

[Table S3 about here.]

[Table S4 about here.]

[Table S5 about here.]

## C.2 Additional controls in the flow misallocation regression

In Table 5, we regress flows onto the misallocation measure and other controls. We find that the magnitudes of the coefficient in front of the misallocation measure is significantly *smaller* than is predicted by the theory. In this subsection, as robustness checks, we investigate the impact of these additional variables on the flow-misallocation relationship.

First, Gallaher et al. (2006) show that advertising matters for flows as the fund families that spend more on advertising tend to attract more flows. More specifically, they show that advertising has a significant positive effect on family flows for relatively high advertisers. We acquire fund family level advertising spending data from Kantar Media. We create a dummy variable taking the value of one if the fund family's advertising spending is above the yearly average, zero otherwise. We include this dummy in the flow-misallocation regression as reported in Table S6, column (2). For the ease of comparison, we reproduce our baseline results from Table 5 in Table S6, column (1). We find that the coefficient in front of the dummy is positive and statistically significant, which confirms the findings of Gallaher et al. (2006). Meanwhile, the coefficient in front of the misallocation measure in column (1) and column (2), is quantitatively unchanged.

In our baseline analysis, we assume that only fund's own performance affects investor flows. However, it is possible that a fund's competitors' performance also drives flows into and out of the given fund (e.g., Brown et al. (1996) and Sirri and Tufano (1998) have shown that *relative* performance matters for fund flows as the winners tend to be rewarded with significant inflows and the losers are punished with outflows; see Christoffersen et al., 2014 for a survey). In order to account for this possibility we use the average alpha of other funds in the same Morningstar category as a measure of a fund's competitors' performance and we create a dummy variable taking the value of one if the given fund's alpha is below the average alpha of other funds in the same Morningstar category, zero otherwise. We include this dummy in our flow-misallocation regression. The results are provided in Table S6, column (3). We find that the coefficient in front of the dummy is negative and statistically significant, consistent with the existing evidence. Importantly, the coefficient in front of the misallocation measure is quantitatively unchanged. As a further robustness check, we replace the simple average category alpha with the AUM-weighted average category alpha, and the results are qualitatively similar.

Lastly, a recent paper by Ben-David et al. (2019) finds that fund flow data are most consistent with investors relying on Morningstar ratings, which are strongly influenced by recent fund returns. Inspired by this interesting finding, we include the fund's average Morningstar ratings and the fund's gross return over the past year into the flow-misallocation regression. The results are provided in Table S6, column (4). We find that the coefficients in front of both the Morningstar ratings and the fund past returns are positive and statistically significant. The adjusted R-squared increases from 5% in the baseline results (column 1) to 8%. These results confirm the findings in Ben-David et al. (2019) that Morningstar ratings and past returns are powerful drivers of fund flows. Meanwhile, the coefficient in front of the misallocation measure stays statistically significant. It indicates that fund flows still respond to the misallocation measure even after controlling for fund ratings and past returns. As a further robustness check, we replace the fund's average Morningstar ratings with the fund's max or min Morningstar ratings over the past year. The results are qualitatively similar.

In Table S7, we repeat the above analysis but replace $\widehat{q}_{j,t}^{BG}$ with $\widetilde{q}_{j,t}^{BG}$, which is the efficient fund size computed using investor's belief (see equation 6). The results are not qualitatively different from Table S6.

[Table S6 about here.]

[Table S7 about here.]

