

A SURVEY OF BEHAVIORAL FINANCE*

NICHOLAS BARBERIS

University of Chicago

RICHARD THALER

University of Chicago

Contents

Abstract	1054
Keywords	1054
1. Introduction	1055
2. Limits to arbitrage	1056
2.1. Market efficiency	1056
2.2. Theory	1058
2.3. Evidence	1061
2.3.1. Twin shares	1061
2.3.2. Index inclusions	1063
2.3.3. Internet carve-outs	1064
3. Psychology	1065
3.1. Beliefs	1065
3.2. Preferences	1069
3.2.1. Prospect theory	1069
3.2.2. Ambiguity aversion	1074
4. Application: The aggregate stock market	1075
4.1. The equity premium puzzle	1078
4.1.1. Prospect theory	1079
4.1.2. Ambiguity aversion	1082
4.2. The volatility puzzle	1083
4.2.1. Beliefs	1084
4.2.2. Preferences	1086
5. Application: The cross-section of average returns	1087
5.1. Belief-based models	1092

* We are very grateful to Markus Brunnermeier, George Constantinides, Kent Daniel, Milt Harris, Ming Huang, Owen Lamont, Jay Ritter, Andrei Shleifer, Jeremy Stein and Tuomo Vuolteenaho for extensive comments.

5.2. Belief-based models with institutional frictions	1095
5.3. Preferences	1097
6. Application: Closed-end funds and comovement	1098
6.1. Closed-end funds	1098
6.2. Comovement	1099
7. Application: Investor behavior	1101
7.1. Insufficient diversification	1101
7.2. Naive diversification	1103
7.3. Excessive trading	1103
7.4. The selling decision	1104
7.5. The buying decision	1105
8. Application: Corporate finance	1106
8.1. Security issuance, capital structure and investment	1106
8.2. Dividends	1109
8.3. Models of managerial irrationality	1111
9. Conclusion	1113
Appendix A	1115
References	1116

Abstract

Behavioral finance argues that some financial phenomena can plausibly be understood using models in which some agents are not fully rational. The field has two building blocks: *limits to arbitrage*, which argues that it can be difficult for rational traders to undo the dislocations caused by less rational traders; and *psychology*, which catalogues the kinds of deviations from full rationality we might expect to see. We discuss these two topics, and then present a number of behavioral finance applications: to the aggregate stock market, to the cross-section of average returns, to individual trading behavior, and to corporate finance. We close by assessing progress in the field and speculating about its future course.

Keywords

behavioral finance, market efficiency, prospect theory, limits to arbitrage, investor psychology, investor behavior

JEL classification: G11, G12, G30

1. Introduction

The traditional finance paradigm, which underlies many of the other articles in this handbook, seeks to understand financial markets using models in which agents are “rational”. Rationality means two things. First, when they receive new information, agents update their beliefs correctly, in the manner described by Bayes’ law. Second, given their beliefs, agents make choices that are normatively acceptable, in the sense that they are consistent with Savage’s notion of Subjective Expected Utility (SEU).

This traditional framework is appealingly simple, and it would be very satisfying if its predictions were confirmed in the data. Unfortunately, after years of effort, it has become clear that basic facts about the aggregate stock market, the cross-section of average returns and individual trading behavior are not easily understood in this framework.

Behavioral finance is a new approach to financial markets that has emerged, at least in part, in response to the difficulties faced by the traditional paradigm. In broad terms, it argues that some financial phenomena can be better understood using models in which some agents are *not* fully rational. More specifically, it analyzes what happens when we relax one, or both, of the two tenets that underlie individual rationality. In some behavioral finance models, agents fail to update their beliefs correctly. In other models, agents apply Bayes’ law properly but make choices that are normatively questionable, in that they are incompatible with SEU.¹

This review essay evaluates recent work in this rapidly growing field. In Section 2, we consider the classic objection to behavioral finance, namely that even if some agents in the economy are less than fully rational, rational agents will prevent them from influencing security prices for very long, through a process known as arbitrage. One of the biggest successes of behavioral finance is a series of theoretical papers showing that in an economy where rational and irrational traders interact, irrationality *can* have a substantial and long-lived impact on prices. These papers, known as the literature on “limits to arbitrage”, form one of the two buildings blocks of behavioral finance.

¹ It is important to note that most models of asset pricing use the Rational Expectations Equilibrium framework (REE), which assumes not only individual rationality but also *consistent beliefs* [Sargent (1993)]. Consistent beliefs means that agents’ beliefs are correct: the subjective distribution they use to forecast future realizations of unknown variables is indeed the distribution that those realizations are drawn from. This requires not only that agents process new information correctly, but that they have *enough* information about the structure of the economy to be able to figure out the correct distribution for the variables of interest.

Behavioral finance departs from REE by relaxing the assumption of individual rationality. An alternative departure is to retain individual rationality but to relax the consistent beliefs assumption: while investors apply Bayes’ law correctly, they lack the information required to know the actual distribution variables are drawn from. This line of research is sometimes referred to as the literature on bounded rationality, or on structural uncertainty. For example, a model in which investors do not know the growth rate of an asset’s cash flows but learn it as best as they can from available data, would fall into this class. Although the literature we discuss also uses the term bounded rationality, the approach is quite different.

To make sharp predictions, behavioral models often need to specify the form of agents' irrationality. How exactly do people misapply Bayes law or deviate from SEU? For guidance on this, behavioral economists typically turn to the extensive experimental evidence compiled by cognitive psychologists on the biases that arise when people form *beliefs*, and on people's *preferences*, or on how they make decisions, given their beliefs. Psychology is therefore the second building block of behavioral finance, and we review the psychology most relevant for financial economists in Section 3.²

In Sections 4–8, we consider specific applications of behavioral finance: to understanding the aggregate stock market, the cross-section of average returns, and the pricing of closed-end funds in Sections 4, 5 and 6 respectively; to understanding how particular groups of investors choose their portfolios and trade over time in Section 7; and to understanding the financing and investment decisions of firms in Section 8. Section 9 takes stock and suggests directions for future research.³

2. Limits to arbitrage

2.1. Market efficiency

In the traditional framework where agents are rational and there are no frictions, a security's price equals its "fundamental value". This is the discounted sum of expected future cash flows, where in forming expectations, investors correctly process all available information, and where the discount rate is consistent with a normatively acceptable preference specification. The hypothesis that actual prices reflect fundamental values is the Efficient Markets Hypothesis (EMH). Put simply, under this hypothesis, "prices are right", in that they are set by agents who understand Bayes' law and have sensible preferences. In an efficient market, there is "no free lunch": no investment strategy can earn excess risk-adjusted average returns, or average returns greater than are warranted for its risk.

Behavioral finance argues that some features of asset prices are most plausibly interpreted as deviations from fundamental value, and that these deviations are brought about by the presence of traders who are not fully rational. A long-standing objection to this view that goes back to Friedman (1953) is that rational traders will quickly undo any dislocations caused by irrational traders. To illustrate the argument, suppose

² The idea, now widely adopted, that behavioral finance rests on the two pillars of limits to arbitrage and investor psychology is originally due to Shleifer and Summers (1990).

³ We draw readers' attention to two other recent surveys of behavioral finance. Shleifer (2000) provides a particularly detailed discussion of the theoretical and empirical work on limits to arbitrage, which we summarize in Section 2. Hirshleifer's (2001) survey is closer to ours in terms of material covered, although we devote less space to asset pricing, and more to corporate finance and individual investor behavior. We also organize the material somewhat differently.

that the fundamental value of a share of Ford is \$20. Imagine that a group of irrational traders becomes excessively pessimistic about Ford's future prospects and through its selling, pushes the price to \$15. Defenders of the EMH argue that rational traders, sensing an attractive opportunity, will buy the security at its bargain price and at the same time, hedge their bet by shorting a "substitute" security, such as General Motors, that has similar cash flows to Ford in future states of the world. The buying pressure on Ford shares will then bring their price back to fundamental value.

Friedman's line of argument is initially compelling, but it has not survived careful theoretical scrutiny. In essence, it is based on two assertions. First, as soon as there is a deviation from fundamental value – in short, a mispricing – an attractive investment opportunity is created. Second, rational traders will immediately snap up the opportunity, thereby correcting the mispricing. Behavioral finance does not take issue with the second step in this argument: when attractive investment opportunities come to light, it is hard to believe that they are not quickly exploited. Rather, it disputes the first step. The argument, which we elaborate on in Sections 2.2 and 2.3, is that even when an asset is wildly mispriced, strategies designed to correct the mispricing can be both risky and costly, rendering them unattractive. As a result, the mispricing can remain unchallenged.

It is interesting to think about common finance terminology in this light. While irrational traders are often known as "noise traders", rational traders are typically referred to as "arbitrageurs". Strictly speaking, an arbitrage is an investment strategy that offers riskless profits at no cost. Presumably, the rational traders in Friedman's fable became known as arbitrageurs because of the belief that a mispriced asset immediately creates an opportunity for riskless profits. Behavioral finance argues that this is *not* true: the strategies that Friedman would have his rational traders adopt are not necessarily arbitrages; quite often, they are very risky.

An immediate corollary of this line of thinking is that "prices are right" and "there is no free lunch" are *not* equivalent statements. While both are true in an efficient market, "no free lunch" can also be true in an inefficient market: just because prices are away from fundamental value does not necessarily mean that there are any excess risk-adjusted average returns for the taking. In other words,

$$\boxed{\text{"prices are right"} \Rightarrow \text{"no free lunch"}}$$

but

$$\boxed{\text{"no free lunch"} \not\Rightarrow \text{"prices are right"}.}$$

This distinction is important for evaluating the ongoing debate on market efficiency. First, many researchers still point to the inability of professional money managers to beat the market as strong evidence of market efficiency [Rubinstein (2001), Ross (2001)]. Underlying this argument, though, is the assumption that "no free lunch" implies "prices are right." If, as we argue in Sections 2.2 and 2.3, this link is broken, the

performance of money managers tells us little about whether prices reflect fundamental value.

Second, while some researchers accept that there is a distinction between “prices are right” and “there is no free lunch”, they believe that the debate should be more about the latter statement than about the former. We disagree with this emphasis. As economists, our ultimate concern is that capital be allocated to the most promising investment opportunities. Whether this is true or not depends much more on whether prices are right than on whether there are any free lunches for the taking.

2.2. Theory

In the previous section, we emphasized the idea that when a mispricing occurs, strategies designed to correct it can be both risky and costly, thereby allowing the mispricing to survive. Here we discuss some of the risks and costs that have been identified. In our discussion, we return to the example of Ford, whose fundamental value is \$20, but which has been pushed down to \$15 by pessimistic noise traders.

Fundamental risk. The most obvious risk an arbitrageur faces if he buys Ford’s stock at \$15 is that a piece of bad news about Ford’s fundamental value causes the stock to fall further, leading to losses. Of course, arbitrageurs are well aware of this risk, which is why they short a substitute security such as General Motors at the same time that they buy Ford. The problem is that substitute securities are rarely perfect, and often highly imperfect, making it impossible to remove all the fundamental risk. Shorting General Motors protects the arbitrageur somewhat from adverse news about the car industry as a whole, but still leaves him vulnerable to news that is specific to Ford – news about defective tires, say.⁴

Noise trader risk. Noise trader risk, an idea introduced by De Long et al. (1990a) and studied further by Shleifer and Vishny (1997), is the risk that the mispricing being exploited by the arbitrageur worsens in the short run. Even if General Motors is a *perfect* substitute security for Ford, the arbitrageur still faces the risk that the pessimistic investors causing Ford to be undervalued in the first place become even more pessimistic, lowering its price even further. Once one has granted the possibility that a security’s price can be different from its fundamental value, then one must also grant the possibility that future price movements will increase the divergence.

Noise trader risk matters because it can force arbitrageurs to liquidate their positions early, bringing them potentially steep losses. To see this, note that most real-world arbitrageurs – in other words, professional portfolio managers – are not managing their

⁴ Another problem is that even if a substitute security exists, it may itself be mispriced. This can happen in situations involving industry-wide mispricing: in that case, the only stocks with similar future cash flows to the mispriced one are themselves mispriced.

own money, but rather managing money for other people. In the words of Shleifer and Vishny (1997), there is “a separation of brains and capital”.

This agency feature has important consequences. Investors, lacking the specialized knowledge to evaluate the arbitrageur’s strategy, may simply evaluate him based on his returns. If a mispricing that the arbitrageur is trying to exploit worsens in the short run, generating negative returns, investors may decide that he is incompetent, and withdraw their funds. If this happens, the arbitrageur will be forced to liquidate his position prematurely. Fear of such premature liquidation makes him less aggressive in combating the mispricing in the first place.

These problems can be severely exacerbated by creditors. After poor short-term returns, creditors, seeing the value of their collateral erode, will call their loans, again triggering premature liquidation.

In these scenarios, the forced liquidation is brought about by the worsening of the mispricing itself. This need not always be the case. For example, in their efforts to remove fundamental risk, many arbitrageurs sell securities short. Should the original owner of the borrowed security want it back, the arbitrageur may again be forced to close out his position if he cannot find other shares to borrow. The risk that this occurs during a temporary worsening of the mispricing makes the arbitrageur more cautious from the start.

Implementation costs. Well-understood transaction costs such as commissions, bid-ask spreads and price impact can make it less attractive to exploit a mispricing. Since shorting is often essential to the arbitrage process, we also include short-sale constraints in the implementation costs category. These refer to anything that makes it less attractive to establish a short position than a long one. The simplest such constraint is the fee charged for borrowing a stock. In general these fees are small – D’Avolio (2002) finds that for most stocks, they range between 10 and 15 basis points – but they can be much larger; in some cases, arbitrageurs may not be able to find shares to borrow at *any* price. Other than the fees themselves, there can be legal constraints: for a large fraction of money managers – many pension fund and mutual fund managers in particular – short-selling is simply not allowed.⁵

We also include in this category the cost of finding and learning about a mispricing, as well as the cost of the resources needed to exploit it [Merton (1987)]. Finding

⁵ The presence of per-period transaction costs like lending fees can expose arbitrageurs to another kind of risk, *horizon risk*, which is the risk that the mispricing takes so long to close that any profits are swamped by the accumulated transaction costs. This applies even when the arbitrageur is certain that no outside party will force him to liquidate early. Abreu and Brunnermeier (2002) study a particular type of horizon risk, which they label *synchronization risk*. Suppose that the elimination of a mispricing requires the participation of a sufficiently large number of separate arbitrageurs. Then in the presence of per-period transaction costs, arbitrageurs may hesitate to exploit the mispricing because they don’t know how many *other* arbitrageurs have heard about the opportunity, and therefore how long they will have to wait before prices revert to correct values.

mispricing, in particular, can be a tricky matter. It was once thought that if noise traders influenced stock prices to any substantial degree, their actions would quickly show up in the form of predictability in returns. Shiller (1984) and Summers (1986) demonstrate that this argument is completely erroneous, with Shiller (1984) calling it "one of the most remarkable errors in the history of economic thought". They show that even if noise trader demand is so strong as to cause a large and persistent mispricing, it may generate so little predictability in returns as to be virtually undetectable.

In contrast, then, to straightforward-sounding textbook arbitrage, real world arbitrage entails both costs and risks, which under some conditions will limit arbitrage and allow deviations from fundamental value to persist. To see what these conditions are, consider two cases.

Suppose first that the mispriced security does *not* have a close substitute. By definition then, the arbitrageur is exposed to fundamental risk. In this case, sufficient conditions for arbitrage to be limited are (i) that arbitrageurs are risk averse and (ii) that the fundamental risk is systematic, in that it cannot be diversified by taking many such positions. Condition (i) ensures that the mispricing will not be wiped out by a single arbitrageur taking a large position in the mispriced security. Condition (ii) ensures that the mispricing will not be wiped out by a large number of investors each adding a *small* position in the mispriced security to their current holdings. The presence of noise trader risk or implementation costs will only limit arbitrage further.

Even if a perfect substitute does exist, arbitrage can still be limited. The existence of the substitute security immunizes the arbitrageur from fundamental risk. We can go further and assume that there are no implementation costs, so that only noise trader risk remains. De Long et al. (1990a) show that noise trader risk is powerful enough, that even with this single form of risk, arbitrage can sometimes be limited. The sufficient conditions are similar to those above, with one important difference. Here arbitrage will be limited if: (i) arbitrageurs are risk averse *and have short horizons* and (ii) the noise trader risk is systematic. As before, condition (i) ensures that the mispricing cannot be wiped out by a single, large arbitrageur, while condition (ii) prevents a large number of small investors from exploiting the mispricing. The central contribution of Shleifer and Vishny (1997) is to point out the real world relevance of condition (i): the possibility of an early, forced liquidation means that many arbitrageurs effectively have short horizons.

In the presence of certain implementation costs, condition (ii) may not even be necessary. If it is costly to learn about a mispricing, or the resources required to exploit it are expensive, that may be enough to explain why a large number of different individuals do not intervene in an attempt to correct the mispricing.

It is also important to note that for particular types of noise trading, arbitrageurs may prefer to trade in the *same* direction as the noise traders, thereby exacerbating the mispricing, rather than against them. For example, De Long et al. (1990b)

consider an economy with positive feedback traders, who buy more of an asset this period if it performed well last period. If these noise traders push an asset's price above fundamental value, arbitrageurs do not sell or short the asset. Rather, they *buy* it, knowing that the earlier price rise will attract more feedback traders next period, leading to still higher prices, at which point the arbitrageurs can exit at a profit.

So far, we have argued that it is not easy for arbitrageurs like hedge funds to exploit market inefficiencies. However, hedge funds are not the only market participants trying to take advantage of noise traders: firm managers also play this game. If a manager believes that investors are overvaluing his firm's shares, he can benefit the firm's existing shareholders by issuing extra shares at attractive prices. The extra supply this generates could potentially push prices back to fundamental value.

Unfortunately, this game entails risks and costs for managers, just as it does for hedge funds. Issuing shares is an expensive process, both in terms of underwriting fees and time spent by company management. Moreover, the manager can rarely be *sure* that investors are overvaluing his firm's shares. If he issues shares, thinking that they are overvalued when in fact they are not, he incurs the costs of deviating from his target capital structure, without getting any benefits in return.

2.3. Evidence

From the theoretical point of view, there is reason to believe that arbitrage is a risky process and therefore that it is only of limited effectiveness. But is there any *evidence* that arbitrage is limited? In principle, any example of persistent mispricing is immediate evidence of limited arbitrage: if arbitrage were not limited, the mispricing would quickly disappear. The problem is that while many pricing phenomena can be interpreted as deviations from fundamental value, it is only in a few cases that the presence of a mispricing can be established beyond any reasonable doubt. The reason for this is what Fama (1970) dubbed the "joint hypothesis problem". In order to claim that the price of a security differs from its properly discounted future cash flows, one needs a model of "proper" discounting. Any test of mispricing is therefore inevitably a *joint* test of mispricing and of a model of discount rates, making it difficult to provide definitive evidence of inefficiency.

In spite of this difficulty, researchers have uncovered a number of financial market phenomena that are almost certainly mispricings, and persistent ones at that. These examples show that arbitrage is indeed limited, and also serve as interesting illustrations of the risks and costs described earlier.

2.3.1. Twin shares

In 1907, Royal Dutch and Shell Transport, at the time completely independent companies, agreed to merge their interests on a 60:40 basis while remaining separate entities. Shares of Royal Dutch, which are primarily traded in the USA and in the

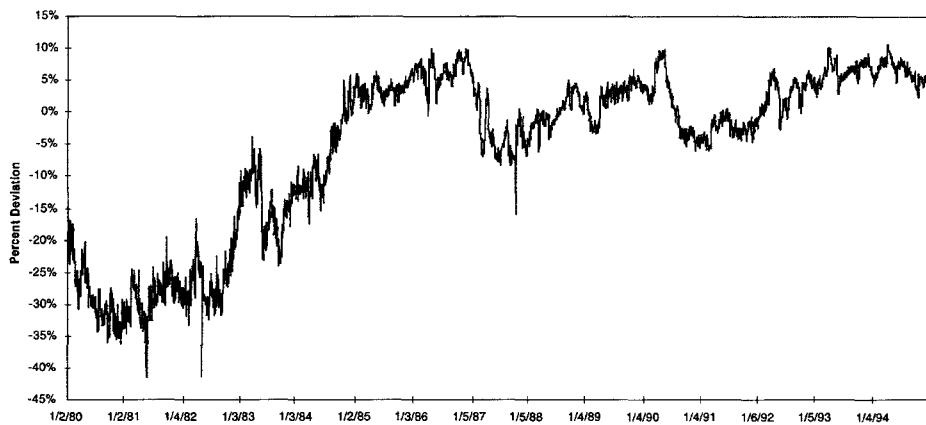


Fig. 1. Log deviations from Royal Dutch/Shell parity. Source: Froot and Dabora (1999).