41

Table S1: Monte Carlo results

| | "True" value | Estimate mean | S.E. for mean |
|---|---|---|---|
| Performance model: | | | |
| $\eta$ | 0.0022 | 0.0022 | (0.0003) |
| $\gamma$ | 0.0003 | 0.0047 | (0.0043) |
| $\mu$ | 0.0138 | 0.0140 | (0.0009) |
| $\kappa$ | 0.0142 | 0.0142 | (0.0009) |
| $\delta$ | 0.0694 | 0.0693 | (0.0004) |
| $\rho$ | 0.9544 | 0.9528 | (0.0173) |
| Flow model: | | | |
| $\widetilde{\eta}$ | 0.0068 | 0.0068 | (0.0003) |
| $\widetilde{\gamma}$ | 0.0301 | 0.0303 | (0.0023) |
| $\widetilde{\mu}$ | 0.0629 | 0.0628 | (0.0017) |
| $\lambda = \widetilde{\kappa}/\widetilde{\delta}$ | 0.6942 | 0.6947 | (0.0136) |
| $\widetilde{\rho}$ | 0.7625 | 0.7622 | (0.0130) |
| $\phi$ | 0.0696 | 0.0698 | (0.0030) |
| $\beta$ | 0.3465 | 0.3479 | (0.0076) |
| $\omega$ | 0.3403 | 0.3405 | (0.0018) |

This table reports the Monte Carlo results. The "true" parameters with which we generate the simulated datasets are set to the values in Table 2. We simulate 100 datasets, and for each dataset, we apply the partial MLE to estimate the parameters. This table reports the means and standard errors of the means of the parameter estimates across the 100 datasets.

Table S2: Parameter estimates with diagnostic expectations

| | (1) | | (2) | | (3) | | (4) | |
|---|---|---|---|---|---|---|---|---|
| | Value | SE | Value | SE | Value | SE | Value | SE |
| $\widetilde{\eta}$ | 0.007 | (0.001) | 0.005 | (4e-4) | 0.003 | (2.0e-4) | 0.007 | (0.001) |
| $\widetilde{\gamma}$ | 0.030 | (0.003) | 0.030 | (0.003) | 0.010 | (0.002) | 0.030 | (0.003) |
| $\widetilde{\mu}$ | 0.063 | (0.005) | 0.049 | (0.003) | 0.034 | (0.001) | 0.063 | (0.005) |
| $\lambda = \widetilde{\kappa}/\widetilde{\delta}$ | 0.694 | (0.035) | 0.190 | (0.010) | 0.203 | fixed | 0.694 | (0.042) |
| $\widetilde{\rho}$ | 0.762 | (0.015) | 1.000 | fixed | 0.862 | (0.014) | 0.762 | (0.018) |
| $\phi$ | 0.070 | (0.005) | 0.063 | (0.004) | 0.057 | (0.003) | 0.070 | (0.005) |
| $\beta$ | 0.347 | (0.005) | 0.363 | (0.004) | 0.355 | (0.004) | 0.347 | (0.005) |
| $\omega$ | 0.340 | (0.001) | 0.349 | (0.001) | 0.345 | (0.001) | 0.340 | (0.001) |
| $\theta$ | - | - | 4.282 | (0.308) | 1.928 | (0.086) | 0.000 | (0.034) |
| Log likelihood | 24,389 | | 23,924 | | 24,091 | | 24,389 | |

This table reports parameter estimates for the model of fund flows that accounts for diagnostic expectations. Column 1 reproduces our baseline model estimates (i.e., $\theta$ fixed at zero). Column 2 estimates a model with diagnostic expectations while fixing $\widetilde{\rho}$ at 1. Column 3 estimates a model with diagnostic expectations while fixing $\lambda$ at the rational benchmark, $\kappa/\delta$. Column 4 estimates a model with diagnostic expectations without fixing any of the parameters. Standard errors are in parentheses. For more details about the definitions of the parameters, please refer to Section 2.

Table S3: Parameter estimates; expanding-window alphas

| | Constant betas | | Expanding betas | | Description |
|---|---|---|---|---|---|
| | Value | SE | Value | SE | |
| **Perf. model:** | | | | | |
| $\eta$ | 2.2e-3 | (0.001) | 1.8e-3 | (0.001) | size of DRS |
| $\gamma$ | 2.9e-4 | (8.7e-4) | 1.0e-4 | (1.1e-3) | shape of DRS |
| $\mu$ | 0.014 | (0.003) | 0.011 | (0.003) | mean of skill prior |
| $\kappa$ | 0.014 | (5.2e-4) | 0.014 | (6.1e-4) | stdv. of skill prior |
| $\delta$ | 0.069 | (2.0e-4) | 0.076 | (2.3e-4) | stdv. of return noise |
| $\rho$ | 0.954 | (0.014) | 0.943 | (0.017) | skill persistence |
| **Flow model:** | | | | | |
| $\widetilde{\eta}$ | 0.007 | (0.001) | 0.005 | (4.5e-4) | size of DRS |
| $\widetilde{\gamma}$ | 0.030 | (0.003) | 0.024 | (0.003) | shape of DRS |
| $\widetilde{\mu}$ | 0.063 | (0.005) | 0.049 | (0.003) | mean of skill prior |
| $\lambda = \widetilde{\kappa}/\widetilde{\delta}$ | 0.694 | (0.035) | 0.544 | (0.023) | ratio of prior and noise stdv. |
| $\widetilde{\rho}$ | 0.762 | (0.015) | 0.697 | (0.013) | skill persistence |
| $\phi$ | 0.070 | (0.005) | 0.057 | (0.004) | flow adjustment rate |
| $\beta$ | 0.347 | (0.005) | 0.348 | (0.004) | serial corr. in flow residual |
| $\omega$ | 0.340 | (0.001) | 0.348 | (0.001) | stdv. of flow residual |

This table replicates the results in Table 2 with expanding-window alphas. The first set of columns reproduces the estimates in Table 2. Standard errors are in parentheses.

Table S4: Sensitivity of net performance to misallocation; expanding-window alphas

| Net performance ($r_{j,t} - p_{j,t}$) | | | | |
|---|---|---|---|---|
| | (1) | (2) | (3) | (4) |
| $q_{j,t} - \widehat{q}_{j,t}^{BG}$ | -0.15*** | | | |
| | (2.63) | | | |
| $\widehat{q}_{j,t}^{BG}$ | | 0.15*** | | |
| | | (2.63) | | |
| $q_{j,t}$ | | -0.15** | | -0.10 |
| | | (2.01) | | (1.18) |
| $q_{j,t} - \widetilde{q}_{j,t}^{BG}$ | | | -0.17 | |
| | | | (1.32) | |
| $\widetilde{q}_{j,t}^{BG}$ | | | | 0.19 |
| | | | | (1.27) |
| Constant | -0.18 | -0.20 | -1.05* | -1.59** |
| | (0.24) | (0.28) | (1.87) | (2.09) |
| N | 25,530 | 25,530 | 25,530 | 25,530 |
| Adj $R^2$ | 0.005 | 0.005 | 0.003 | 0.003 |

This table replicates Table 3 with expanding-window alphas. The numbers in parentheses are $t$-statistics.

Table S5: Sensitivity of flows to misallocation; expanding-window alphas

| | (1) | (2) | (3) | (4) |
|---|---|---|---|---|
| $q_{j,t} - \widehat{q}_{j,t}^{BG}$ | -0.01*** | -0.01*** | | |
| | (7.16) | (4.73) | | |
| $q_{j,t} - \widetilde{q}_{j,t}^{BG}$ | | | -0.03*** | -0.02*** |
| | | | (8.92) | (6.22) |
| Lag expense ratio | | 2.52 | | 1.62 |
| | | (1.52) | | (1.27) |
| Lag load dummy | | 0.02** | | 0.01* |
| | | (2.06) | | (1.82) |
| Lag flow | | 0.15*** | | 0.14*** |
| | | (4.70) | | (4.27) |
| Lag annual alpha vol | | -0.97** | | -1.04** |
| | | (2.19) | | (1.98) |
| Lag log fundsize | | -0.02*** | | -0.01* |
| | | (4.48) | | (1.83) |
| Lag age | | 2e-4 | | 2e-4 |
| | | (0.43) | | (0.48) |
| Constant | 0.07*** | 0.18*** | 6e-3 | 0.07 |
| | (5.49) | (3.48) | (0.61) | (1.08) |
| N | 23,018 | 20,716 | 23,018 | 20,716 |
| Adj $R^2$ | 0.01 | 0.05 | 0.04 | 0.06 |

This table replicates Table 5 with expanding-window alphas. The numbers in parentheses are $t$-statistics.