Netherlands, are a claim to 60% of the total cash flow of the two companies, while Shell, which trades primarily in the UK, is a claim to the remaining 40%. If prices equal fundamental value, the market value of Royal Dutch equity should always be 1.5 times the market value of Shell equity. Remarkably, it isn't.

Figure 1, taken from Froot and Dabora's (1999) analysis of this case, shows the ratio of Royal Dutch equity value to Shell equity value relative to the efficient markets benchmark of 1.5. The picture provides strong evidence of a persistent inefficiency. Moreover, the deviations are not small. Royal Dutch is sometimes 35% underpriced relative to parity, and sometimes 15% overpriced.

This evidence of mispricing is simultaneously evidence of limited arbitrage, and it is not hard to see why arbitrage might be limited in this case. If an arbitrageur wanted to exploit this phenomenon – and several hedge funds, Long-Term Capital Management included, did try to – he would buy the relatively undervalued share and short the other. Table 1 summarizes the risks facing the arbitrageur. Since one share is a good substitute for the other, fundamental risk is nicely hedged: news about fundamentals should affect the two shares equally, leaving the arbitrageur immune. Nor are there

Table 1
Arbitrage costs and risks that arise in exploiting mispricing

Example	Fundamental risk (FR)	Noise trader risk (NTR)	Implementation costs (IC)
Royal Dutch/Shell	✗	✓	✗
Index Inclusions	✓	✓	✗
Palm/3Com	✗	✗	✓

any major implementation costs to speak of: shorting shares of either company is an easy matter.

The one risk that remains is noise trader risk. Whatever investor sentiment is causing one share to be undervalued relative to the other could also cause that share to become *even more* undervalued in the short term. The graph shows that this danger is very real: an arbitrageur buying a 10% undervalued Royal Dutch share in March 1983 would have seen it drop still further in value over the next six months. As discussed earlier, when a mispriced security has a perfect substitute, arbitrage can still be limited if (i) arbitrageurs are risk averse and have short horizons and (ii) the noise trader risk is systematic, or the arbitrage requires specialized skills, or there are costs to learning about such opportunities. It is very plausible that both (i) and (ii) are true, thereby explaining why the mispricing persisted for so long. It took until 2001 for the shares to finally sell at par.

This example also provides a nice illustration of the distinction between “prices are right” and “no free lunch” discussed in Section 2.1. While prices in this case are clearly *not* right, there are no easy profits for the taking.

2.3.2. Index inclusions

Every so often, one of the companies in the S&P 500 is taken out of the index because of a merger or bankruptcy, and is replaced by another firm. Two early studies of such index inclusions, Harris and Gurel (1986) and Shleifer (1986), document a remarkable fact: when a stock is added to the index, it jumps in price by an average of 3.5%, and much of this jump is permanent. In one dramatic illustration of this phenomenon, when Yahoo was added to the index, its shares jumped by 24% in a single day.

The fact that a stock jumps in value upon inclusion is once again clear evidence of mispricing: the price of the share changes even though its fundamental value does not. Standard and Poor’s emphasizes that in selecting stocks for inclusion, they are simply trying to make their index representative of the U.S. economy, not to convey any information about the level or riskiness of a firm’s future cash flows.⁶

This example of a deviation from fundamental value is also evidence of limited arbitrage. When one thinks about the risks involved in trying to exploit the anomaly, its persistence becomes less surprising. An arbitrageur needs to short the included security and to buy as good a substitute security as he can. This entails considerable

⁶ After the initial studies on index inclusions appeared, some researchers argued that the price increase might be rationally explained through information or liquidity effects. While such explanations cannot be completely ruled out, the case for mispricing was considerably strengthened by Kaul, Mehrotra and Morck (2000). They consider the case of the TS300 index of Canadian equities, which in 1996 changed the weights of some of its component stocks to meet an innocuous regulatory requirement. The reweighting was accompanied by significant price effects. Since the affected stocks were *already* in the index at the time of the event, information and liquidity explanations for the price jumps are extremely implausible.

fundamental risk because individual stocks rarely have good substitutes. It also carries substantial noise trader risk: whatever caused the initial jump in price – in all likelihood, buying by S&P 500 index funds – may continue, and cause the price to rise still further in the short run; indeed, Yahoo went from \$115 prior to its S&P inclusion announcement to \$210 a month later.

Wurgler and Zhuravskaya (2002) provide additional support for the limited arbitrage view of S&P 500 inclusions. They hypothesize that the jump upon inclusion should be particularly large for those stocks with the worst substitute securities, in other words, for those stocks for which the arbitrage is riskiest. By constructing the best possible substitute portfolio for each included stock, they are able to test this, and find strong support. Their analysis also shows just how hard it is to find good substitute securities for individual stocks. For most regressions of included stock returns on the returns of the best substitute securities, the R^2 is below 25%.

2.3.3. Internet carve-outs

In March 2000, 3Com sold 5% of its wholly owned subsidiary Palm Inc. in an initial public offering, retaining ownership of the remaining 95%. After the IPO, a shareholder of 3Com indirectly owned 1.5 shares of Palm. 3Com also announced its intention to spin off the remainder of Palm within 9 months, at which time they would give each 3Com shareholder 1.5 shares of Palm.

At the close of trading on the first day after the IPO, Palm shares stood at \$95, putting a lower bound on the value of 3Com at \$142. In fact, 3Com's price was \$81, implying a market valuation of 3Com's substantial businesses outside of Palm of about -\$60 per share!

This situation surely represents a severe mispricing, and it persisted for several weeks. To exploit it, an arbitrageur could buy one share of 3Com, short 1.5 shares of Palm, and wait for the spin-off, thus earning certain profits at no cost. This strategy entails no fundamental risk and no noise trader risk. Why, then, is arbitrage limited? Lamont and Thaler (2003), who analyze this case in detail, argue that implementation costs played a major role. Many investors who tried to borrow Palm shares to short were either told by their broker that no shares were available, or else were quoted a very high borrowing price. This barrier to shorting was not a legal one, but one that arose endogenously in the marketplace: such was the demand for shorting Palm, that the supply of Palm shorts was unable to meet it. Arbitrage was therefore limited, and the mispricing persisted.⁷

Some financial economists react to these examples by arguing that they are simply isolated instances with little broad relevance.⁸ We think this is an overly complacent

⁷ See also Mitchell, Pulvino and Stafford (2002) and Ofek and Richardson (2003) for further discussion of such “negative stub” situations, in which the market value of a company is less than the sum of its publicly traded parts.

⁸ During a discussion of these issues at a University of Chicago seminar, one economist argued that these examples are “the tip of the iceberg”, to which another retorted that “they *are* the iceberg”.

view. The “twin shares” example illustrates that in situations where arbitrageurs face only one type of risk – noise trader risk – securities can become mispriced by almost 35%. This suggests that if a typical stock trading on the NYSE or NASDAQ becomes subject to investor sentiment, the mispricing could be an order of magnitude larger. Not only would arbitrageurs face noise trader risk in trying to correct the mispricing, but fundamental risk as well, not to mention implementation costs.

3. Psychology

The theory of limited arbitrage shows that if irrational traders cause deviations from fundamental value, rational traders will often be powerless to do anything about it. In order to say more about the structure of these deviations, behavioral models often assume a specific form of irrationality. For guidance on this, economists turn to the extensive experimental evidence compiled by cognitive psychologists on the systematic biases that arise when people form *beliefs*, and on people’s *preferences*.⁹

In this section, we summarize the psychology that may be of particular interest to financial economists. Our discussion of each finding is necessarily brief. For a deeper understanding of the phenomena we touch on, we refer the reader to the surveys of Camerer (1995) and Rabin (1998) and to the edited volumes of Kahneman, Slovic and Tversky (1982), Kahneman and Tversky (2000) and Gilovich, Griffin and Kahneman (2002).

3.1. Beliefs

A crucial component of any model of financial markets is a specification of how agents form expectations. We now summarize what psychologists have learned about how people appear to form beliefs in practice.

Overconfidence. Extensive evidence shows that people are overconfident in their judgments. This appears in two guises. First, the confidence intervals people assign to their estimates of quantities – the level of the Dow in a year, say – are far too narrow. Their 98% confidence intervals, for example, include the true quantity only about 60% of the time [Alpert and Raiffa (1982)]. Second, people are poorly calibrated when estimating probabilities: events they think are certain to occur actually occur only

⁹ We emphasize, however, that behavioral models do not *need* to make extensive psychological assumptions in order to generate testable predictions. In Section 6, we discuss Lee, Shleifer and Thaler’s (1991) theory of closed-end fund pricing. That theory makes numerous crisp predictions using only the assumptions that there are noise traders with correlated sentiment in the economy, and that arbitrage is limited.

around 80% of the time, and events they deem impossible occur approximately 20% of the time [Fischhoff, Slovic and Lichtenstein (1977)].¹⁰

Optimism and wishful thinking. Most people display unrealistically rosy views of their abilities and prospects [Weinstein (1980)]. Typically, over 90% of those surveyed think they are above average in such domains as driving skill, ability to get along with people and sense of humor. They also display a systematic planning fallacy: they predict that tasks (such as writing survey papers) will be completed much sooner than they actually are [Buehler, Griffin and Ross (1994)].

Representativeness. Kahneman and Tversky (1974) show that when people try to determine the probability that a data set A was generated by a model B, or that an object A belongs to a class B, they often use the representativeness heuristic. This means that they evaluate the probability by the degree to which A reflects the essential characteristics of B.

Much of the time, representativeness is a helpful heuristic, but it can generate some severe biases. The first is *base rate neglect*. To illustrate, Kahneman and Tversky present this description of a person named Linda:

Linda is 31 years old, single, outspoken, and very bright. She majored in philosophy. As a student, she was deeply concerned with issues of discrimination and social justice, and also participated in anti-nuclear demonstrations.

When asked which of “Linda is a bank teller” (statement A) and “Linda is a bank teller and is active in the feminist movement” (statement B) is more likely, subjects typically assign greater probability to B. This is, of course, impossible. Representativeness provides a simple explanation. The description of Linda *sounds* like the description of a feminist – it is representative of a feminist – leading subjects to pick B. Put differently, while Bayes law says that

$$p(\text{statement B} \mid \text{description}) = \frac{p(\text{description} \mid \text{statement B}) p(\text{statement B})}{p(\text{description})},$$

people apply the law incorrectly, putting too much weight on $p(\text{description} \mid \text{statement B})$, which captures representativeness, and too little weight on the base rate, $p(\text{statement B})$.

¹⁰ Overconfidence may in part stem from two other biases, self-attribution bias and hindsight bias. Self-attribution bias refers to people’s tendency to ascribe any success they have in some activity to their own talents, while blaming failure on bad luck, rather than on their ineptitude. Doing this repeatedly will lead people to the pleasing but erroneous conclusion that they are very talented. For example, investors might become overconfident after several quarters of investing success [Gervais and Odean (2001)]. Hindsight bias is the tendency of people to believe, after an event has occurred, that they predicted it before it happened. If people think they predicted the past better than they actually did, they may also believe that they can predict the future better than they actually can.

Representativeness also leads to another bias, *sample size neglect*. When judging the likelihood that a data set was generated by a particular model, people often fail to take the size of the sample into account: after all, a small sample can be just as representative as a large one. Six tosses of a coin resulting in three heads and three tails are as representative of a fair coin as 500 heads and 500 tails are in a total of 1000 tosses. Representativeness implies that people will find the two sets of tosses equally informative about the fairness of the coin, even though the second set is much more so.

Sample size neglect means that in cases where people do not initially know the data-generating process, they will tend to infer it too quickly on the basis of too few data points. For instance, they will come to believe that a financial analyst with four good stock picks is talented because four successes are not representative of a bad or mediocre analyst. It also generates a “hot hand” phenomenon, whereby sports fans become convinced that a basketball player who has made three shots in a row is on a hot streak and will score again, even though there is no evidence of a hot hand in the data [Gilovich, Vallone and Tversky (1985)]. This belief that even small samples will reflect the properties of the parent population is sometimes known as the “law of small numbers” [Rabin (2002)].

In situations where people *do* know the data-generating process in advance, the law of small numbers leads to a gambler’s fallacy effect. If a fair coin generates five heads in a row, people will say that “tails are due”. Since they believe that even a short sample should be representative of the fair coin, there have to be more tails to balance out the large number of heads.

Conservatism. While representativeness leads to an underweighting of base rates, there are situations where base rates are *over-emphasized* relative to sample evidence. In an experiment run by Edwards (1968), there are two urns, one containing 3 blue balls and 7 red ones, and the other containing 7 blue balls and 3 red ones. A random draw of 12 balls, with replacement, from one of the urns yields 8 reds and 4 blues. What is the probability the draw was made from the first urn? While the correct answer is 0.97, most people estimate a number around 0.7, apparently overweighting the base rate of 0.5.

At first sight, the evidence of conservatism appears at odds with representativeness. However, there may be a natural way in which they fit together. It appears that if a data sample is representative of an underlying model, then people overweight the data. However, if the data is not representative of any salient model, people react too little to the data and rely too much on their priors. In Edwards’ experiment, the draw of 8 red and 4 blue balls is not particularly representative of either urn, possibly leading to an overreliance on prior information.¹¹

¹¹ Mullainathan (2001) presents a formal model that neatly reconciles the evidence on underweighting sample information with the evidence on overweighting sample information.

Belief perseverance. There is much evidence that once people have formed an opinion, they cling to it too tightly and for too long [Lord, Ross and Lepper (1979)]. At least two effects appear to be at work. First, people are reluctant to search for evidence that contradicts their beliefs. Second, even if they find such evidence, they treat it with excessive skepticism. Some studies have found an even stronger effect, known as confirmation bias, whereby people misinterpret evidence that goes against their hypothesis as actually being in their favor. In the context of academic finance, belief perseverance predicts that if people start out believing in the Efficient Markets Hypothesis, they may continue to believe in it long after compelling evidence to the contrary has emerged.

Anchoring. Kahneman and Tversky (1974) argue that when forming estimates, people often start with some initial, possibly arbitrary value, and then adjust away from it. Experimental evidence shows that the adjustment is often insufficient. Put differently, people “anchor” too much on the initial value.

In one experiment, subjects were asked to estimate the percentage of United Nations’ countries that are African. More specifically, before giving a percentage, they were asked whether their guess was higher or lower than a randomly generated number between 0 and 100. Their subsequent estimates were significantly affected by the initial random number. Those who were asked to compare their estimate to 10, subsequently estimated 25%, while those who compared to 60, estimated 45%.

Availability biases. When judging the probability of an event – the likelihood of getting mugged in Chicago, say – people often search their memories for relevant information. While this is a perfectly sensible procedure, it can produce biased estimates because not all memories are equally retrievable or “available”, in the language of Kahneman and Tversky (1974). More recent events and more salient events – the mugging of a close friend, say – will weigh more heavily and distort the estimate.

Economists are sometimes wary of this body of experimental evidence because they believe (i) that people, through repetition, will learn their way out of biases; (ii) that experts in a field, such as traders in an investment bank, will make fewer errors; and (iii) that with more powerful incentives, the effects will disappear.

While all these factors can attenuate biases to some extent, there is little evidence that they wipe them out altogether. The effect of learning is often muted by errors of application: when the bias is explained, people often understand it, but then immediately proceed to violate it again in specific applications. Expertise, too, is often a hindrance rather than a help: experts, armed with their sophisticated models, have been found to exhibit *more* overconfidence than laymen, particularly when they receive only limited feedback about their predictions. Finally, in a review of dozens of studies on the topic, Camerer and Hogarth (1999, p. 7) conclude that while incentives can

sometimes reduce the biases people display, “no replicated study has made rationality violations disappear purely by raising incentives”.

3.2. Preferences

3.2.1. Prospect theory

An essential ingredient of any model trying to understand asset prices or trading behavior is an assumption about investor preferences, or about how investors evaluate risky gambles. The vast majority of models assume that investors evaluate gambles according to the expected utility framework, EU henceforth. The theoretical motivation for this goes back to Von Neumann and Morgenstern (1944), VNM henceforth, who show that if preferences satisfy a number of plausible axioms – completeness, transitivity, continuity, and independence – then they can be represented by the expectation of a utility function.

Unfortunately, experimental work in the decades after VNM has shown that people systematically violate EU theory when choosing among risky gambles. In response to this, there has been an explosion of work on so-called non-EU theories, all of them trying to do a better job of matching the experimental evidence. Some of the better known models include weighted-utility theory [Chew and MacCrimmon (1979), Chew (1983)], implicit EU [Chew (1989), Dekel (1986)], disappointment aversion [Gul (1991)], regret theory [Bell (1982), Loomes and Sugden (1982)], rank-dependent utility theories [Quiggin (1982), Segal (1987, 1989), Yaari (1987)], and prospect theory [Kahneman and Tversky (1979), Tversky and Kahneman (1992)].

Should financial economists be interested in any of these alternatives to expected utility? It may be that EU theory is a good approximation to how people evaluate a risky gamble like the stock market, even if it does not explain attitudes to the kinds of gambles studied in experimental settings. On the other hand, the difficulty the EU approach has encountered in trying to explain basic facts about the stock market suggests that it may be worth taking a closer look at the experimental evidence. Indeed, recent work in behavioral finance has argued that some of the lessons we learn from violations of EU are central to understanding a number of financial phenomena.

Of all the non-EU theories, prospect theory may be the most promising for financial applications, and we discuss it in detail. The reason we focus on this theory is, quite simply, that it is the most successful at capturing the experimental results. In a way, this is not surprising. Most of the other non-EU models are what might be called quasi-normative, in that they try to capture some of the anomalous experimental evidence by slightly weakening the VNM axioms. The difficulty with such models is that in trying to achieve two goals – normative and descriptive – they end up doing an unsatisfactory job at both. In contrast, prospect theory has no aspirations as a normative theory: it simply tries to capture people’s attitudes to risky gambles as parsimoniously as possible. Indeed, Tversky and Kahneman (1986) argue convincingly that normative approaches are doomed to failure, because people routinely make choices that are

simply impossible to justify on normative grounds, in that they violate dominance or invariance.

Kahneman and Tversky (1979), KT henceforth, lay out the original version of prospect theory, designed for gambles with at most two non-zero outcomes. They propose that when offered a gamble

$$(x, p; y, q),$$

to be read as “get outcome x with probability p , outcome y with probability q ”, where $x \leq 0 < y$ or $y \leq 0 < x$, people assign it a value of

$$\pi(p)v(x) + \pi(q)v(y), \quad (1)$$

where v and π are shown in Figure 2. When choosing between different gambles, they pick the one with the highest value.

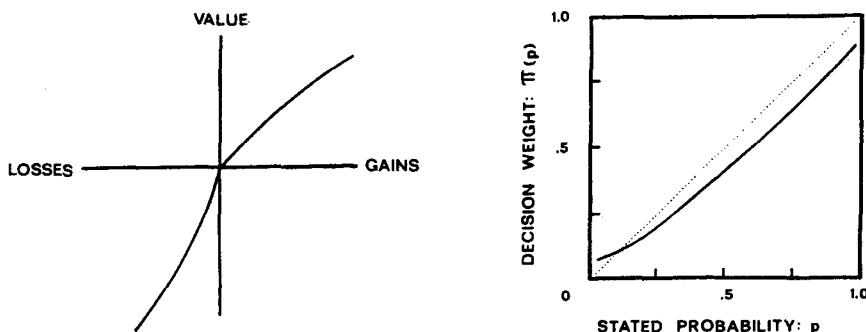


Fig. 2. Kahneman and Tversky's (1979) proposed value function v and probability weighting function π .

This formulation has a number of important features. First, utility is defined over gains and losses rather than over final wealth positions, an idea first proposed by Markowitz (1952). This fits naturally with the way gambles are often presented and discussed in everyday life. More generally, it is consistent with the way people perceive attributes such as brightness, loudness, or temperature relative to earlier levels, rather than in absolute terms. Kahneman and Tversky (1979) also offer the following violation of EU as evidence that people focus on gains and losses. Subjects are asked:¹²

¹² All the experiments in Kahneman and Tversky (1979) are conducted in terms of Israeli currency. The authors note that at the time of their research, the median monthly family income was about 3000 Israeli lira.