Table S6: Sensitivity of flows to misallocation; Additional controls

| | (1) | (2) | (3) | (4) |
|---|---|---|---|---|
| $q_{j,t} - \widehat{q}_{j,t}^{BG}$ | -0.01*** | -0.01*** | -0.01** | -0.01** |
| | (4.38) | (4.43) | (2.34) | (2.16) |
| Lag expense ratio | 2.54 | 2.69 | 0.27 | 1.19 |
| | (1.56) | (1.64) | (0.16) | (0.66) |
| Lag load dummy | 0.01 | 0.01 | 0.01 | 0.03*** |
| | (1.61) | (1.61) | (1.35) | (4.01) |
| Lag flow | 0.14*** | 0.14*** | 0.15*** | 0.12*** |
| | (4.55) | (4.53) | (4.69) | (4.01) |
| Lag annual alpha vol | -1.04** | -1.05** | -0.92** | -0.46 |
| | (2.27) | (2.30) | (2.17) | (1.07) |
| Lag log fundsize | -0.02*** | -0.02*** | -0.02*** | -0.03*** |
| | (4.27) | (4.35) | (4.78) | (5.56) |
| Lag age | 2e-4 | 2e-4 | 2e-4 | 1e-3 |
| | (0.40) | (0.41) | (0.38) | (1.18) |
| Lag dummy above ave fund family ad spending | | 0.02** | | |
| | | (2.01) | | |
| Lag dummy below ave perf. in same MS style | | | -0.12*** | |
| | | | (13.10) | |
| Lag MS ratings | | | | 0.05*** |
| | | | | (6.21) |
| Lag annual gross ret | | | | 0.19*** |
| | | | | (2.99) |
| Constant | 0.17*** | 0.17*** | 0.24*** | -0.01 |
| | (3.26) | (3.26) | (4.50) | (0.20) |
| N | 20,716 | 20,716 | 20,716 | 20,692 |
| Adj $R^2$ | 0.05 | 0.05 | 0.08 | 0.08 |

This table replicates Table 5 with additional control variables. The numbers in parentheses are $t$-statistics.

Table S7: Sensitivity of flows to (investor) misallocation; Additional controls

| | (1) | (2) | (3) | (4) |
|---|---|---|---|---|
| $q_{j,t} - \widetilde{q}_{j,t}^{BG}$ | -0.03*** | -0.03*** | -0.02*** | -0.02*** |
| | (6.78) | (6.84) | (4.39) | (4.83) |
| Lag expense ratio | 1.32 | 1.44 | 0.04 | 0.91 |
| | (1.01) | (1.09) | (0.03) | (0.67) |
| Lag load dummy | 0.01 | 0.01 | 0.01 | 0.03*** |
| | (1.45) | (1.46) | (1.32) | (3.70) |
| Lag flow | 0.13*** | 0.13*** | 0.14*** | 0.11*** |
| | (3.99) | (3.97) | (4.24) | (3.68) |
| Lag annual alpha vol | -1.20* | -1.21* | -1.05* | -0.73 |
| | (1.91) | (1.94) | (1.95) | (1.36) |
| Lag log fundsize | -0.01* | -0.01* | -0.01*** | -0.02*** |
| | (1.77) | (1.85) | (2.73) | (3.20) |
| Lag age | 2e-4 | 2e-4 | 2e-4 | 1e-3 |
| | (0.43) | (0.44) | (0.40) | (1.08) |
| Lag dummy above ave fund family ad spending | | 0.02** | | |
| | | (2.04) | | |
| Lag dummy below ave perf. in same MS style | | | -0.09*** | |
| | | | (9.16) | |
| Lag MS ratings | | | | 0.04*** |
| | | | | (5.30) |
| Lag annual gross ret | | | | 0.14** |
| | | | | (2.34) |
| Constant | 0.08 | 0.08 | 0.16** | -0.05 |
| | (1.19) | (1.18) | (2.52) | (0.81) |
| N | 20,716 | 20,716 | 20,716 | 20,692 |
| Adj $R^2$ | 0.07 | 0.07 | 0.09 | 0.09 |

This table replicates Table 5 with additional control variables. The numbers in parentheses are $t$-statistics.