In addition to whatever you own, you have been given 1000. Now choose between

$$A = (1000, 0.5)$$

$$B = (500, 1).$$

B was the more popular choice. The same subjects were then asked:

In addition to whatever you own, you have been given 2000. Now choose between

$$C = (-1000, 0.5)$$

$$D = (-500, 1).$$

This time, *C* was more popular.

Note that the two problems are identical in terms of their final wealth positions and yet people choose differently. The subjects are apparently focusing only on gains and losses. Indeed, when they are not given any information about prior winnings, they choose *B* over *A* and *C* over *D*.

The second important feature is the shape of the value function v , namely its concavity in the domain of gains and convexity in the domain of losses. Put simply, people are risk averse over gains, and risk-seeking over losses. Simple evidence for this comes from the fact just mentioned, namely that in the absence of any information about prior winnings¹³

$$B \succ A, \quad C \succ D.$$

The v function also has a kink at the origin, indicating a greater sensitivity to losses than to gains, a feature known as *loss aversion*. Loss aversion is introduced to capture aversion to bets of the form:

$$E = \left(110, \frac{1}{2}; -100, \frac{1}{2}\right).$$

It may seem surprising that we need to depart from the expected utility framework in order to understand attitudes to gambles as simple as *E*, but it is nonetheless true. In a remarkable paper, Rabin (2000) shows that if an expected utility maximizer rejects gamble *E* at all wealth levels, then he will also reject

$$\left(20000000, \frac{1}{2}; -1000, \frac{1}{2}\right),$$

an utterly implausible prediction. The intuition is simple: if a smooth, increasing, and concave utility function defined over final wealth has sufficient local curvature to reject

¹³ In this section $G_1 \succ G_2$ should be read as “a statistically significant fraction of Kahneman and Tversky’s subjects preferred G_1 to G_2 .”

E over a wide range of wealth levels, it must be an extraordinarily concave function, making the investor extremely risk averse over large stakes gambles.

The final piece of prospect theory is the nonlinear probability transformation. Small probabilities are overweighted, so that $\pi(p) > p$. This is deduced from KT's finding that

$$(5000, 0.001) \succ (5, 1),$$

and

$$(-5, 1) \succ (-5000, 0.001),$$

together with the earlier assumption that v is concave (convex) in the domain of gains (losses). Moreover, people are more sensitive to differences in probabilities at higher probability levels. For example, the following pair of choices,

$$(3000, 1) \succ (4000, 0.8; 0, 0.2),$$

and

$$(4000, 0.2; 0, 0.8) \succ (3000, 0.25),$$

which violate EU theory, imply

$$\frac{\pi(0.25)}{\pi(0.2)} < \frac{\pi(1)}{\pi(0.8)}.$$

The intuition is that the 20% jump in probability from 0.8 to 1 is more striking to people than the 20% jump from 0.2 to 0.25. In particular, people place much more weight on outcomes that are certain relative to outcomes that are merely probable, a feature sometimes known as the "certainty effect".

Along with capturing experimental evidence, prospect theory also simultaneously explains preferences for insurance and for buying lottery tickets. Although the concavity of v in the region of gains generally produces risk aversion, for lotteries which offer a small chance of a large gain, the overweighting of small probabilities in Figure 2 dominates, leading to risk-seeking. Along the same lines, while the convexity of v in the region of losses typically leads to risk-seeking, the same overweighting of small probabilities induces risk aversion over gambles which have a small chance of a large loss.

Based on additional evidence, Tversky and Kahneman (1992) propose a generalization of prospect theory which can be applied to gambles with more than two

outcomes. Specifically, if a gamble promises outcome x_i with probability p_i , Tversky and Kahneman (1992) propose that people assign the gamble the value

$$\sum_i \pi_i v(x_i), \quad (2)$$

where

$$v = \begin{cases} x^\alpha & \text{if } x \geq 0 \\ -\lambda(-x)^\alpha & \text{if } x < 0 \end{cases}$$

and

$$\pi_i = w(P_i) - w(P_i^*),$$

$$w(P) = \frac{P^\gamma}{(P^\gamma + (1-P)^\gamma)^{1/\gamma}}.$$

Here, P_i (P_i^*) is the probability that the gamble will yield an outcome at least as good as (strictly better than) x_i . Tversky and Kahneman (1992) use experimental evidence to estimate $\alpha = 0.88$, $\lambda = 2.25$, and $\gamma = 0.65$. Note that λ is the coefficient of loss aversion, a measure of the relative sensitivity to gains and losses. Over a wide range of experimental contexts λ has been estimated in the neighborhood of 2.

Earlier in this section, we saw how prospect theory could explain why people made different choices in situations with identical final wealth levels. This illustrates an important feature of the theory, namely that it can accommodate the effects of problem description, or of *framing*. Such effects are powerful. There are numerous demonstrations of a 30 to 40% shift in preferences depending on the wording of a problem. No normative theory of choice can accommodate such behavior since a first principle of rational choice is that choices should be independent of the problem description or representation.

Framing refers to the way a problem is posed for the decision maker. In many actual choice contexts the decision maker also has flexibility in how to think about the problem. For example, suppose that a gambler goes to the race track and wins \$200 in his first bet, but then loses \$50 on his second bet. Does he code the outcome of the second bet as a loss of \$50 or as a reduction in his recently won gain of \$200? In other words, is the utility of the second loss $v(-50)$ or $v(150) - v(200)$? The process by which people formulate such problems for themselves is called *mental accounting* [Thaler (2000)]. Mental accounting matters because in prospect theory, v is nonlinear.

One important feature of mental accounting is *narrow framing*, which is the tendency to treat individual gambles separately from other portions of wealth. In other words, when offered a gamble, people often evaluate it as if it is the only gamble they face in the world, rather than merging it with pre-existing bets to see if the new bet is a worthwhile addition.

Redelmeier and Tversky (1992) provide a simple illustration, based on the gamble

$$F = \left(2000, \frac{1}{2}; -500, \frac{1}{2}\right).$$

Subjects in their experiment were asked whether they were willing to take this bet; 57% said they would not. They were then asked whether they would prefer to play F five times or six times; 70% preferred the six-fold gamble. Finally they were asked:

Suppose that you have played F five times but you don't yet know your wins and losses. Would you play the gamble a sixth time?

60% rejected the opportunity to play a sixth time, reversing their preference from the earlier question. This suggests that some subjects are framing the sixth gamble narrowly, segregating it from the other gambles. Indeed, the 60% rejection level is very similar to the 57% rejection level for the one-off play of F .

3.2.2. Ambiguity aversion

Our discussion so far has centered on understanding how people act when the outcomes of gambles have known objective probabilities. In reality, probabilities are rarely objectively known. To handle these situations, Savage (1964) develops a counterpart to expected utility known as subjective expected utility, SEU henceforth. Under certain axioms, preferences can be represented by the expectation of a utility function, this time weighted by the individual's subjective probability assessment.

Experimental work in the last few decades has been as unkind to SEU as it was to EU. The violations this time are of a different nature, but they may be just as relevant for financial economists.

The classic experiment was described by Ellsberg (1961). Suppose that there are two urns, 1 and 2. Urn 2 contains a total of 100 balls, 50 red and 50 blue. Urn 1 also contains 100 balls, again a mix of red and blue, but the subject does not know the proportion of each.

Subjects are asked to choose one of the following two gambles, each of which involves a possible payment of \$100, depending on the color of a ball drawn at random from the relevant urn

- a_1 : a ball is drawn from Urn 1, \$100 if red, \$0 if blue,
- a_2 : a ball is drawn from Urn 2, \$100 if red, \$0 if blue.

Subjects are then also asked to choose between the following two gambles:

- b_1 : a ball is drawn from Urn 1, \$100 if blue, \$0 if red,
- b_2 : a ball is drawn from Urn 2, \$100 if blue, \$0 if red.

a_2 is typically preferred to a_1 , while b_2 is chosen over b_1 . These choices are inconsistent with SEU: the choice of a_2 implies a subjective probability that *fewer* than 50% of the balls in Urn 1 are red, while the choice of b_2 implies the opposite.

The experiment suggests that people do not like situations where they are uncertain about the probability distribution of a gamble. Such situations are known as situations of ambiguity, and the general dislike for them, as ambiguity aversion.¹⁴ SEU does not allow agents to express their degree of confidence about a probability distribution and therefore cannot capture such aversion.

Ambiguity aversion appears in a wide variety of contexts. For example, a researcher might ask a subject for his estimate of the probability that a certain team will win its upcoming football match, to which the subject might respond 0.4. The researcher then asks the subject to imagine a chance machine, which will display 1 with probability 0.4 and 0 otherwise, and asks whether the subject would prefer to bet on the football game – an ambiguous bet – or on the machine, which offers no ambiguity. In general, people prefer to bet on the machine, illustrating aversion to ambiguity.

Heath and Tversky (1991) argue that in the real world, ambiguity aversion has much to do with how competent an individual feels he is at assessing the relevant distribution. Ambiguity aversion over a bet can be strengthened by highlighting subjects' feelings of incompetence, either by showing them other bets in which they have more expertise, or by mentioning other people who are more qualified to evaluate the bet [Fox and Tversky (1995)].

Further evidence that supports the competence hypothesis is that in situations where people feel especially competent in evaluating a gamble, the opposite of ambiguity aversion, namely a “preference for the familiar”, has been observed. In the example above, people chosen to be especially knowledgeable about football often prefer to bet on the outcome of the game than on the chance machine. Just as with ambiguity aversion, such behavior cannot be captured by SEU.

4. Application: The aggregate stock market

Researchers studying the aggregate U.S. stock market have identified a number of interesting facts about its behavior. Three of the most striking are:

The Equity Premium. The stock market has historically earned a high excess rate of return. For example, using annual data from 1871–1993, Campbell and Cochrane (1999) report that the average log return on the S&P 500 index is 3.9% higher than the average log return on short-term commercial paper.

Volatility. Stock returns and price–dividend ratios are both highly variable. In the same data set, the annual standard deviation of excess log returns on the S&P 500 is 18%, while the annual standard deviation of the log price–dividend ratio is 0.27.

¹⁴ An early discussion of this aversion can be found in Knight (1921), who defines risk as a gamble with known distribution and uncertainty as a gamble with unknown distribution, and suggests that people dislike uncertainty more than risk.

Predictability. Stock returns are forecastable. Using monthly, real, equal-weighted NYSE returns from 1941–1986, Fama and French (1988) show that the dividend–price ratio is able to explain 27% of the variation of cumulative stock returns over the subsequent four years.¹⁵

All three of these facts can be labelled puzzles. The first fact has been known as the equity premium puzzle since the work of Mehra and Prescott (1985) [see also Hansen and Singleton (1983)]. Campbell (1999) calls the second fact the volatility puzzle and we refer to the third fact as the predictability puzzle. The reason they are called puzzles is that they are hard to rationalize in a simple consumption-based model.

To see this, consider the following endowment economy, which we come back to a number of times in this section. There are an infinite number of identical investors, and two assets: a risk-free asset in zero net supply, with gross return $R_{f,t}$ between time t and $t+1$, and a risky asset – the stock market – in fixed positive supply, with gross return R_{t+1} between time t and $t+1$. The stock market is a claim to a perishable stream of dividends $\{D_t\}$, where

$$\frac{D_{t+1}}{D_t} = \exp [g_D + \sigma_D \varepsilon_{t+1}], \quad (3)$$

and where each period's dividend can be thought of as one component of a consumption endowment C_t , where

$$\frac{C_{t+1}}{C_t} = \exp [g_C + \sigma_C \eta_{t+1}], \quad (4)$$

and

$$\begin{pmatrix} \varepsilon_t \\ \eta_t \end{pmatrix} \sim N \left(\begin{pmatrix} 0 \\ 0 \end{pmatrix}, \begin{pmatrix} 1 & \omega \\ \omega & 1 \end{pmatrix} \right), \quad \text{i.i.d. over time.} \quad (5)$$

Investors choose consumption C_t and an allocation S_t to the risky asset to maximize

$$E_0 \sum_{t=0}^{\infty} \rho^t \frac{C_t^{1-\gamma}}{1-\gamma}, \quad (6)$$

subject to the standard budget constraint.¹⁶ Using the Euler equation of optimality,

$$1 = \rho E_t \left[\left(\frac{C_{t+1}}{C_t} \right)^{-\gamma} R_{t+1} \right], \quad (7)$$

it is straightforward to derive expressions for stock returns and prices. The details are in the Appendix.

¹⁵ These three facts are widely agreed on, but they are not completely uncontroversial. A large literature has debated the statistical significance of the time series predictability, while others have argued that the equity premium is overstated due to survivorship bias [Brown, Goetzmann and Ross (1995)].

¹⁶ For $\gamma = 1$, we replace $C_t^{1-\gamma}/1-\gamma$ with $\log(C_t)$.

Table 2
Parameter values for a simple consumption-based model

Parameter	g_C	σ_C	g_D	σ_D	ω	γ	ρ
Value	1.84%	3.79%	1.5%	12.0%	0.15	1.0	0.98

We can now examine the model's quantitative predictions for the parameter values in Table 2. The endowment process parameters are taken from U.S. data spanning the 20th century, and are standard in the literature. It is also standard to start out by considering *low* values of γ . The reason is that when one computes, for various values of γ , how much wealth an individual would be prepared to give up to avoid a large-scale timeless wealth gamble, low values of γ match best with introspection as to what the answers should be [Mankiw and Zeldes (1991)]. We take $\gamma = 1$, which corresponds to log utility.

In an economy with these parameter values, the average log return on the stock market would be just 0.1% higher than the risk-free rate, not the 3.9% observed historically. The standard deviation of log stock returns would be only 12%, not 18%, and the price–dividend ratio would be constant (implying, of course, that the dividend–price ratio has no forecast power for future returns).

It is useful to recall the intuition for these results. In an economy with power utility preferences, the equity premium is determined by risk aversion γ and by risk, measured as the covariance of stock returns and consumption growth. Since consumption growth is very smooth in the data, this covariance is very low, thus predicting a very low equity premium. Stocks simply do not appear risky to investors with the preferences in Equation (6) and with low γ , and therefore do not warrant a large premium. Of course, the equity premium predicted by the model can be increased by using higher values of γ . However, other than making counterintuitive predictions about individuals' attitudes to large-scale gambles, this would also predict a counterfactually high risk-free rate, a problem known as the risk-free rate puzzle [Weil (1989)].

To understand the volatility puzzle, note that in the simple economy described above, both discount rates and expected dividend growth are constant over time. A direct application of the present value formula implies that the price–dividend ratio, P/D henceforth, is constant. Since

$$R_{t+1} = \frac{D_{t+1} + P_{t+1}}{P_t} = \frac{1 + P_{t+1}/D_{t+1}}{P_t/D_t} \frac{D_{t+1}}{D_t}, \quad (8)$$

it follows that

$$r_{t+1} = \Delta d_{t+1} + \text{const.} \equiv d_{t+1} - d_t + \text{const.}, \quad (9)$$

where lower case letters indicate log variables. The standard deviation of log returns will therefore only be as high as the standard deviation of log dividend growth, namely 12%.

The particular volatility puzzle seen here illustrates a more general point, first made by Shiller (1981) and LeRoy and Porter (1981), namely that it is difficult to explain the historical volatility of stock returns with *any* model in which investors are rational and discount rates are constant.

To see the intuition, consider the identity in Equation (8) again. Since the volatility of log dividend growth is only 12%, the only way for a model to generate an 18% volatility of log returns is to introduce variation in the P/D ratio. But if discount rates are constant, a quick glance at a present-value formula shows that the only way to do that is to introduce variation in investors' forecasts of the dividend growth rate: a higher forecast raises the P/D ratio, a lower forecast brings it down. There is a catch here, though: if investors are rational, their expectations for dividend growth must, on average, be confirmed. In other words, times of higher (lower) P/D ratios should, on average, be followed by higher (lower) cash-flow growth. Unfortunately, price–dividend ratios are *not* reliable forecasters of dividend growth, neither in the USA nor in most international markets (see Campbell (1999), for recent evidence).

Shiller and LeRoy and Porter's results shocked the profession when they first appeared. At the time, most economists felt that discount rates *were* close to constant over time, apparently implying that stock market volatility could only be fully explained by appealing to investor irrationality. Today, it is well understood that rational variation in discount rates can help explain the volatility puzzle, although we argue later that models with irrational beliefs also offer a plausible way of thinking about the data.

Both the rational and behavioral approaches to finance have made progress in understanding the three puzzles singled out at the start of this section. The advances on the rational side are well described in other articles in this handbook. Here, we discuss the behavioral approaches, starting with the equity premium puzzle and then turning to the volatility puzzle.

We do not consider the predictability puzzle separately, because in any model with a stationary P/D ratio, a resolution of the volatility puzzle is simultaneously a resolution of the predictability puzzle. To see this, recall from Equation (8) that any model which captures the empirical volatility of returns must involve variation in the P/D ratio. Moreover, for a model to be a *satisfactory* resolution of the volatility puzzle, it should not make the counterfactual prediction that P/D ratios forecast subsequent dividend growth. Now suppose that the P/D ratio is higher than average. The only way it can return to its mean is if cash flows D subsequently go up, or if prices P fall. Since the P/D ratio is not allowed to forecast cash flows, it must forecast lower returns, thereby explaining the predictability puzzle.

4.1. The equity premium puzzle

The core of the equity premium puzzle is that even though stocks appear to be an attractive asset – they have high average returns and a low covariance with consumption

growth – investors appear very unwilling to hold them. In particular, they appear to demand a substantial risk premium in order to hold the market supply.

To date, behavioral finance has pursued two approaches to this puzzle. Both are based on preferences: one relies on prospect theory, the other on ambiguity aversion. In essence, both approaches try to understand what it is that is missing from the popular preference specification in Equation (6) that makes investors fear stocks so much, leading them to charge a high premium in equilibrium.

4.1.1. Prospect theory

One of the earliest papers to link prospect theory to the equity premium is Benartzi and Thaler (1995), BT henceforth. They study how an investor with prospect theory-type preferences allocates his financial wealth between T-Bills and the stock market. Prospect theory argues that when choosing between gambles, people compute the gains and losses for each one and select the one with the highest prospective utility. In a financial context, this suggests that people may choose a portfolio allocation by computing, for each allocation, the potential gains and losses in the value of their holdings, and then taking the allocation with the highest prospective utility. In other words, they choose ω , the fraction of financial wealth in stocks, to maximize

$$E_\pi v [(1 - \omega) R_{f,t+1} + \omega R_{t+1} - 1], \quad (10)$$

where π and v are defined in Equation (2). In particular, v captures loss aversion, the experimental finding that people are more sensitive to losses than to gains. $R_{f,t+1}$ and R_{t+1} are the gross returns on T-Bills and the stock market between t and $t + 1$, respectively, making the argument of v the return on financial wealth.

In order to implement this model, BT need to stipulate how often investors evaluate their portfolios. In other words, how long is the time interval between t and $t + 1$? To see why this matters, compare two investors: energetic Nick who calculates the gains and losses in his portfolio every day, and laid-back Dick who looks at his portfolio only once per decade. Since, on a daily basis, stocks go down in value almost as often as they go up, the loss aversion built into v makes stocks appear unattractive to Nick. In contrast, loss aversion does not have much effect on Dick's perception of stocks since, at ten year horizons, stocks offer only a small risk of losing money.

Rather than simply pick an evaluation interval, BT calculate how often investors would have to evaluate their portfolios to make them indifferent between stocks and T-Bills: in other words, given historical U.S. data on stocks and T-Bills, for what evaluation interval would substituting $\omega = 0$ and $\omega = 1$ into Equation (10) give the same prospective utility? Roughly speaking, this calculation can be thought of as asking what kind of equity premium might be sustainable in equilibrium: how often would investors need to evaluate their gains and losses so that even in the face of the large historical equity premium, they would still be happy to hold the market supply of T-Bills.

BT find that for the parametric forms for π and v estimated in experimental settings, the answer is one year, and they argue that this is indeed a natural evaluation period for investors to use. The way people frame gains and losses is plausibly influenced by the way information is presented to them. Since we receive our most comprehensive mutual fund reports once a year, and do our taxes once a year, it is not unreasonable that gains and losses might be expressed as annual changes in value.

The BT calculation therefore suggests a simple way of understanding the high historical equity premium. If investors get utility from annual changes in financial wealth and are loss averse over these changes, their fear of a major drop in financial wealth will lead them to demand a high premium as compensation. BT call the combination of loss aversion and frequent evaluations *myopic loss aversion*.

BT's result is only *suggestive* of a solution to Mehra and Prescott's equity premium puzzle. As emphasized at the start of this section, that puzzle is in large part a consumption puzzle: given the low volatility of consumption growth, why are investors so reluctant to buy a high return asset, stocks, especially when that asset's covariance with consumption growth is so low? Since BT do not consider an intertemporal model with consumption choice, they cannot address this issue directly.

To see if prospect theory can in fact help with the equity premium puzzle, Barberis, Huang and Santos (2001), BHS henceforth, make a first attempt at building it into a dynamic equilibrium model of stock returns. A simple version of their model, an extension of which we consider later, examines an economy with the same structure as the one described at the start of Section 4, but in which investors have the preferences

$$E_0 \sum_{t=0}^{\infty} \left[\rho' \frac{C_t^{1-\gamma}}{1-\gamma} + b_0 \bar{C}_t^{-\gamma} \hat{v}(X_{t+1}) \right]. \quad (11)$$

The investor gets utility from consumption, but over and above that, he gets utility from changes in the value of his holdings of the risky asset between t and $t+1$, denoted here by X_{t+1} . Motivated by BT's findings, BHS define the unit of time to be a year, so that gains and losses are measured annually.

The utility from these gains and losses is determined by \hat{v} where

$$\hat{v}(X) = \begin{cases} X & \text{for } X \geq 0, \\ 2.25X & \text{for } X < 0. \end{cases} \quad (12)$$

The 2.25 factor comes from Tversky and Kahneman's (1992) experimental study of attitudes to timeless gambles. This functional form is simpler than the one used by BT, v . It captures loss aversion, but ignores other elements of prospect theory, such as the concavity (convexity) over gains (losses) and the probability transformation. In part this is because it is difficult to incorporate all these features into a fully dynamic framework; but also, it is based on BT's observation that it is mainly loss aversion that drives their results.¹⁷

¹⁷ The $b_0 \bar{C}_t^{-\gamma}$ coefficient on the loss aversion term is a scaling factor which ensures that risk premia in the economy remain stationary even as aggregate wealth increases over time. It involves per capita

BHS show that loss aversion can indeed provide a partial explanation of the high Sharpe ratio on the aggregate stock market. However, how much of the Sharpe ratio it can explain depends heavily on the importance of the second source of utility in Equation (11), or in short, on b_0 . As a way of thinking about this parameter, BHS note that when $b_0 = 0.7$, the psychological pain of losing \$100 in the stock market, captured by the second term, is roughly equal to the consumption-related pain of having to consume \$100 less, captured by the first term. For this b_0 , the Sharpe ratio of the risky asset is 0.11, about a third of its historical value.

BT and BHS are both effectively assuming that investors engage in narrow framing, both cross-sectionally and temporally. Even if they have many forms of wealth, both financial and non-financial, they still get utility from changes in the value of one specific component of their total wealth: financial wealth in the case of BT, and stock holdings in the case of BHS. And even if investors have long investment horizons, they still evaluate their portfolio returns on an annual basis.

The assumption about cross-sectional narrow framing can be motivated in a number of ways. The simplest possibility is that it captures non-consumption utility, such as regret. Regret is the pain we feel when we realize that we would be better off if we had not taken a certain action in the past. If the investor's stock holdings fall in value, he may regret the specific decision he made to invest in stocks. Such feelings are naturally captured by defining utility directly over changes in the investors' financial wealth or in the value of his stock holdings.

Another possibility is that while people actually care only about consumption-related utility, they are boundedly rational. For example, suppose that they are concerned that their consumption might fall below some habit level. They know that the right thing to do when considering a stock market investment is to merge the stock market risk with other pre-existing risks that they face – labor income risk, say – and then to compute the likelihood of consumption falling below habit. However, this calculation may be too complex. As a result, people may simply focus on gains and losses in stock market wealth alone, rather than on gains and losses in total wealth.

What about temporal narrow framing? We suggested above that the way information is presented may lead investors to care about annual changes in financial wealth even if they have longer investment horizons. To provide further evidence for this, Thaler, Tversky, Kahneman and Schwartz (1997) provide an *experimental* test of the idea that the manner in which information is presented affects the frame people adopt in their decision-making.¹⁸

consumption \bar{C}_t , which is exogenous to the investor, and so does not affect the intuition of the model. The constant b_0 controls the importance of the loss aversion term in the investor's preferences; setting $b_0 = 0$ reduces the model to the much studied case of power utility over consumption. As $b_0 \rightarrow \infty$, the investor's decisions are driven primarily by concern about gains and losses in financial wealth, as assumed by Benartzi and Thaler.

¹⁸ See also Gneezy and Potters (1997) for a similar experiment.

In their experiment, subjects are asked to imagine that they are portfolio managers for a small college endowment. One group of subjects – Group I, say – is shown monthly observations on two funds, Fund A and Fund B. Returns on Fund A (B) are drawn from a normal distribution calibrated to mimic bond (stock) returns as closely as possible, although subjects are not given this information. After each monthly observation, subjects are asked to allocate their portfolio between the two funds over the next month. They are then shown the realized returns over that month, and asked to allocate once again.

A second group of investors – Group II – is shown exactly the same series of returns, except that it is aggregated at the annual level; in other words, these subjects do not see the monthly fund fluctuations, but only cumulative annual returns. After each annual observation, they are asked to allocate their portfolio between the two funds over the next year.

A final group of investors – Group III – is shown exactly the same data, this time aggregated at the five-year level, and they too are asked to allocate their portfolio after each observation.

After going through a total of 200 months worth of observations, each group is asked to make one final portfolio allocation, which is to apply over the next 400 months. Thaler et al. (1997) find that the average final allocation chosen by subjects in Group I is much lower than that chosen by people in Groups II and III. This result is consistent with the idea that people code gains and losses based on how information is presented to them. Subjects in Group I see monthly observations and hence more frequent losses. If they adopt the monthly distribution as a frame, they will be more wary of stocks and will allocate less to them.

4.1.2. Ambiguity aversion

In Section 3, we presented the Ellsberg paradox as evidence that people dislike ambiguity, or situations where they are not sure what the probability distribution of a gamble is. This is potentially very relevant for finance, as investors are often uncertain about the distribution of a stock's return.

Following the work of Ellsberg, many models of how people react to ambiguity have been proposed; Camerer and Weber (1992) provide a comprehensive review. One of the more popular approaches is to suppose that when faced with ambiguity, people entertain a range of possible probability distributions and act to maximize the minimum expected utility under any candidate distribution. In effect, people behave as if playing a game against a malevolent opponent who picks the actual distribution of the gamble so as to leave them as worse off as possible. Such a decision rule was first axiomatized by Gilboa and Schmeidler (1989). Epstein and Wang (1994) showed how such an approach could be incorporated into a dynamic asset pricing model, although they did not try to assess the quantitative implications of ambiguity aversion for asset prices.

Quantitative implications *have* been derived using a closely related framework known as robust control. In this approach, the agent has a reference probability

distribution in mind, but wants to ensure that his decisions are good ones even if the reference model is misspecified to some extent. Here too, the agent essentially tries to guard against a “worst-case” misspecification. Anderson, Hansen and Sargent (1998) show how such a framework can be used for portfolio choice and pricing problems, even when state equations and objective functions are nonlinear.

Maenhou (1999) applies the Anderson et al. framework to the specific issue of the equity premium. He shows that if investors are concerned that their model of stock returns is misspecified, they will charge a substantially higher equity premium as compensation for the perceived ambiguity in the probability distribution. He notes, however, that to explain the full 3.9% equity premium requires an unreasonably high concern about misspecification. At best then, ambiguity aversion is only a partial resolution of the equity premium puzzle.

4.2. The volatility puzzle

Before turning to behavioral work on the volatility puzzle, it is worth thinking about how rational approaches to this puzzle might proceed. Since, in the data, the volatility of returns is higher than the volatility of dividend growth, Equation (8) makes it clear that we have to make up the gap by introducing variation in the price–dividend ratio. What are the different ways we might do this? A useful framework for thinking about this is a version of the present value formula originally derived by Campbell and Shiller (1988). Starting from

$$R_{t+1} = \frac{P_{t+1} + D_{t+1}}{P_t}, \quad (13)$$

where P_t is the value of the stock market at time t , they use a log-linear approximation to show that the log price–dividend ratio can be written

$$p_t - d_t = E_t \sum_{j=0}^{\infty} \rho^j \Delta d_{t+1+j} - E_t \sum_{j=0}^{\infty} \rho^j r_{t+1+j} + E_t \lim_{j \rightarrow \infty} \rho^j (p_{t+j} - d_{t+j}) + \text{const.}, \quad (14)$$

where lower case letters represent log variables – $p_t = \log P_t$, for example – and where $\Delta d_{t+1} = d_{t+1} - d_t$.

If the price–dividend ratio is stationary, so that the third term on the right is zero, this equation shows clearly that there are just two reasons price–dividend ratios can move around: changing expectations of future dividend growth or changing discount rates. Discount rates, in turn, can change because of changing expectations of future risk-free rates, changing forecasts of risk or changing risk aversion.

While there appear to be many ways of introducing variation in the P/D ratio, it has become clear that most of them cannot form the basis of a rational explanation of the volatility puzzle. We cannot use changing forecasts of dividend growth to drive the P/D ratio: restating the argument of Shiller (1981) and LeRoy and Porter (1981), if

these forecasts are indeed rational, it must be that P/D ratios predict cash-flow growth in the time series, which they do not.¹⁹ Nor can we use changing forecasts of future risk-free rates: again, if the forecasts are rational, P/D ratios must predict interest rates in the time series, which they do not. Even changing forecasts of risk cannot work, as there is little evidence that P/D ratios predict changes in risk in the time series. The only story that remains is therefore one about changing risk aversion, and this is the idea behind the Campbell and Cochrane (1999) model of aggregate stock market behavior. They propose a habit formation framework in which changes in consumption relative to habit lead to changes in risk aversion and hence variation in P/D ratios. This variation helps to plug the gap between the volatility of dividend growth and the volatility of returns.

Some rational approaches try to introduce variation in the P/D ratio through the third term on the right in Equation (14). Since this requires investors to expect explosive growth in P/D ratios forever, they are known as models of rational bubbles. The idea is that prices are high today because they are expected to be higher next period; and they are higher next period because they are expected to be higher the period after that, and so on, forever. While such a model might initially seem appealing, a number of papers, most recently Santos and Woodford (1997), show that the conditions under which rational bubbles can survive are extremely restrictive.²⁰

We now discuss some of the behavioral approaches to the volatility puzzle, grouping them by whether they focus on beliefs or on preferences.

4.2.1. Beliefs

One possible story is that investors believe that the mean dividend growth rate is more variable than it actually is. When they see a surge in dividends, they are too quick to believe that the mean dividend growth rate has increased. Their exuberance pushes prices up relative to dividends, adding to the volatility of returns.

A story of this kind can be derived as a direct application of representativeness and in particular, of the version of representativeness known as the law of small numbers, whereby people expect even short samples to reflect the properties of the parent population. If the investor sees many periods of good earnings, the law of small numbers leads him to believe that earnings growth has gone up, and hence that earnings

¹⁹ There is an important caveat to the statement that changing cash-flow forecasts cannot be the basis of a satisfactory solution to the volatility puzzle. A large literature on structural uncertainty and learning, in which investors do not know the parameters of the cash-flow process but learn them over time, has had some success in matching the empirical volatility of returns [Brennan and Xia (2001), Veronesi (1999)]. In these models, variation in price-dividend ratios comes precisely from changing forecasts of cash-flow growth. While these forecasts are not subsequently confirmed in the data, investors are not considered irrational – they simply don't have enough data to infer the correct model. In related work, Barsky and De Long (1993) generate return volatility in an economy where investors forecast cash flows using a model that is wrong, but not easily rejected with available data.

²⁰ Brunnermeier (2001) provides a comprehensive review of this literature.

will continue to be high in the future. After all, the earnings growth rate cannot be “average”. If it were, then according to the law of small numbers, earnings should appear average, even in short samples: some good earnings news, some bad earnings news, but not several good pieces of news in a row.

Another belief-based story relies more on private, rather than public information, and in particular, on overconfidence about private information. Suppose that an investor has seen public information about the economy, and has formed a prior opinion about future cash-flow growth. He then does some research on his own and becomes overconfident about the information he gathers: he overestimates its accuracy and puts too much weight on it relative to his prior. If the private information is positive, he will push prices up too high relative to current dividends, again adding to return volatility.²¹

Price–dividend ratios and returns might also be excessively volatile because investors extrapolate *past returns* too far into the future when forming expectations of future returns. Such a story might again be based on representativeness and the law of small numbers. The same argument for why investors might extrapolate past cash flows too far into the future can be applied here to explain why they might do the same thing with past returns.

The reader will have noticed that we do not cite any specific papers in connection with these behavioral stories. This is because these ideas were originally put forward in papers whose primary focus is explaining *cross-sectional* anomalies such as the value premium, even though they also apply here in a natural way. In brief, many of those papers – which we discuss in detail in Section 5 – generate certain cross-sectional anomalies by building excessive time series variation into the price–earnings ratios of individual stocks. It is therefore not surprising that the mechanisms proposed there might also explain the substantial time series variation in *aggregate*-level price–earnings ratios. In fact, it is perhaps satisfying that these behavioral theories simultaneously address both aggregate and firm-level evidence.

We close this section with a brief mention of “money illusion”, the confusion between real and nominal values first discussed by Fisher (1928), and more recently investigated by Shafir et al. (1997). In financial markets, Modigliani and Cohn (1979) and more recently, Ritter and Warr (2002), have argued that part of the variation in *P/D* ratios and returns may be due to investors mixing real and nominal quantities when forecasting future cash flows. The value of the stock market can be determined

²¹ Campbell (2000), among others, notes that behavioral models based on cash-flow forecasts often ignore potentially important interest rate effects. If investors are forecasting excessively high cash-flow growth, pushing up prices, interest rates should also rise, thereby dampening the price rise. One response is that interest rates are governed by expectations about *consumption* growth, and in the short run, consumption and dividends can be somewhat delinked: even if dividend growth is expected to be high, this need not necessarily trigger an immediate interest rate response. Alternatively, one can try to specify investors’ expectations in such a way that interest rate effects become less important. Cecchetti, Lam and Mark (2000) take a step in this direction.

by discounted real cash flows at real rates, or nominal cash flows at nominal rates. At times of especially high or especially low inflation though, it is possible that some investors mistakenly discount *real* cash flows at *nominal* rates. If inflation increases, so will the nominal discount rate. If investors then discount the *same* set of cash flows at this higher rate, they will push the value of the stock market down. Of course, this calculation is incorrect: the same inflation which pushes up the discount rate should also push up future cash flows. On net, inflation should have little effect on market value. Such real vs. nominal confusion may therefore cause excessive variation in *P/D* ratios and returns and seems particularly relevant to understanding the low market valuations during the high inflation years of the 1970s, as well as the high market valuations during the low inflation 1990s.

4.2.2. Preferences

Barberis, Huang and Santos (2001) show that a straightforward extension of the version of their model discussed in Section 4.1 can explain both the equity premium and volatility puzzles. To do this, they appeal to experimental evidence about dynamic aspects of loss aversion. This evidence suggests that the degree of loss aversion is not the same in all circumstances but depends on prior gains and losses. In particular, Thaler and Johnson (1990) find that after prior gains, subjects take on gambles they normally do not, and that after prior losses, they refuse gambles that they normally accept. The first finding is sometimes known as the “house money effect”, reflecting gamblers’ increasing willingness to bet when ahead. One interpretation of this evidence is that losses are less painful after prior gains because they are cushioned by those gains. However, after being burned by a painful loss, people may become more wary of additional setbacks.²²

To capture these ideas, Barberis, Huang and Santos (2001) modify the utility function in Equation (11) to

$$E_0 \sum_{t=0}^{\infty} \left[\rho' \frac{C_t^{1-\gamma}}{1-\gamma} + b_0 \bar{C}_t^{-\gamma} \tilde{v}(X_{t+1}, z_t) \right]. \quad (15)$$

Here, z_t is a state variable that tracks past gains and losses on the stock market. For any fixed z_t , the function \tilde{v} is a piecewise linear function similar in form to \hat{v} , defined in Equation (12). However, the investors’ sensitivity to losses is no longer constant at

²² It is important to distinguish Thaler and Johnson’s (1990) evidence from other evidence presented by Kahneman and Tversky (1979) and discussed in Section 3, showing that people are risk averse over gains and risk seeking over losses. One set of evidence pertains to one-shot gambles, the other to sequences of gambles. Kahneman and Tversky’s (1979) evidence suggests that people are willing to take risks in order to avoid a loss; Thaler and Johnson’s (1990) evidence suggests that if these efforts are unsuccessful and the investor suffers an unpleasant loss, he will subsequently act in a more risk averse manner.

2.25, but is determined by z_t , in a way that reflects the experimental evidence described above.

A model of this kind can help explain the volatility puzzle. Suppose that there is some good cash-flow news. This pushes the stock market up, generating prior gains for investors, who are now less scared of stocks: any losses will be cushioned by the accumulated gains. They therefore discount future cash flows at a lower rate, pushing prices up still further relative to current dividends and adding to return volatility.

5. Application: The cross-section of average returns

While the behavior of the aggregate stock market is not easy to understand from the rational point of view, promising rational models have nonetheless been developed and can be tested against behavioral alternatives. Empirical studies of the behavior of *individual* stocks have unearthed a set of facts which is altogether more frustrating for the rational paradigm. Many of these facts are about the *cross-section* of average returns: they document that one group of stocks earns higher average returns than another. These facts have come to be known as “anomalies” because they cannot be explained by the simplest and most intuitive model of risk and return in the financial economist’s toolkit, the Capital Asset Pricing Model, or CAPM.

We now outline some of the more salient findings in this literature and then consider some of the rational and behavioral approaches in more detail.

The size premium. This anomaly was first documented by Banz (1981). We report the more recent findings of Fama and French (1992). Every year from 1963 to 1990, Fama and French group all stocks traded on the NYSE, AMEX, and NASDAQ into deciles based on their market capitalization, and then measure the average return of each decile over the next year. They find that for this sample period, the average return of the smallest stock decile is 0.74% per month higher than the average return of the largest stock decile. This is certainly an anomaly relative to the CAPM: while stocks in the smallest decile do have higher betas, the difference in risk is not enough to explain the difference in average returns.²³

Long-term reversals. Every three years from 1926 to 1982, De Bondt and Thaler (1985) rank all stocks traded on the NYSE by their prior three year cumulative return and form two portfolios: a “winner” portfolio of the 35 stocks with the best prior record and a “loser” portfolio of the 35 worst performers. They then measure the average return of these two portfolios over the three years subsequent to their formation. They

²³ The last decade of data has served to reduce the size premium considerably. Gompers and Metrick (2001) argue that this is due to demand pressure for large stocks resulting from the growth of institutional investors, who prefer such stocks.

find that over the whole sample period, the average annual return of the loser portfolio is higher than the average return of the winner portfolio by almost 8% per year.

The predictive power of scaled-price ratios. These anomalies, which are about the cross-sectional predictive power of variables like the book-to-market (B/M) and earnings-to-price (E/P) ratios, where some measure of fundamentals is scaled by price, have a long history in finance going back at least to Graham (1949), and more recently Dreman (1977), Basu (1983) and Rosenberg, Reid and Lanstein (1985). We concentrate on Fama and French's (1992) more recent evidence.

Every year, from 1963 to 1990, Fama and French group all stocks traded on the NYSE, AMEX and NASDAQ into deciles based on their book-to-market ratio, and measure the average return of each decile over the next year. They find that the average return of the highest-B/M-ratio decile, containing so called "value" stocks, is 1.53% per month higher than the average return on the lowest-B/M-ratio decile, "growth" or "glamour" stocks, a difference much higher than can be explained through differences in beta between the two portfolios. Repeating the calculations with the earnings–price ratio as the ranking measure produces a difference of 0.68% per month between the two extreme decile portfolios, again an anomalous result.²⁴

Momentum. Every month from January 1963 to December 1989, Jegadeesh and Titman (1993) group all stocks traded on the NYSE into deciles based on their prior six month return and compute average returns of each decile over the six months after portfolio formation. They find that the decile of biggest prior winners outperforms the decile of biggest prior losers by an average of 10% on an annual basis.

Comparing this result to De Bondt and Thaler's (1985) study of prior winners and losers illustrates the crucial role played by the length of the prior ranking period. In one case, prior winners continue to win; in the other, they perform poorly.²⁵ A challenge to both behavioral and rational approaches is to explain why extending the formation period switches the result in this way.

There is some evidence that tax-loss selling creates seasonal variation in the momentum effect. Stocks with poor performance during the year may later be subject to selling by investors keen to realize losses that can offset capital gains elsewhere. This selling pressure means that prior losers continue to lose, enhancing the momentum effect. At the turn of the year, though, the selling pressure eases off, allowing prior losers to rebound and weakening the momentum effect. A careful analysis by Grinblatt

²⁴ Ball (1978) and Berk (1995) point out that the size premium and the scaled-price ratio effects emerge naturally in any model where investors apply different discount rates to different stocks: if investors discount a stock's cash flows at a higher rate, that stock will typically have a lower market capitalization and a lower price-earnings ratio, but also higher returns. Note, however, that this view does not shed any light on whether the variation in discount rates is rationally justifiable or not.

²⁵ In fact, De Bondt and Thaler (1985) also report that one-year big winners outperform one-year big losers over the following year, but do not make much of this finding.

and Moskowitz (1999) finds that on net, tax-loss selling may explain part of the momentum effect, but by no means all of it. In any case, while selling a stock for tax purposes is rational, a model of predictable price movements based on such behavior is not. Roll (1983) calls such explanations “stupid” since investors would have to be stupid not to buy in December if prices were going to increase in January.

A number of studies have examined stock returns following important corporate announcements, a type of analysis known as an event study. Chapter 5 in this Handbook discusses many of these studies in detail; here, we summarize them briefly.

Event studies of earnings announcements. Every quarter from 1974 to 1986, Bernard and Thomas (1989) group all stocks traded on the NYSE and AMEX into deciles based on the size of the surprise in their most recent earnings announcement. “Surprise” is measured relative to a simple random walk model of earnings. They find that on average, over the 60 days after the earnings announcement, the decile of stocks with surprisingly good news outperforms the decile with surprisingly bad news by an average of about 4%, a phenomenon known as post-earnings announcement drift. Once again, this difference in returns is not explained by differences in beta between the two portfolios. A later study by Chan, Jegadeesh and Lakonishok (1996) measures surprise in other ways – relative to analyst expectations, and by the stock price reaction to the news – and obtains similar results.²⁶

Event studies of dividend initiations and omissions. Michaely, Thaler and Womack (1995) study firms which announced initiation or omission of a dividend payment between 1964 and 1988. They find that on average, the shares of firms initiating (omitting) dividends significantly outperform (underperform) the market portfolio over the year after the announcement.

Event studies of stock repurchases. Ikenberry, Lakonishok and Vermaelen (1995) look at firms which announced a share repurchase between 1980 and 1990, while Mitchell and Stafford (2001) study firms which did either self-tenders or share repurchases between 1960 and 1993. The latter study finds that on average, the shares of these firms outperform a control group matched on size and book-to-market by a substantial margin over the four year period following the event.

Event studies of primary and secondary offerings. Loughran and Ritter (1995) study firms which undertook primary or secondary equity offerings between 1970 and 1990.

²⁶ Vuolteenaho (2002) combines a clean-surplus accounting version of the present value formula with Campbell's (1991) log-linear decomposition of returns to estimate a measure of cash-flow news that is potentially more accurate than earnings announcements. Analogous to the post-earnings announcement studies, he finds that stocks with good cash-flow news subsequently have higher average returns than stocks with disappointing cash-flow news.

They find that the average return of shares of these firms over the five-year period after the issuance is markedly below the average return of shares of non-issuing firms matched to the issuing firms on size. Brav and Gompers (1997) and Brav, Geczy and Gompers (2000) argue that this anomaly may not be distinct from the scaled-price anomaly listed above: when the returns of event firms are compared to the returns of firms matched on both size and book-to-market, there is very little difference.

Long-term event studies like the last three analyses summarized above raise some thorny statistical problems. In particular, conducting statistical inference with long-term buy-and-hold post-event returns is a treacherous business. Barber and Lyon (1997), Lyon, Barber and Tsai (1999), Brav (2000), Fama (1998), Loughran and Ritter (2000) and Mitchell and Stafford (2001) are just a few of the papers that discuss this topic. Cross-sectional correlation is one important issue: if a certain firm announces a share repurchase shortly after another firm does, their four-year post event returns will overlap and cannot be considered independent. Although the problem is an obvious one, it is not easy to deal with effectively. Some recent attempts to do so, such as Brav (2000), suggest that the anomalous evidence in the event studies on dividend announcements, repurchase announcements, and equity offerings is statistically weaker than initially thought, although how much weaker remains controversial.

A more general concern with *all* the above empirical evidence is data-mining. After all, if we sort and rank stocks in enough different ways, we are bound to discover striking – but completely spurious – cross-sectional differences in average returns.

A first response to the data-mining critique is to note that the above studies do not use the kind of obscure firm characteristics or marginal corporate announcements that would suggest data-mining. Indeed, it is hard to think of an important class of corporate announcements that has *not* been associated with a claim about anomalous post-event returns. A more direct check is to perform out-of-sample tests. Interestingly, a good deal of the above evidence *has* been replicated in other data sets. Fama, French and Davis (2000) show that there is a value premium in the subsample of U.S. data that precedes the data set used in Fama and French (1992), while Fama and French (1998) document a value premium in international stock markets. Rouwenhorst (1998) shows that the momentum effect is alive and well in international stock market data.

If the empirical results are taken at face value, then the challenge to the rational paradigm is to show that the above cross-sectional evidence emerges naturally from a model with fully rational investors. In special cases, models of this form reduce to the CAPM, and we know that this does not explain the evidence. More generally, rational models predict a multifactor pricing structure,

$$\bar{r}_i - r_f = \beta_{i,1} (\bar{F}_1 - r_f) + \cdots + \beta_{i,K} (\bar{F}_K - r_f), \quad (16)$$

where the factors proxy for marginal utility growth and where the loadings $\beta_{i,k}$ come from a time series regression of excess stock returns on excess factor returns,

$$r_{i,t} - r_{f,t} = \alpha_i + \beta_{i,1} (F_{1,t} - r_{f,t}) + \cdots + \beta_{i,K} (F_{K,t} - r_{f,t}) + \varepsilon_{i,t}. \quad (17)$$

To date, it has proved difficult to derive a multi-factor model which explains the cross-sectional evidence, although this remains a major research direction.

Alternatively, one can skip the step of *deriving* a factor model, and simply try a specific model to see how it does. This is the approach of Fama and French (1993, 1996). They show that a certain three factor model does a good job explaining the average returns of portfolios formed on size and book-to-market rankings. Put differently, the α_i intercepts in regression (17) are typically close to zero for these portfolios and for their choice of factors. The specific factors they use are the return on the market portfolio, the return on a portfolio of small stocks minus the return on a portfolio of large stocks – the “size” factor – and the return on a portfolio of value stocks minus the return on a portfolio of growth stocks – the “book-to-market” factor. By constructing these last two factors, Fama and French are isolating common factors in the returns of small stocks and value stocks, and their three factor model can be loosely motivated by the idea that this comovement is a systematic risk that is priced in equilibrium.

The low α_i intercepts obtained by Fama and French (1993, 1996) are not necessarily cause for celebration. After all, as Roll (1977) emphasizes, in any specific sample, it is always possible to mechanically construct a one factor model that prices average returns *exactly*.²⁷ This sounds a cautionary note: just because a factor model happens to work well does not necessarily mean that we are learning anything about the economic drivers of average returns. To be fair, Fama and French (1996) themselves admit that their results can only have their full impact once it is explained what it is about investor preferences and the structure of the economy that leads people to price assets according to their model.

One general feature of the rational approach is that it is loadings or betas, and not firm characteristics, that determine average returns. For example, a risk-based approach would argue that value stocks earn high returns not because they have high book-to-market ratios, but because such stocks happen to have a high loading on the book-to-market factor. Daniel and Titman (1997) cast doubt on this specific prediction by performing double sorts of stocks on both book-to-market ratios and loadings on book-to-market factors, and showing that stocks with different loadings but the same book-to-market ratio do *not* differ in their average returns. These results appear quite damaging to the rational approach, but they have also proved controversial. Using a longer data set and a different methodology, Fama, French and Davis (2000) claim to reverse Daniel and Titman’s findings.

More generally, rational approaches to the cross-sectional evidence face a number of other obstacles. First, rational models typically measure risk as the covariance of

²⁷ For any sample of observations on individual returns, choose any one of the ex-post mean-variance efficient portfolios. Roll (1977) shows that there is an exact linear relationship between the sample mean returns of the individual assets and their betas, computed with respect to the mean-variance efficient portfolio.

returns with marginal utility of consumption. Stocks are risky if they fail to pay out at times of high marginal utility – in “bad” times – and instead pay out when marginal utility is low – in “good” times. The problem is that for many of the above findings, there is little evidence that the portfolios with anomalously *high* average returns do poorly in bad times, whatever plausible measure of bad times is used. For example, Lakonishok, Shleifer and Vishny (1994) show that in their 1968 to 1989 sample period, value stocks do well relative to growth stocks even when the economy is in recession. Similarly, De Bondt and Thaler (1987) find that their loser stocks have higher betas than winners in up markets and lower betas in down markets – an attractive combination that no one would label “risky”.

Second, some of the portfolios in the above studies – the decile of stocks with the lowest book-to-market ratios for example – earn average returns below the risk-free rate. It is not easy to explain why a rational investor would willingly accept a lower return than the T-Bill rate on a volatile portfolio.

Third, Chopra, Lakonishok and Ritter (1992) and La Porta et al. (1997) show that a large fraction of the high (low) average returns to prior losers (winners) documented by De Bondt and Thaler (1985), and of the high (low) returns to value (growth) stocks, is earned over a very small number of days around earnings announcements. It is hard to tell a rational story for why the premia should be concentrated in this way, given that there is no evidence of changes in *systematic* risk around earnings announcements.

Finally, in some of the examples given above, it is not just that one portfolio outperforms another on average. In some cases, the outperformance is present in almost every period of the sample. For example, in Bernard and Thomas’ (1989) study, firms with surprisingly good earnings outperform those with surprisingly poor earnings in 46 out of the 50 quarters studied. It is not easy to see any risk here that might justify the outperformance.

5.1. Belief-based models

There are a number of behavioral models which try to explain some of the above phenomena. We classify them based on whether their mechanism centers on beliefs or on preferences.

Barberis, Shleifer and Vishny (1998), BSV henceforth, argue that much of the above evidence is the result of systematic errors that investors make when they use public information to form expectations of future cash flows. They build a model that incorporates two of the updating biases from Section 3: conservatism, the tendency to underweight new information relative to priors; and representativeness, and in particular the version of representativeness known as the law of small numbers, whereby people expect even short samples to reflect the properties of the parent population.

When a company announces surprisingly good earnings, conservatism means that investors react insufficiently, pushing the price up too little. Since the price is too low, subsequent returns will be higher on average, thereby generating both post-earnings

announcement drift and momentum. After a *series* of good earnings announcements, though, representativeness causes people to overreact and push the price up too high. The reason is that after many periods of good earnings, the law of small numbers leads investors to believe that this is a firm with particularly high earnings growth, and hence to forecast high earnings in the future. After all, the firm cannot be “average”. If it were, then according to the law of small numbers, its earnings should *appear* average, even in short samples. Since the price is now too high, subsequent returns are too low on average, thereby generating long-term reversals and a scaled-price ratio effect.

To capture these ideas mathematically, BSV consider a model with a representative risk-neutral investor in which the true earnings process for all assets is a random walk. Investors, however, do not use the random-walk model to forecast future earnings. They think that at any time, earnings are being generated by one of two regimes: a “mean-reverting” regime, in which earnings are more mean-reverting than in reality, and a “trending” regime in which earnings trend more than in reality. The investor believes that the regime generating earnings changes exogenously over time and sees his task as trying to figure out which of the two regimes is currently generating earnings.

This framework offers one way of modelling the updating biases described above. Including a “trending” regime in the model captures the effect of representativeness by allowing investors to put more weight on trends than they should. Conservatism suggests that people may put too little weight on the latest piece of earnings news relative to their prior beliefs. In other words, when they get a good piece of earnings news, they effectively act as if part of the shock will be reversed in the next period, in other words, as if they believe in a “mean-reverting” regime. BSV confirm that for a wide range of parameter values, this model does indeed generate post-earnings announcement drift, momentum, long-term reversals and cross-sectional forecasting power for scaled-price ratios.²⁸

Daniel, Hirshleifer and Subrahmanyam (1998, 2001), DHS henceforth, stress biases in the interpretation of *private*, rather than public information. Imagine that the investor does some research on his own to try to determine a firm’s future cash flows. DHS assume that he is overconfident about this information; in particular, they argue that investors are more likely to be overconfident about private information they have worked hard to generate than about public information. If the private information is positive, overconfidence means that investors will push prices up too far relative to fundamentals. Future public information will slowly pull prices back to their correct value, thus generating long-term reversals and a scaled-price effect. To get momentum and a post-earnings announcement effect, DHS assume that the public information alters the investor’s confidence in his original private information in

²⁸ Potoshman (2001) finds evidence of a BSV-type expectations formation process in the options market. He shows that when pricing options, traders appear to underreact to individual daily changes in instantaneous variance, while overreacting to longer sequences of increasing or decreasing changes in instantaneous variance.

an asymmetric fashion, a phenomenon known as self-attribution bias: public news which confirms the investor's research strongly increases the confidence he has in that research. Disconfirming public news, though, is given less attention, and the investor's confidence in the private information remains unchanged. This asymmetric response means that initial overconfidence is on average followed by even greater overconfidence, generating momentum.

If, as BSV and DHS argue, long-term reversals and the predictive power of scaled-price ratios are driven by excessive optimism or pessimism about future cash flows followed by a correction, then most of the correction should occur at those times when investors find out that their initial beliefs were too extreme, in other words, at earnings announcement dates. The findings of Chopra, Lakonishok and Ritter (1992) and La Porta et al. (1997), who show that a large fraction of the premia to prior losers and to value stocks is earned around earnings announcement days, strongly confirm this prediction.

Perhaps the simplest way of capturing much of the cross-sectional evidence is positive feedback trading, where investors buy more of an asset that has recently gone up in value [De Long et al. (1990b), Barberis and Shleifer (2003)]. If a company's stock price goes up this period on good earnings, positive feedback traders buy the stock in the following period, causing a further price rise. On the one hand, this generates momentum and post-earnings announcement drift. On the other hand, since the price has now risen above what is justified by fundamentals, subsequent returns will on average be too low, generating long-term reversals and a scaled-price ratio effect.

The simplest way of motivating positive feedback trading is extrapolative expectations, where investors' expectations of future returns are based on past returns. This in turn, may be due to representativeness and to the law of small numbers in particular. The same argument made by BSV as to why investors might extrapolate past cash flows too far into the future can be applied here to explain why they might extrapolate past *returns* too far into the future. De Long et al. (1990b) note that institutional features such as portfolio insurance or margin calls can also generate positive feedback trading.

Positive feedback trading also plays a central role in the model of Hong and Stein (1999), although in this case it emerges endogenously from more primitive assumptions. In this model, two boundedly rational groups of investors interact, where bounded rationality means that investors are only able to process a subset of available information. "Newswatchers" make forecasts based on private information, but do not condition on past prices. "Momentum traders" condition only on the most recent price change.

Hong and Stein also assume that private information diffuses slowly through the population of newswatchers. Since these investors are unable to extract each others' private information from prices, the slow diffusion generates momentum. Momentum traders are then added to the mix. Given what they are allowed to condition on, their optimal strategy is to engage in positive feedback trading: a price increase last period is a sign that good private information is diffusing through the economy. By buying,

momentum traders hope to profit from the continued diffusion of information. This behavior preserves momentum, but also generates price reversals: since momentum traders cannot observe the extent of news diffusion, they keep buying even after price has reached fundamental value, generating an overreaction that is only later reversed.

These four models differ most in their explanation of momentum. In two of the models – BSV and Hong and Stein (1999) – momentum is due to an initial underreaction followed by a correction. In De Long et al. (1990b) and DHS, it is due to an initial overreaction followed by even more overreaction. Within each pair, the stories are different again.²⁹

Hong, Lim and Stein (2000) present supportive evidence for the view of Hong and Stein (1999) that momentum is due simply to slow diffusion of private information through the economy. They argue that the diffusion of information will be particularly slow among small firms and among firms with low analyst coverage, and that the momentum effect should therefore be more prominent there, a prediction they confirm in the data. They also find that among firms with low analyst coverage, momentum is almost entirely driven by prior losers continuing to lose. They argue that this, too, is consistent with a diffusion story. If a firm not covered by analysts is sitting on good news, it will do its best to convey the news to as many people as possible, and as quickly as possible; bad news, however, will be swept under the carpet, making its diffusion much slower.

5.2. Belief-based models with institutional frictions

Some authors have argued that models which combine mild assumptions about investor irrationality with institutional frictions may offer a fruitful way of thinking about some of the anomalous cross-sectional evidence.

The institutional friction that has attracted the most attention is short-sale constraints. As mentioned in Section 2.2, these can be thought of as anything which makes investors less willing to establish a short position than a long one. They include the direct cost of shorting, namely the lending fee; the risk that the loan is recalled by the lender at an inopportune moment; as well as legal restrictions: a large fraction of mutual funds are not allowed to short stocks.

Several papers argue that when investors differ in their beliefs, the existence of short-sale constraints can generate deviations from fundamental value and in particular, explain why stocks with high price–earnings ratios earn lower average returns in the cross-section. The simplest way of motivating the assumption of heterogeneous beliefs is overconfidence, which is why that assumption is often thought of as capturing a mild form of irrationality. In the absence of overconfidence, investors’ beliefs converge

²⁹ In particular, the models make different predictions about how individual investors would trade following certain sequences of past returns. Armed with transaction-level data, Hvidkjaer (2001) exploits this to provide initial evidence that may distinguish the theories.

rapidly as they hear each other's opinions and hence deduce each other's private information.

There are at least two mechanisms through which differences of opinion and short-sale constraints can generate price-earnings ratios that are too high, and thereby explain why price-earnings ratios predict returns in the cross-section.

Miller (1977) notes that when investors hold different views about a stock, those with bullish opinions will, of course, take long positions. Bearish investors, on the other hand, want to short the stock, but being unable to do so, they sit out of the market. Stock prices therefore reflect only the opinions of the most optimistic investors which, in turn, means that they are too high and that they will be followed by lower returns.

Harrison and Kreps (1978) and Scheinkman and Xiong (2003) argue that in a dynamic setting, a second, speculation-based mechanism arises. They show that when there are differences in beliefs, investors will be happy to buy a stock for more than its fundamental value in anticipation of being able to sell it later to other investors even more optimistic than themselves. Note that short-sale constraints are essential to this story: in their absence, an investor can profit from another's greater optimism by simply shorting the stock. With short-sale constraints, the only way to do so is to buy the stock first, and then sell it on later.

Both types of models make the intriguing prediction that stocks which investors disagree about more will have higher price-earnings ratios and lower subsequent returns. Three recent papers test this prediction, each using a different measure of differences of opinion.

Diether, Malloy and Scherbina (2002) use IBES data on analyst forecasts to obtain a direct measure of heterogeneity of opinion. They group stocks into quintiles based on the level of dispersion in analysts' forecasts of current year earnings and confirm that the highest dispersion portfolio earns lower average returns than the lowest dispersion portfolio.

Chen, Hong and Stein (2002) use "breadth of ownership" – defined roughly as the fraction of mutual funds that hold a particular stock – as a proxy for divergence of opinion about the stock. The more dispersion in opinions there is, the more mutual funds will need to sit out the market due to short sales constraints, leading to lower breadth. Chen et al. predict, and confirm in the data, that stocks experiencing a decrease in breadth subsequently have lower average returns compared to stocks whose breadth increases.

Jones and Lamont (2002) use the cost of short-selling a stock – in other words, the lending fee – to measure differences of opinion about that stock. The idea is that if there is a lot of disagreement about a stock's prospects, many investors will want to short the stock, thereby pushing up the cost of doing so. Jones and Lamont confirm that stocks with higher lending fees have higher price-earnings ratios and earn lower subsequent returns. It is interesting to note that their data set spans the years from 1926 to 1933. At that time, there existed a centralized market for borrowing stocks and lending fees were published daily in the Wall Street Journal. Today, by contrast, stock lending is an over-the-counter market, and data on lending fees is harder to come by.

In other related work, Hong and Stein (2003) show that short-sale constraints and differences of opinion also have implications for higher order moments, in that they can lead to skewness. The intuition is that when a stock's price goes down, more information is revealed: by seeing at what point they enter the market, we learn the valuations of those investors whose pessimistic views could not initially be reflected in the stock price, because of short-sale constraints. When the stock market goes up, the sidelined investors stay out of the market and there is less information revelation. This increase in volatility after a downturn is the source of the skewness.

One prediction of this idea is that stocks which investors disagree about more should exhibit greater skewness. Chen, Hong and Stein (2001) test this idea using increases in turnover as a sign of investor disagreement. They show that stocks whose turnover increases subsequently display greater skewness.

5.3. Preferences

Earlier, we discussed Barberis, Huang and Santos (2001), which tries to explain aggregate stock market behavior by combining loss aversion and narrow framing with an assumption about how the degree of loss aversion changes over time. Barberis and Huang (2001) show that applying the same ideas to individual stocks can generate the evidence on long-term reversals and on scaled-price ratios. The key idea is that when investors hold a number of different stocks, narrow framing may induce them to derive utility from gains and losses in the value of *individual* stocks. The specification of this additional source of utility is exactly the same as in BHS, except that it is now applied at the individual stock level instead of at the portfolio level: the investor is loss averse over individual stock fluctuations and the pain of a loss on a specific stock depends on that stock's past performance.

To see how this model generates a value premium, consider a stock which has had poor returns several periods in a row. Precisely because the investor focuses on individual stock gains and losses, he finds this painful and becomes especially sensitive to the possibility of further losses on the stock. In effect, he perceives the stock as riskier, and discounts its future cash flows at a higher rate: this lowers its price-earnings ratio and leads to higher subsequent returns, generating a value premium. In one sense, this model is narrower than those in the "beliefs" section, Section 5.1, as it does not claim to address momentum. In another sense, it is broader, in that it simultaneously explains the equity premium and derives the risk-free rate endogenously.

The models we describe in Sections 5.1, 5.2 and 5.3 focus primarily on momentum, long-term reversals, the predictive power of scaled-price ratios and post-earnings announcement drift. What about the other examples of anomalous evidence with which we began Section 5? In Section 7, we argue that the long-run return patterns following equity issuance and repurchases may be the result of rational managers responding to the kinds of noise traders analyzed in the preceding behavioral models. In short, if investors cause prices to swing away from fundamental value, managers may try to

time these cycles, issuing equity when it is overpriced, and repurchasing it when it is cheap. In such a world, equity issues will indeed be followed by low returns, and repurchases by high returns. The models we have discussed so far do not, however, shed light on the size anomaly, nor on the dividend announcement event study.

6. Application: Closed-end funds and comovement

6.1. *Closed-end funds*

Closed-end funds differ from more familiar open-end funds in that they only issue a fixed number of shares. These shares are then traded on exchanges: an investor who wants to buy a share of a closed-end fund must go to the exchange and buy it from another investor at the prevailing price. By contrast, should he want to buy a share of an open-end fund, the fund would create a new share and sell it to him at its net asset value, or NAV, the per share market value of its asset holdings.

The central puzzle about closed-end funds is that fund share prices differ from NAV. The typical fund trades at a discount to NAV of about 10% on average, although the difference between price and NAV varies substantially over time. When closed-end funds are created, the share price is typically above NAV; when they are terminated, either through liquidation or open-ending, the gap between price and NAV closes.

A number of rational explanations for the average closed-end fund discount have been proposed. These include expenses, expectations about future fund manager performance, and tax liabilities. These factors can go some way to explaining certain aspects of the closed-end fund puzzle. However, none of them can satisfactorily explain *all* aspects of the evidence. For example, management fees can explain why funds usually sell at discounts, but not why they typically initially sell at a premium, nor why discounts tend to vary from week to week.

Lee, Shleifer and Thaler (1991), LST henceforth, propose a simple behavioral view of these closed-end fund puzzles. They argue that some of the individual investors who are the primary owners of closed-end funds are noise traders, exhibiting irrational swings in their expectations about future fund returns. Sometimes they are too optimistic, while at other times, they are too pessimistic. Changes in their sentiment affect fund share prices and hence also the difference between prices and net asset values.³⁰

This view provides a clean explanation of all aspects of the closed-end fund puzzle. Owners of closed-end funds have to contend with two sources of risk: fluctuations

³⁰ For the noise traders to affect the *difference* between price and NAV rather than just price, it must be that they are more active traders of closed-end fund shares than they are of assets owned by the funds. As evidence for this, LST point out that while funds are primarily owned by individual investors, the funds' assets are not.

in the value of the funds' assets, and fluctuations in noise trader sentiment. If this second risk is systematic – we return to this issue shortly – rational investors will demand compensation for it. In other words, they will require that the fund's shares trade at a discount to NAV.

This also explains why new closed-end funds are often sold at a premium. Entrepreneurs will choose to create closed-end funds at times of investor exuberance, when they know that they can sell fund shares for more than they are worth. On the other hand, when a closed-end fund is liquidated, rational investors no longer have to worry about changes in noise trader sentiment because they know that at liquidation, the fund price will equal NAV. They therefore no longer demand compensation for this risk, and the fund price rises towards NAV.

An immediate prediction of the LST view is that prices of closed-end funds should comove strongly, even if the cash-flow fundamentals of the assets held by the funds do not: if noise traders become irrationally pessimistic, they will sell closed-end funds across the board, depressing their prices regardless of cash-flow news. LST confirm in the data that closed-end fund discounts are highly correlated.

The LST story depends on noise trader risk being systematic. There is good reason to think that it is. If the noise traders who hold closed-end funds also hold other assets, then negative changes in sentiment, say, will drive down the prices of closed-end funds *and* of their other holdings, making the noise trader risk systematic. To check this, LST compute the correlation of closed-end fund discounts with another group of assets primarily owned by individuals, small stocks. Consistent with the noise trader risk being systematic, they find a significant positive correlation.

6.2. Comovement

The LST model illustrates that behavioral models can make interesting predictions not only about the *average* level of returns, but also about patterns of comovement. In particular, it explains why the prices of closed-end funds comove so strongly, and also why closed-end funds as a class comove with small stocks. This raises the hope that behavioral models might be able to explain other puzzling instances of comovement as well.

Before studying this in more detail, it is worth setting out the traditional view of return comovement. This view, derived from economies without frictions and with rational investors, holds that comovement in prices reflects comovement in fundamental values. Since, in a frictionless economy with rational investors, price equals fundamental value – an asset's rationally forecasted cash flows discounted at a rate appropriate for their risk – any comovement in prices must be due to comovement in fundamentals. There is little doubt that many instances of return comovement can be explained by fundamentals: stocks in the automotive industry move together primarily because their earnings are correlated.

The closed-end fund evidence shows that the fundamentals-based view of comovement is at best, incomplete: in that case, the prices of closed-end funds comove even

though their fundamentals do not.³¹ Other evidence is just as puzzling. Froot and Dabora (1999) study “twin stocks”, which are claims to the same cash-flow stream, but are traded in different locations. The Royal Dutch/Shell pair, discussed in Section 2, is perhaps the best known example. If return comovement is simply a reflection of comovement in fundamentals, these two stocks should be perfectly correlated. In fact, as Froot and Dabora show, Royal Dutch comoves strongly with the S&P 500 index of U.S. stocks, while Shell comoves with the FTSE index of UK stocks.

Fama and French (1993) uncover salient common factors in the returns of small stocks, as well as in the returns of value stocks. In order to test the rational view of comovement, Fama and French (1995) investigate whether these strong common factors can be traced to common factors in news about the earnings of these stocks. While they do uncover a common factor in the earnings news of small stocks, as well as in the earnings news of value stocks, these cash-flow factors are weaker than the factors in returns and there is little evidence that the return factors are driven by the cash-flow factors. Once again, there appears to be comovement in returns that has little to do with fundamentals-based comovement.³²

In response to this evidence, researchers have begun to posit behavioral theories of comovement. LST is one such theory. To state their argument more generally, they start by observing that many investors choose to trade only a subset of all available securities. As these investors’ risk aversion or sentiment changes, they alter their exposure to the particular securities they hold, thereby inducing a common factor in the returns of these securities. Put differently, this “habitat” view of comovement predicts that there will be a common factor in the returns of securities that are the primary holdings of a specific subset of investors, such as individual investors. This story seems particularly appropriate for thinking about closed-end funds, and also for Froot and Dabora’s evidence.

A second behavioral view of comovement was recently proposed by Barberis and Shleifer (2003). They argue that to simplify the portfolio allocation process, many investors first group stocks into categories such as small-cap stocks or automotive industry stocks, and then allocate funds across these various categories. If these categories are also adopted by noise traders, then as these traders move funds from

³¹ Bodurtha et al. (1993) and Hardouvelis et al. (1994) provide further interesting examples of a delinking between fundamentals-based comovement and return comovement in the closed-end fund market. They study closed-end *country* funds, whose assets trade in a different location from the funds themselves and find that the funds comove as much with the national stock market in the country where they are traded as with the national stock market in the country where their *assets* are traded. For example, a closed-end fund invested in German equities but traded in the USA typically comoves as much with the U.S. stock market as with the German stock market.

³² In principle, comovement can also be rationally generated through changes in discount rates. However, changes in interest rates or risk aversion induce a common factor in the returns on *all* stocks, and do not explain why a particular group of stocks comoves. A common factor in news about the risk of certain assets may also be a source of comovement for those assets, but there is little direct evidence to support such a mechanism in the case of small stocks or value stocks.

one category to another, the price pressure from their coordinated demand will induce common factors in the returns of stocks that happen to be classified into the same category, even if those stocks' cash flows are largely uncorrelated. In particular, this view predicts that when an asset is added to a category, it should begin to comove more with that category than before.

Barberis, Shleifer and Wurgler (2001) test this "category" view of comovement by taking a sample of stocks that have been added to the S&P 500, and computing the betas of these stocks with the S&P 500 both before and after inclusion. Based on both univariate and multivariate regressions, they show that upon inclusion, a stock's beta with the S&P 500 rises significantly, as does the fraction of its variance that is explained by the S&P 500, while its beta with stocks outside the index falls.³³ This result does not sit well with the cash-flow view of comovement – addition to the S&P 500 is not intended to carry any information about the covariance of a stock's cash flows with other stocks' cash flows – but emerges naturally from a model where prices are affected by category-level demand shocks.

7. Application: Investor behavior

Behavioral finance has also had some success in explaining how certain groups of investors behave, and in particular, what kinds of portfolios they choose to hold and how they trade over time. The goal here is less controversial than in the previous three sections: it is simply to explain the actions of certain investors, and not necessarily to claim that these actions also affect prices. Two factors make this type of research increasingly important. First, now that the costs of entering the stock market have fallen, more and more individuals are investing in equities. Second, the worldwide trend toward defined contribution retirement savings plans, and the possibility of individual accounts in social security systems mean that individuals are more responsible for their own financial well-being in retirement. It is therefore natural to ask how well they are handling these tasks.

We now describe some of the evidence on the actions of investors and the behavioral ideas that have been used to explain it.

7.1. Insufficient diversification

A large body of evidence suggests that investors diversify their portfolio holdings much less than is recommended by normative models of portfolio choice.

First, investors exhibit a pronounced "home bias". French and Poterba (1991) report that investors in the USA, Japan and the UK allocate 94%, 98%, and 82% of their overall equity investment, respectively, to *domestic* equities. It has not been easy to

³³ Similar results from univariate regressions can also be found in earlier work by Vlijh (1994).

explain this fact on rational grounds [Lewis (1999)]. Indeed, normative portfolio choice models that take human capital into account typically advise investors to *short* their national stock market, because of its high correlation with their human capital [Baxter and Jermann (1997)].

Some studies have found an analog to home bias *within* countries. Using an especially detailed data set from Finland, Grinblatt and Keloharju (2001) find that investors in that country are much more likely to hold and trade stocks of Finnish firms which are located close to them geographically, which use their native tongue in company reports, and whose chief executive shares their cultural background. Huberman (2001) studies the geographic distribution of shareholders of U.S. Regional Bell Operating Companies (RBOCs) and finds that investors are much more likely to hold shares in their local RBOC than in out-of-state RBOCs. Finally, studies of allocation decisions in 401(k) plans find a strong bias towards holding own company stock: over 30% of defined contribution plan assets in large U.S. companies are invested in employer stock, much of this representing voluntary contributions by employees [Benartzi (2001)].

In Section 3, we discussed evidence showing that people dislike ambiguous situations, where they feel unable to specify a gamble's probability distribution. Often, these are situations where they feel that they have little competence in evaluating a certain gamble. On the other hand, people show an excessive liking for familiar situations, where they feel they are in a better position than others to evaluate a gamble.

Ambiguity and familiarity offer a simple way of understanding the different examples of insufficient diversification. Investors may find their national stock markets more familiar – or less ambiguous – than foreign stock indices; they may find firms situated close to them geographically more familiar than those located further away; and they may find their employer's stock more familiar than other stocks.³⁴ Since familiar assets are attractive, people invest heavily in those, and invest little or nothing at all in ambiguous assets. Their portfolios therefore appear undiversified relative to the predictions of standard models that ignore the investor's degree of confidence in the probability distribution of a gamble.

Not all evidence of home bias should be interpreted as a preference for the familiar. Coval and Moskowitz (1999) show that U.S. mutual fund managers tend to hold stocks whose company headquarters are located close to their funds' headquarters. However, Coval and Moskowitz's (2001) finding that these local holdings subsequently perform well suggests that an information story is at work here, not a preference for the familiar. It is simply less costly to research local firms and so fund managers do indeed focus on those firms, picking out the stocks with higher expected returns. There is no obvious information-based explanation for the results of French and Poterba (1991), Huberman

³⁴ Particularly relevant to this last point is survey data showing that people consider their own company stock less risky than a diversified index [Driscoll et al. (1995)].

(2001) or Benartzi (2001), while Grinblatt and Keloharju (2001) argue against such an interpretation of their findings.

7.2. *Naive diversification*

Benartzi and Thaler (2001) find that when people *do* diversify, they do so in a naive fashion. In particular, they provide evidence that in 401(k) plans, many people seem to use strategies as simple as allocating $1/n$ of their savings to each of the n available investment options, whatever those options are. Some evidence that people think in this way comes from the laboratory. Benartzi and Thaler ask subjects to make an allocation decision in each of the following three conditions: first, between a stock fund and a bond fund; next, between a stock fund and a balanced fund, which invests 50% in stocks and 50% in bonds; and finally, between a bond fund and a balanced fund. They find that in all three cases, a 50:50 split across the two funds is a popular choice, although of course this leads to very different effective choices between stocks and bonds: the average allocation to stocks in the three conditions was 54%, 73% and 35%, respectively.

The $1/n$ diversification heuristic and other similar naive diversification strategies predict that in 401(k) plans which offer predominantly stock funds, investors will allocate more to stocks. Benartzi and Thaler test this in a sample of 170 large retirement savings plans. They divide the plans into three groups based on the fraction of funds – low, medium, or high – they offer that are stock funds. The allocation to stocks increases across the three groups, from 49% to 60% to 64%, confirming the initial prediction.

7.3. *Excessive trading*

One of the clearest predictions of rational models of investing is that there should be very little trading. In a world where rationality is common knowledge, I am reluctant to buy if you are ready to sell. In contrast to this prediction, the volume of trading on the world's stock exchanges is very high. Furthermore, studies of individuals and institutions suggest that both groups trade more than can be justified on rational grounds.

Barber and Odean (2000) examine the trading activity from 1991 to 1996 in a large sample of accounts at a national discount brokerage firm. They find that after taking trading costs into account, the average return of investors in their sample is well below the return of standard benchmarks. Put simply, these investors would do a lot better if they traded less. The underperformance in this sample is largely due to transaction costs. However, there is also some evidence of poor security selection: in a similar data set covering the 1987 to 1993 time period, Odean (1999) finds that the average gross return of stocks that investors buy, over the year after they buy them, is lower than the average gross return of stocks that they sell, over the year after they sell them.

The most prominent behavioral explanation of such excessive trading is overconfidence: people believe that they have information strong enough to justify a trade, whereas in fact the information is too weak to warrant any action. This hypothesis immediately predicts that people who are more overconfident will trade more and, because of transaction costs, earn lower returns. Consistent with this, Barber and Odean (2000) show that the investors in their sample who trade the most earn by far the lowest average returns. Building on evidence that men are more overconfident than women, and using the same data as in their earlier study, Barber and Odean (2001) predict and confirm that men trade more and earn lower returns on average.

Working with the same data again, Barber and Odean (2002a) study the subsample of individual investors who switch from phone-based to online trading. They argue that for a number of reasons, the switch should be accompanied by an increase in overconfidence. First, better access to information and a greater degree of control – both features of an online trading environment – have been shown to increase overconfidence. Moreover, the investors who switch have often earned high returns prior to switching, which may only increase their overconfidence further. If this is indeed the case, they should trade more actively after switching and perform worse. Barber and Odean confirm these predictions.

7.4. The selling decision

Several studies find that investors are reluctant to sell assets trading at a loss relative to the price at which they were purchased, a phenomenon labelled the “disposition effect” by Shefrin and Statman (1985). Working with the same discount brokerage data used in the Odean (1999) study from above, Odean (1998) finds that the individual investors in his sample are more likely to sell stocks which have gone up in value relative to their purchase price, rather than stocks which have gone down.

It is hard to explain this behavior on rational grounds. Tax considerations point to the selling of losers, not winners.³⁵ Nor can one argue that investors rationally sell the winners because of information that their future performance will be poor. Odean reports that the average performance of stocks that people sell is better than that of stocks they hold on to.

Two behavioral explanations of these findings have been suggested. First, investors may have an irrational belief in mean-reversion. A second possibility relies on prospect theory and narrow framing. We have used these ingredients before, but this time it is not loss aversion that is central, but rather the concavity (convexity) of the value function in the region of gains (losses).

To see the argument, suppose that a stock that was originally bought at \$50 now sells for \$55. Should the investor sell it at this point? Suppose that the gains and losses of

³⁵ Odean (1998) does find that in December, investors prefer to sell past losers rather than past winners, but overall, this effect is swamped by a strong preference for selling past winners in the remaining 11 months.

prospect theory refer to the sale price minus the purchase price. In that case, the utility from selling the stock now is $v(5)$. Alternatively, the investor can wait another period, whereupon we suppose that the stock could go to \$50 or \$60 with equal probability; in other words, we abstract from belief-based trading motives by saying that the investor expects the stock price to stay flat. The expected value of waiting and selling next period is then $\frac{1}{2}v(0) + \frac{1}{2}v(10)$. Since the value function v is concave in the region of gains, the investor sells now. In a different scenario, the stock may currently be trading at \$45. This time, the comparison is between $v(-5)$ and $\frac{1}{2}v(-10) + \frac{1}{2}v(0)$, assuming a second period distribution of \$40 and \$50 with equal probability. Convexity of v pushes the investor to wait. Intuitively, by not selling, he is gambling that the stock will eventually break even, saving him from having to experience a painful loss.

The disposition effect is not confined to individual stocks. In an innovative study, Genesove and Mayer (2001) find evidence of a reluctance to sell at a loss in the housing market. They show that sellers whose expected selling price is below their original purchase price, set an asking price that exceeds the asking price of other sellers with comparable houses. Moreover, this is not simply wishful thinking on the sellers' part that is later corrected by the market: sellers facing a possible loss do actually transact at considerably higher prices than other sellers.

Coval and Shumway (2000) study the behavior of professional traders in the Treasury Bond futures pit at the CBOT. If the gains and losses of prospect theory are taken to be daily profits and losses, the curvature of the value function implies that traders with profits (losses) by the middle of the trading day will take less (more) risk in their afternoon trading. This prediction is borne out in the data.

Grinblatt and Han (2001) argue that the investor behavior inherent in the disposition effect may be behind a puzzling feature of the cross-section of average returns, namely momentum in stock returns. Due to the concavity of the value function in the region of gains, investors will be keen to sell a stock which has earned them capital gains on paper. The selling pressure that results may initially depress the stock price, generating higher returns later. On the other hand, if the holders of a stock are facing capital losses, convexity in the region of losses means that they will only sell if offered a price premium; the price is therefore initially inflated, generating lower returns later. Grinblatt and Han provide supportive evidence for their story by regressing, in the cross-section, a stock's return on its past 12-month return as well as on a measure of the capital gain or loss faced by its holders. This last variable is computed as the current stock price minus investors' average cost basis, itself inferred from past volume. They find that the capital gain or loss variable steals a substantial amount of explanatory power from the past return.

7.5. The buying decision

Odean (1999) presents useful information about the stocks the individual investors in his sample choose to buy. Unlike "sells", which are mainly prior winners, "buys" are evenly split between prior winners and losers. Conditioning on the stock being a prior

winner (loser) though, the stock is a big prior winner (loser). In other words, a good deal of the action is in the extremes.

Odean argues that the results for stock purchases are in part due to an attention effect. When buying a stock, people do not tend to systematically sift through the thousands of listed shares until they find a good “buy”. They typically buy a stock that has caught their attention and perhaps the best attention draw is extreme past performance, whether good or bad.

Among individual investors, attention is less likely to matter for stock sales because of a fundamental way in which the selling decision differs from the buying decision. Due to short-sale constraints, when individuals are looking for a stock to sell, they limit their search to those stocks that they currently own. When buying stocks, though, people have a much wider range of possibilities to choose from, and factors related to attention may enter the decision more.

Using the same discount brokerage data as in their earlier papers, Barber and Odean (2002b) test the idea that for individual investors, buying decisions are more driven by attention than are selling decisions. On any particular day, they create portfolios of “attention-getting” stocks using a number of different criteria: stocks with abnormally high trading volume, stocks with abnormally high or low returns, and stocks with news announcements. They find that the individual investors in their sample are more likely, on the following day, to be purchasers of these high-attention stocks than sellers.

8. Application: Corporate finance

8.1. Security issuance, capital structure and investment

An important strand of research in behavioral finance asks whether irrational investors such as those discussed in earlier sections affect the financing and investment decisions of firms.

We first address this question theoretically, and ask how a rational manager interested in maximizing true firm value – in other words, the stock price that will prevail once any mispricing has worked its way out of valuations – should act in the face of irrational investors. Stein (1996) provides a useful framework for thinking about this, as well as about other issues that arise in this section. He shows that when a firm’s stock price is too high, the rational manager should issue more shares so as to take advantage of investor exuberance. Conversely, when the price is too low, the manager should repurchase shares. We refer to this model of security issuance as the “market timing” view.

What evidence there is to date on security issuance appears remarkably consistent with this framework. First, at the aggregate level, the share of new equity issues among total new issues – the “equity share” – is higher when the overall stock market is more highly valued. In fact, Baker and Wurgler (2000) show that the equity share is a reliable predictor of future stock returns: a high share predicts low, and sometimes negative,

stock returns. This is consistent with managers timing the market, issuing more equity at its peaks, just before it sinks back to more realistic valuation levels.

At the individual firm level, a number of papers have shown that the book-to-market ratio of a firm is a good cross-sectional predictor of new equity issuance [see Korajczyk, Lucas and McDonald (1991), Jung, Kim and Stulz (1996), Loughran, Ritter and Rydqvist (1994), Pagano, Panetta and Zingales (1998), Baker and Wurgler (2002a)]. Firms with high valuations issue more equity while those with low valuations repurchase their shares. Moreover, long-term stock returns after an IPO or SEO are low [Loughran and Ritter (1995)], while long-term returns after the announcement of a repurchase are high [Ikenberry, Lakonishok and Vermaelen (1995)]. Once again, this evidence is consistent with managers timing the market in their own securities.

More support for the market-timing view comes from survey evidence. Graham and Harvey (2001) report that 67% of surveyed CFOs said that “the amount by which our stock is undervalued or overvalued” was an important consideration when issuing common stock.

The success of the market-timing framework in predicting patterns of equity issuance offers the hope that it might also be the basis of a successful theory of capital structure. After all, a firm’s capital structure simply represents its cumulative financing decisions over time. Consider, for example, two firms which are similar in terms of characteristics like firm size, profitability, fraction of tangible assets, and current market-to-book ratio, which have traditionally been thought to affect capital structure. Suppose, however, that in the past, the market-to-book ratio of firm A has reached much higher levels than that of firm B. Since, under the market timing theory, managers of firm A may have issued more shares at that time to take advantage of possible overvaluation, firm A may have more equity in its capital structure today.

In an intriguing recent paper, Baker and Wurgler (2002a) confirm this prediction. They show that all else equal, a firm’s weighted-average historical market-to-book ratio, where more weight is placed on years in which the firm made an issuance of some kind, whether debt or equity, is a good cross-sectional predictor of the fraction of equity in the firm’s capital structure today.

There is some evidence, then, that irrational investor sentiment affects financing decisions. We now turn to the more critical question of whether this sentiment affects actual investment decisions. Once again, we consider the benchmark case in Stein’s (1996) model, in which the manager is both rational and interested in maximizing the firm’s true value.

Suppose that a firm’s stock price is too high. As discussed above, the manager should issue more equity at this point. More subtly, though, Stein shows that he should *not* channel the fresh capital into any actual new investment, but instead keep it in cash or in another fairly priced capital market security. While investors’ exuberance means that, in *their* view, the firm has many positive net present value (NPV) projects it could undertake, the rational manager knows that these projects are not, in fact, positive NPV and that in the interest of true firm value, they should be avoided. Conversely, if the manager thinks that his firm’s stock price is irrationally low, he should repurchase

shares at the advantageously low price but not scale back actual investment. In short, irrational investors may affect the timing of security issuance, but they should not affect the firm's investment plans.

Once we move beyond this simple benchmark case, though, there emerge several channels through which sentiment might affect investment after all. First, the above argument properly applies only to *non-equity dependent* firms; in other words, to firms which because of their ample internal funds and borrowing capacity do not need the equity markets to finance their marginal investments.

For equity-dependent firms, however, investor sentiment and, in particular, excessive investor pessimism, may distort investment: when investors are excessively pessimistic, such firms may have to forgo attractive investment opportunities because it is too costly to finance them with undervalued equity. This thinking leads to a cross-sectional prediction, namely that the investment of equity-dependent firms should be more sensitive to gyrations in stock price than the investment of non-equity dependent firms.

Other than this equity-dependence mechanism, there are other channels through which investor sentiment might distort investment. Consider the case where investors are excessively optimistic about a firm's prospects. Even if a manager is in principle interested in maximizing true value, he faces the danger that if he refuses to undertake projects investors perceive as profitable, they may depress stock prices, exposing him to the risk of a takeover, or more simply, try to have him fired.³⁶

Even if the manager is rational, this does not mean he will choose to maximize the firm's true value. The agency literature has argued that some managers may maximize other objectives – the size of their firm, say – as a way of enhancing their prestige. This suggests another channel for investment distortion: managers might use investor exuberance as a cover for doing negative NPV “empire building” projects.

Finally, investor sentiment can also affect investment if managers put some weight on investors' opinions, perhaps because they think investors know something they don't. Managers may then mistake excessive optimism for well-founded optimism and get drawn into making negative NPV investments.

An important goal of empirical research, then, is to try to understand whether sentiment does affect investment, and if so, through which channel. Early studies produced little evidence of investment distortion. In aggregate data, Blanchard, Rhee and Summers (1993) find that movements in price apparently unrelated to movements in fundamentals have only weak forecasting power for future investment: the effects are marginally statistically significant and weak in economic terms. To pick out two particular historical episodes: the rise in stock prices through the 1920s did not lead to

³⁶ Shleifer and Vishny (2004) argue that in a situation such as this, where the manager feels forced to undertake some kind of investment, the best investment of all may be an acquisition of a less overvalued firm, in other words, one more likely to retain its value in the long run. This observation leads to a parsimonious theory of takeover waves, which predicts, among other things, an increase in stock-financed acquisitions at times of high dispersion in valuations.

a commensurate rise in investment, nor did the crash of 1987 slow investment down appreciably. Morck, Shleifer and Vishny (1990) reach similar conclusions using firm level data, as do Baker and Wurgler (2002a): in their work on capital structure, they show that not only do firms with higher market-to-book ratios in their past have more equity in their capital structure today, but also that the equity funds raised are typically used to increase cash balances and *not* to finance new investment.

More recently though, Polk and Sapienza (2001) report stronger evidence of investment distortion. They identify overvalued firms as firms with high accruals, defined as earnings minus actual cash flow, and as firms with high net issuance of equity. Firms with high accruals may become overvalued if investors fail to understand that earnings are overstating actual cash flows, and Chan et al. (2001) confirm that such firms indeed earn low returns. Overvalued firms may also be identified through their opportunistic issuance of equity, and we have already discussed the evidence that such firms earn low long-run returns. Controlling for actual investment opportunities as accurately as possible, Polk and Sapienza find that the firms they identify as overvalued appear to invest more than other firms, suggesting that sentiment does influence investment.

Further evidence of distortion comes from Baker, Stein and Wurgler's (2003) test of the cross-sectional prediction that equity-dependent firms will be more sensitive to stock price gyrations than will non-equity dependent firms. They identify equity-dependent firms on the basis of their low cash balances, among other measures, and find that these firms have an investment sensitivity to stock prices about three times as high as that of non-equity dependent firms. This study therefore provides initial evidence that for some firms at least, sentiment may distort investment, and that it does so through the equity-dependence channel.

8.2. Dividends

A major open question in corporate finance asks why firms pay dividends. Historically, dividends have been taxed at a higher rate than capital gains. This means that stockholders who pay taxes would always prefer that the firm repurchase shares rather than pay a dividend. Since the tax exempt shareholders would be indifferent between the dividend payment and the share repurchase, the share repurchase is a Pareto improving action. Why then, do investors seem perfectly happy to accept a substantial part of their return in the form of dividends? Or, using behavioral language, why do firms choose to frame part of their return as an explicit payment to stockholders, and in so doing, apparently make some of their shareholders worse off?

Shefrin and Statman (1984) propose a number of behavioral explanations for why investors exhibit a preference for dividends. Their first idea relies on the notion of self-control. Many people exhibit self-control problems. On the one hand, we want to deny ourselves an indulgence, but on the other hand, we quickly give in to temptation: today, we tell ourselves that tomorrow we will not overeat, and yet, when tomorrow arrives, we again eat too much. To deal with self-control problems, people often set rules, such

as “bank the wife’s salary, and only spend from the husband’s paycheck”. Another very natural rule people might create to prevent themselves from overconsuming their wealth is “only consume the dividend, but don’t touch the portfolio capital”. In other words, people may like dividends because dividends help them surmount self-control problems through the creation of simple rules.

A second rationale for dividends is based on mental accounting: by designating an explicit dividend payment, firms make it easier for investors to segregate gains from losses and hence to increase their utility. To see this, consider the following example. Over the course of a year, the value of a firm has increased by \$10 per share. The firm could choose *not* to pay a dividend and return this increase in value to investors as a \$10 capital gain. Alternatively, it could pay a \$2 dividend, leaving an \$8 capital gain. In the language of prospect theory, investors will code the first option as $v(10)$. They may also code the second option as $v(10)$, but the explicit segregation performed by the firm may encourage them to code it as $v(2) + v(8)$. This will, of course, result in a higher perceived utility, due to the concavity of v in the domain of gains.

This manipulation is equally useful in the case of losses. A firm whose value has declined by \$10 per share over the year can offer investors a \$10 capital loss or a \$12 capital loss combined with a \$2 dividend gain. While the first option will be coded as $v(-10)$, the second is more likely to be coded as $v(2) + v(-12)$, again resulting in a higher perceived utility, this time because of the convexity of v in the domain of losses.

The utility enhancing trick in these examples depends on investors segregating the overall gain or loss into different components. The key insight of Shefrin and Statman is that by paying dividends, firms make it easier for investors to perform this segregation.

Finally, Shefrin and Statman argue that by paying dividends, firms help investors avoid regret. Regret is a frustration that people feel when they imagine having taken an action that would have led to a more desirable outcome. It is stronger for errors of commission – cases where people suffer because of an action they took – than for errors of omission – where people suffer because of an action they *failed* to take.

Consider a company which does not pay a dividend. In order to finance consumption, an investor has to sell stock. If the stock subsequently goes up in value, the investor feels substantial regret because the error is one of commission: he can readily imagine how not selling the stock would have left him better off. If the firm had paid a dividend and the investor was able to finance his consumption out of it, a rise in the stock price would not have caused so much regret. This time, the error would have been one of omission: to be better off, the investor would have had to reinvest the dividend.

Shefrin and Statman try to explain why firms pay dividends at all. Another question asks how dividend paying firms decide on the size of their dividend. The classic paper on this subject is Lintner (1956). His treatment is based on extensive interviews with executives of large American companies in which he asked the respondent, often the CFO, how the firm set dividend policy. Based on these interviews Lintner proposed what we would now call a behavioral model. In his model, firms first establish a target

dividend payout rate based on notions of fairness, in other words, on what portion of the earnings it is fair to return to shareholders. Then, as earnings increase and the dividend payout ratio falls below the target level, firms increase dividends only when they are confident that they will not have to reduce them in the future.

There are several behavioral aspects to this model. First, the firm is not setting the dividend to maximize firm value or shareholder after-tax wealth. Second, perceptions of fairness are used to set the target payout rate. Third, the asymmetry between an increase in dividends and a decrease is explicitly considered. Although fewer firms now decide to start paying dividends, for those that do Lintner's model appears to be valid to this day [Benartzi, Michaely and Thaler (1997), Fama and French (2001)].

Baker and Wurgler (2002b) argue that changes in dividend policy may also reflect changing investor sentiment about dividend-paying firms relative to their sentiment about non-paying firms. They argue that for some investors, dividend-paying firms and non-paying firms represent salient categories and that these investors exhibit changing sentiment about the categories. For instance, when investors become more risk averse, they may prefer dividend-paying stocks because of a confused notion that these firms are less risky (the well-known "bird in the hand" fallacy). If managers are interested in maximizing short-run value, perhaps because it is linked to their compensation, they may be tempted to change their dividend policy in the direction favored by investors.

Baker and Wurgler find some supportive evidence for their theory. They measure relative investor sentiment about dividend-paying firms as the log market-to-book ratio of paying firms minus the log market-to-book ratio of non-paying firms, and find that in the time series, a high value of this measure one year predicts that in the following year, a higher fraction of non-paying firms initiate a dividend and a larger fraction of newly-listed firms choose to pay one. Similar results obtain for other measures of sentiment about dividend-paying firms.

8.3. Models of managerial irrationality

The theories we have discussed so far interpret the data as reflecting actions taken by rational managers in response to irrationality on the part of investors. Other papers have argued that some aspects of managerial behavior are the result of irrationality on the part of managers themselves.

Much of Section 2 was devoted to thinking about whether rational agents might be able to correct dislocations caused by irrational traders. Analogously, before we consider models of irrational managers, we should ask to what extent rational agents can undo their effects.

On reflection, it doesn't seem any easier to deal with irrational managers than irrational investors. It is true that many firms have mechanisms in place designed to solve agency problems and to keep the manager's mind focused on maximizing firm value: giving him stock options for example, or saddling him with debt. The problem is that these mechanisms are unlikely to have much of an effect on irrational managers. These managers *think* that they are maximizing firm value, even if in reality, they are

not. Since they think that they are already doing the right thing, stock options or debt are unlikely to change their behavior.

In the best known paper on managerial irrationality, Roll (1986) argues that much of the evidence on takeover activity is consistent with an economy in which there are *no* overall gains to takeovers, but in which managers are overconfident, a theory he terms the “hubris hypothesis”. When managers think about taking over another firm, they conduct a valuation analysis of that firm, taking synergies into account. If managers are overconfident about the accuracy of their analysis, they will be too quick to launch a bid when their valuation exceeds the market price of the target. Just as overconfidence among individual investors may lead to excessive trading, so overconfidence among managers may lead to excessive takeover activity.

The main predictions of the hubris hypothesis are that there will be a large amount of takeover activity, but that the total combined gain to bidder and target will be zero; and that on the announcement of a bid, the price of the target will rise and the value of the bidder will fall by a similar amount. Roll examines the available evidence and concludes that it is impossible to reject any of these predictions.

Heaton (2002) analyses the consequences of managerial optimism whereby managers overestimate the probability that the future performance of their firm will be good. He shows that it can explain pecking order rules for capital structure: since managers are optimistic relative to the capital markets, they believe their equity is undervalued, and are therefore reluctant to issue it unless they have exhausted internally generated funds or the debt market. Managerial optimism can also explain the puzzlingly high correlation of investment and cash flow: when cash flow is low, managers’ reluctance to use external markets for financing means that they forgo an unusually large number of projects, lowering investment at the same time.

Malmendier and Tate (2001) test Heaton’s model by investigating whether firms with excessively optimistic CEOs display a greater sensitivity of investment to cash flow. They detect excessive optimism among CEOs by examining at what point they exercise their stock options: CEOs who hold on to their options longer than recommended by normative models of optimal exercise are deemed to have an overly optimistic forecast of their stock’s future price. Malmendier and Tate find that the investment of these CEOs’ firms is indeed more sensitive to cash flow than the investment of other firms.³⁷

³⁷ Another paper which can be included in the managerial irrationality category is Loughran and Ritter’s (2002) explanation for why managers issuing shares appear to leave significant amounts of money “on the table”, as evidenced by the high average return of IPOs on their first day of trading. The authors note that the IPOs with good first day performance are often those IPOs in which the price has risen far above its filing range, giving the managers a sizeable wealth gain. One explanation is therefore that since managers are already enjoying a major windfall, they do not care too much about the fact that they could have been even wealthier.

9. Conclusion

Behavioral finance is a young field, with its formal beginnings in the 1980s. Much of the research we have discussed was completed in the past five years. Where do we stand? Substantial progress has been made on numerous fronts.

Empirical investigation of apparently anomalous facts. When De Bondt and Thaler's (1985) paper was published, many scholars thought that the best explanation for their findings was a programming error. Since then their results have been replicated numerous times by authors both sympathetic to their view and by those with alternative views. At this stage, we think that most of the empirical facts are agreed upon by most of the profession, although the interpretation of those facts is still in dispute. This is progress. If we all agree that the planets do orbit the sun, we can focus on understanding why.

Limits to arbitrage. Twenty years ago, many financial economists thought that the Efficient Markets Hypothesis had to be true because of the forces of arbitrage. We now understand that this was a naive view, and that the limits to arbitrage can permit substantial mispricing. It is now also understood by most that the absence of a profitable investment strategy does not imply the absence of mispricing. Prices can be very wrong without creating profit opportunities.

Understanding bounded rationality. Thanks largely to the work of cognitive psychologists such as Daniel Kahneman and Amos Tversky, we now have a long list of robust empirical findings that catalogue some of the ways in which actual humans form expectations and make choices. There has also been progress in writing down formal models of these processes, with prospect theory being the most notable. Economists once thought that behavior was either rational or impossible to formalize. We now know that models of bounded rationality are both possible and also much more accurate descriptions of behavior than purely rational models.

Behavioral finance theory building. In the past few years there has been a burst of theoretical work modelling financial markets with less than fully rational agents. These papers relax the assumption of individual rationality either through the belief formation process or through the decision-making process. Like the work of psychologists discussed above, these papers are important existence proofs, showing that it is possible to think coherently about asset pricing while incorporating salient aspects of human behavior.

Investor behavior. We have now begun the important job of trying to document and understand how investors, both amateurs and professionals, make their portfolio choices. Until recently such research was notably absent from the repertoire of financial economists.

This is a lot of accomplishment in a short period of time, but we are still much closer to the beginning of the research agenda than we are to the end. We know enough about the perils of forecasting to realize that most of the future progress of the field is unpredictable. Still, we cannot resist venturing a few observations on what may be coming next.

First, much of the work we have summarized is narrow. Models typically capture something about investors' beliefs, or their preferences, or the limits to arbitrage, but not all three. This comment applies to most research in economics, and is a natural implication of the fact that researchers are boundedly rational too. Still, as progress is made, we expect theorists to begin to incorporate more than one strand into their models.

An example can, perhaps, illustrate the point. The empirical literature repeatedly finds that the asset pricing anomalies are more pronounced in small and mid-cap stocks than in the large cap sector. It seems likely that this finding reflects limits to arbitrage: the costs of trading smaller stocks are higher, keeping many potential arbitrageurs uninterested. While this observation may be an obvious one, it has not found its way into formal models. We expect investigation of the interplay between limits to arbitrage and cognitive biases to be an important research area in the coming years.

Second, there are obviously competing behavioral explanations for some of the empirical facts. Some critics view this as a weakness of the field. It is sometimes said that the long list of cognitive biases summarized in Section 3 offer behavioral modelers so many degrees of freedom that anything can be explained. We concede that there are numerous degrees of freedom, but note that rational modelers have just as many options to choose from. As Arrow (1986) has forcefully argued, rationality *per se* does not yield many predictions. The predictions in rational models often come from auxiliary assumptions.

There is really only one scientific way to compare alternative theories, behavioral or rational, and that is with empirical tests. One kind of test looks for novel predictions the theory makes. For example, Lee, Shleifer and Thaler (1991) test their model's prediction that small firm returns will be correlated with closed-end fund discounts, while Hong, Lim and Stein (2000) test the implication of the Hong and Stein (1999) model that momentum will be stronger among stocks with thinner analyst coverage.

Another sort of test is to look for evidence that agents actually behave the way a model claims they do. The Odean (1998) and Genesove and Mayer (2001) investigations of the disposition effect using actual market behavior fall into this category. Bloomfield and Hales (2002) offers an experimental test of the behavior theorized by Barberis, Shleifer and Vishny (1998). Of course, such tests are never airtight, but we should be skeptical of theories based on behavior that is undocumented empirically. Since behavioral theories claim to be grounded in realistic assumptions about behavior, we hope behavioral finance researchers will continue to give their

assumptions empirical scrutiny. We would urge the same upon authors of rational theories.³⁸

We have two predictions about the outcome of direct tests of the assumptions of economic models. First, we will find that most of our current theories, both rational and behavioral, are wrong. Second, substantially better theories will emerge.

Appendix A

We show that for the economy laid out in Equations (3–6), there is an equilibrium in which the risk-free rate is constant and given by

$$R_f = \frac{1}{\rho} \exp \left[\gamma g_C - \frac{1}{2} \gamma^2 \sigma_C^2 \right], \quad (18)$$

and in which the price–dividend ratio is a constant f , and satisfies

$$1 = \rho \frac{1+f}{f} \exp \left[g_D - \gamma g_C + \frac{1}{2} (\sigma_D^2 + \gamma^2 \sigma_C^2 - 2\gamma \sigma_C \sigma_D \omega) \right]. \quad (19)$$

In this equilibrium, returns are therefore given by

$$R_{t+1} = \frac{D_{t+1} + P_{t+1}}{P_t} = \frac{1 + P_{t+1}/D_{t+1}}{P_t/D_t} \cdot \frac{D_{t+1}}{D_t} = \frac{1+f}{f} \exp [g_D + \sigma_D \varepsilon_{t+1}]. \quad (20)$$

To see this, start from the Euler equations of optimality, obtained through the usual perturbation arguments,

$$1 = \rho R_f E_t \left[\left(\frac{C_{t+1}}{C_t} \right)^{-\gamma} \right], \quad (21)$$

$$1 = \rho E_t \left[R_{t+1} \left(\frac{C_{t+1}}{C_t} \right)^{-\gamma} \right]. \quad (22)$$

Computing the expectation in Equation (21) gives Equation (18). We conjecture that in this economy, there is an equilibrium in which the price–dividend ratio is a constant f , so that returns are given by Equation (20). Substituting this into

³⁸ Directly testing the validity of a model's assumptions is not common practice in economics, perhaps because of Milton Friedman's influential argument that one should evaluate theories based on the validity of their predictions rather than the validity of their assumptions. Whether or not this is sound scientific practice, we note that much of the debate over the past 20 years has occurred precisely because the evidence has not been consistent with the theories, so it may be a good time to start worrying about the assumptions. If a theorist wants to claim that fact X can be explained by behavior Y , it seems prudent to check whether people actually do Y .

Equation (22) and computing the expectation gives Equation (19), as required. For given parameter values, the quantitative implications for *P/D* ratios and returns are now easily computed.

References

- Abreu, D., and M. Brunnermeier (2002), "Synchronization risk and delayed arbitrage", *Journal of Financial Economics* 66:341–360.
- Alpert, M., and H. Raiffa (1982), "A progress report on the training of probability assessors", in: D. Kahneman, P. Slovic and A. Tversky, eds., *Judgment Under Uncertainty: Heuristics and Biases* (Cambridge University Press, Cambridge) pp. 294–305.
- Anderson, E., L. Hansen and T. Sargent (1998), "Risk and robustness in equilibrium", Working Paper (University of Chicago).
- Arrow, K. (1986), "Rationality of self and others", in: R. Hogarth and M. Reder, eds., *Rational Choice* (University of Chicago Press, Chicago) pp. 201–215.
- Baker, M., and J. Wurgler (2000), "The equity share in new issues and aggregate stock returns", *Journal of Finance* 55:2219–2257.
- Baker, M., and J. Wurgler (2002a), "Market timing and capital structure", *Journal of Finance* 57:1–32.
- Baker, M., and J. Wurgler (2002b), "A catering theory of dividends", Working Paper (Harvard University).
- Baker, M., J. Stein and J. Wurgler (2003), "When does the market matter? Stock prices and the investment of equity dependent firms", *Quarterly Journal of Economics*, forthcoming.
- Ball, R. (1978), "Anomalies in relations between securities' yields and yield surrogates", *Journal of Financial Economics* 6:103–126.
- Banz, R. (1981), "The relation between return and market value of common stocks", *Journal of Financial Economics* 9:3–18.
- Barber, B., and J. Lyon (1997), "Detecting long-run abnormal stock returns: the empirical power and specification of test statistics", *Journal of Financial Economics* 43:341–372.
- Barber, B., and T. Odean (2000), "Trading is hazardous to your wealth: the common stock performance of individual investors", *Journal of Finance* 55:773–806.
- Barber, B., and T. Odean (2001), "Boys will be boys: gender, overconfidence, and common stock investment", *Quarterly Journal of Economics* 141:261–292.
- Barber, B., and T. Odean (2002a), "Online investors: do the slow die first?", *Review of Financial Studies* 15:455–487.
- Barber, B., and T. Odean (2002b), "All that glitters: the effect of attention and news on the buying behavior of individual and institutional investors", Working Paper (University of California, Berkeley, CA).
- Barberis, N., and M. Huang (2001), "Mental accounting, loss aversion and individual stock returns", *Journal of Finance* 56:1247–1292.
- Barberis, N., and A. Shleifer (2003), "Style investing", *Journal of Financial Economics* 68:161–199.
- Barberis, N., A. Shleifer and R. Vishny (1998), "A model of investor sentiment", *Journal of Financial Economics* 49:307–345.
- Barberis, N., M. Huang and T. Santos (2001), "Prospect theory and asset prices", *Quarterly Journal of Economics* 116:1–53.
- Barberis, N., A. Shleifer and J. Wurgler (2001), "Comovement", Working Paper (University of Chicago).
- Barsky, R., and B. De Long (1993), "Why does the stock market fluctuate?", *Quarterly Journal of Economics* 107:291–311.
- Basu, S. (1983), "The relationship between earnings yield, market value and return for NYSE common stocks: further evidence", *Journal of Financial Economics* 12:129–156.

- Baxter, M., and U. Jermann (1997), "The international diversification puzzle is worse than you think", *American Economic Review* 87:170–180.
- Bell, D. (1982), "Regret in decision making under uncertainty", *Operations Research* 30:961–981.
- Benartzi, S. (2001), "Excessive extrapolation and the allocation of 401(k) accounts to company stock", *Journal of Finance* 56:1747–1764.
- Benartzi, S., and R. Thaler (1995), "Myopic loss aversion and the equity premium puzzle", *Quarterly Journal of Economics* 110:75–92.
- Benartzi, S., and R. Thaler (2001), "Naïve diversification strategies in defined contribution savings plans", *American Economic Review* 91:79–98.
- Benartzi, S., R. Michaely and R. Thaler (1997), "Do changes in dividends signal the future or the past?", *Journal of Finance* 52:1007–1034.
- Berk, J. (1995), "A critique of size related anomalies", *Review of Financial Studies* 8:275–286.
- Bernard, V., and J. Thomas (1989), "Post-earnings announcement drift: delayed price response or risk premium?", *Journal of Accounting Research (Supplement)*, pp. 1–36.
- Blanchard, O., C. Rhee and L. Summers (1993), "The stock market, profit, and investment", *Quarterly Journal of Economics* 108:115–136.
- Bloomfield, R., and J. Hales (2002), "Predicting the next step of a random walk: experimental evidence of regime-shifting beliefs", *Journal of Financial Economics* 65:397–414.
- Bodurtha, J., D. Kim and C.M. Lee (1993), "Closed-end country funds and U.S. market sentiment", *Review of Financial Studies* 8:879–918.
- Brav, A. (2000), "Inference in long-horizon event studies", *Journal of Finance* 55:1979–2016.
- Brav, A., and P. Gompers (1997), "Myth or reality? The long-run underperformance of initial public offerings: evidence from venture and non-venture-backed companies", *Journal of Finance* 52:1791–1821.
- Brav, A., C. Geczy and P. Gompers (2000), "Is the abnormal return following equity issuances anomalous?", *Journal of Financial Economics* 56:209–249.
- Brennan, M., and Y. Xia (2001), "Stock return volatility and the equity premium", *Journal of Monetary Economics* 47:249–283.
- Brown, S., W. Goetzmann and S. Ross (1995), "Survival", *Journal of Finance* 50:853–873.
- Brunnermeier, M. (2001), *Asset Pricing under Asymmetric Information – Bubbles, Crashes, Technical Analysis, and Herding* (Oxford University Press).
- Buehler, R., D. Griffin and M. Ross (1994), "Exploring the planning fallacy: why people underestimate their task completion times", *Journal of Personality and Social Psychology* 67:366–381.
- Camerer, C. (1995), "Individual decision making", in: J. Kagel and A. Roth, eds., *Handbook of Experimental Economics* (Princeton University Press).
- Camerer, C., and R. Hogarth (1999), "The effects of financial incentives in experiments: a review and capital-labor production framework", *Journal of Risk and Uncertainty* 19:7–42.
- Camerer, C., and M. Weber (1992), "Recent developments in modeling preferences: uncertainty and ambiguity", *Journal of Risk and Uncertainty* 5:325–70.
- Campbell, J.Y. (1991), "A variance decomposition for stock returns", *Economic Journal* 101:157–179.
- Campbell, J.Y. (1999), "Asset prices, consumption and the business cycle", in: J. Taylor and M. Woodford, eds., *Handbook of Macroeconomics* (Elsevier, Amsterdam) pp. 1231–1303.
- Campbell, J.Y. (2000), "Asset pricing at the millennium", *Journal of Finance* 55:1515–1567.
- Campbell, J.Y., and J. Cochrane (1999), "By force of habit: a consumption-based explanation of aggregate stock market behavior", *Journal of Political Economy* 107:205–251.
- Campbell, J.Y., and R. Shiller (1988), "Stock prices, earnings and expected dividends", *Journal of Finance* 43:661–676.
- Cecchetti, S., P. Lam and N. Mark (2000), "Asset pricing with distorted beliefs: are equity returns too good to be true?", *American Economic Review* 90:787–805.
- Chan, K., L. Chan, N. Jegadeesh and J. Lakonishok (2001), "Earnings quality and stock returns", Working Paper (University of Illinois, Urbana, IL).

- Chan, L., N. Jegadeesh and J. Lakonishok (1996), "Momentum strategies", *Journal of Finance* 51: 1681–1713.
- Chen, J., H. Hong and J. Stein (2001), "Forecasting crashes: trading volume, past returns and conditional skewness in stock prices", *Journal of Financial Economics* 61:345–381.
- Chen, J., H. Hong and J. Stein (2002), "Breadth of ownership and stock returns", *Journal of Financial Economics* 66:171–205.
- Chew, S. (1983), "A generalization of the quasilinear mean with applications to the measurement of income inequality and decision theory resolving the allais paradox", *Econometrica* 51:1065–1092.
- Chew, S. (1989), "Axiomatic utility theories with the betweenness property", *Annals of Operations Research* 19:273–98.
- Chew, S., and K. MacCrimmon (1979), "Alpha-nu choice theory: an axiomatization of expected utility", Working Paper (University of British Columbia, Vancouver, BC).
- Chopra, N., J. Lakonishok and J. Ritter (1992), "Measuring abnormal performance: do stocks overreact?", *Journal of Financial Economics* 31:235–268.
- Coval, J., and T. Moskowitz (1999), "Home bias at home: local equity preference in domestic portfolios", *Journal of Finance* 54:2045–2073.
- Coval, J., and T. Moskowitz (2001), "The geography of investment: informed trading and asset prices", *Journal of Political Economy* 109:811–841.
- Coval, J., and T. Shumway (2000), "Do behavioral biases affect prices?", Working Paper (University of Michigan, Ann Arbor, MI).
- Daniel, K., and S. Titman (1997), "Evidence on the characteristics of cross-sectional variation in stock returns", *Journal of Finance* 52:1–33.
- Daniel, K., D. Hirshleifer and A. Subrahmanyam (1998), "Investor psychology and security market under- and overreactions", *Journal of Finance* 53:1839–1885.
- Daniel, K., D. Hirshleifer and A. Subrahmanyam (2001), "Overconfidence, arbitrage and equilibrium asset pricing", *Journal of Finance* 56:921–965.
- D'Avolio, G. (2002), "The market for borrowing stock", *Journal of Financial Economics* 66:271–306.
- De Bondt, W., and R. Thaler (1985), "Does the stock market overreact?", *Journal of Finance* 40:793–808.
- De Bondt, W., and R. Thaler (1987), "Further evidence on investor overreaction and stock market seasonality", *Journal of Finance* 42:557–581.
- De Long, J.B., A. Shleifer, L. Summers and R. Waldmann (1990a), "Noise trader risk in financial markets", *Journal of Political Economy* 98:703–738.
- De Long, J.B., A. Shleifer, L. Summers and R. Waldmann (1990b), "Positive feedback investment strategies and destabilizing rational speculation", *Journal of Finance* 45:375–395.
- Dekel, E. (1986), "An axiomatic characterization of preferences under uncertainty: weakening the independence axiom", *Journal of Economic Theory* 40:304–18.
- Diether, K., C. Malloy and A. Scherbina (2002), "Stock prices and differences of opinion: empirical evidence that stock prices reflect optimism", *Journal of Finance* 57:2113–2141.
- Dreman, D. (1977), *Psychology and the Stock Market: Investment Strategy Beyond Random Walk* (Warner Books, New York).
- Driscoll, K., J. Malcolm, M. Sirul and P. Slotter (1995), *Gallup Survey of Defined Contribution Plan Participants* (John Hancock Financial Services).
- Edwards, W. (1968), "Conservatism in human information processing", in: B. Kleinmuntz, ed., *Formal Representation of Human Judgment* (Wiley, New York) pp. 17–52.
- Ellsberg, D. (1961), "Risk, ambiguity, and the savage axioms", *Quarterly Journal of Economics* 75: 643–69.
- Epstein, L., and T. Wang (1994), "Intertemporal asset pricing under Knightian uncertainty", *Econometrica* 62:283–322.
- Fama, E. (1970), "Efficient capital markets: a review of theory and empirical work", *Journal of Finance* 25:383–417.

- Fama, E. (1998), "Market efficiency, long-term returns and behavioral finance", *Journal of Financial Economics* 49:283–307.
- Fama, E., and K. French (1988), "Dividend yields and expected stock returns", *Journal of Financial Economics* 22:3–25.
- Fama, E., and K. French (1992), "The cross-section of expected stock returns", *Journal of Finance* 47:427–465.
- Fama, E., and K. French (1993), "Common risk factors in the returns of bonds and stocks", *Journal of Financial Economics* 33:3–56.
- Fama, E., and K. French (1995), "Size and book-to-market factors in earnings and returns", *Journal of Finance* 50:131–155.
- Fama, E., and K. French (1996), "Multifactor explanations of asset pricing anomalies", *Journal of Finance* 51:55–84.
- Fama, E., and K. French (1998), "Value vs. growth: the international evidence", *Journal of Finance* 53:1975–1999.
- Fama, E., and K. French (2001), "Disappearing dividends: changing firm characteristics or lower propensity to pay?", *Journal of Financial Economics* 60:3–43.
- Fama, E., K. French and J. Davis (2000), "Characteristics, covariances and average returns 1929–1997", *Journal of Finance* 55:389–406.
- Fischhoff, B., P. Slovic and S. Lichtenstein (1977), "Knowing with certainty: the appropriateness of extreme confidence", *Journal of Experimental Psychology: Human Perception and Performance* 3: 552–564.
- Fisher, I. (1928), *Money Illusion* (Adelphi, New York).
- Fox, C., and A. Tversky (1995), "Ambiguity aversion and comparative ignorance", *Quarterly Journal of Economics* 110:585–603.
- French, K., and J. Poterba (1991), "Investor diversification and international equity markets", *American Economic Review* 81:222–226.
- Friedman, M. (1953), "The case for flexible exchange rates", in: *Essays in Positive Economics* (University of Chicago Press) pp. 157–203.
- Froot, K., and E. Dabora (1999), "How are stock prices affected by the location of trade?", *Journal of Financial Economics* 53:189–216.
- Genesove, D., and C. Mayer (2001), "Loss aversion and seller behavior: evidence from the housing market", *Quarterly Journal of Economics* 116:1233–1260.
- Gervais, S., and T. Odean (2001), "Learning to be overconfident", *Review of Financial Studies* 14:1–27.
- Gilboa, I., and D. Schmeidler (1989), "Maxmin expected utility with a non-unique prior", *Journal of Mathematical Economics* 18:141–153.
- Gilovich, T., R. Vallone and A. Tversky (1985), "The hot hand in basketball: on the misperception of random sequences", *Cognitive Psychology* 17:295–314.
- Gilovich, T., D. Griffin and D. Kahneman, eds (2002), *Heuristics and Biases: The Psychology of Intuitive Judgment* (Cambridge University Press).
- Gneezy, U., and J. Potters (1997), "An experiment on risk taking and evaluation periods", *Quarterly Journal of Economics* 112:631–645.
- Gompers, P., and A. Metrick (2001), "Institutional investors and equity prices", *Quarterly Journal of Economics* 116:229–259.
- Graham, B. (1949), *The Intelligent Investor: A Book of Practical Counsel* (Harper and Row, New York).
- Graham, J., and C. Harvey (2001), "The theory and practice of corporate finance: evidence from the field", *Journal of Financial Economics* 60:187–243.
- Grinblatt, M., and B. Han (2001), "The disposition effect and momentum", Working Paper (University of California, Los Angeles, CA).
- Grinblatt, M., and M. Keloharju (2001), "How distance, language, and culture influence stockholdings and trades", *Journal of Finance* 56:1053–1073.

- Grinblatt, M., and T. Moskowitz (1999), "The cross-section of expected returns and its relation to past returns", Working Paper (University of Chicago).
- Gul, F. (1991), "A theory of disappointment in decision making under uncertainty", *Econometrica* 59:667–686.
- Hansen, L., and K. Singleton (1983), "Stochastic consumption, risk aversion and the temporal behavior of asset returns", *Journal of Political Economy* 91:249–268.
- Hardouvelis, G., R. La Porta and T. Wizman (1994), "What moves the discount on country equity funds?", in: J. Frankel, ed., *The Internationalization of Equity Markets* (University of Chicago Press) pp. 345–397.
- Harris, L., and E. Gurel (1986), "Price and volume effects associated with changes in the S&P 500: new evidence for the existence of price pressure", *Journal of Finance* 41:851–860.
- Harrison, J.M., and D. Kreps (1978), "Speculative investor behavior in a stock market with heterogeneous expectations", *Quarterly Journal of Economics* 92:323–336.
- Heath, C., and A. Tversky (1991), "Preference and belief: ambiguity and competence in choice under uncertainty", *Journal of Risk and Uncertainty* 4:5–28.
- Heaton, J.B. (2002), "Managerial optimism and corporate finance", *Financial Management* (Summer), pp. 33–45.
- Hirshleifer, D. (2001), "Investor psychology and asset pricing", *Journal of Finance* 56:1533–1597.
- Hong, H., and J. Stein (1999), "A unified theory of underreaction, momentum trading, and overreaction in asset markets", *Journal of Finance* 54:2143–2184.
- Hong, H., and J. Stein (2003), "Differences of opinion, short-sale constraints and market crashes", *Review of Financial Studies* 16:487–525.
- Hong, H., T. Lim and J. Stein (2000), "Bad news travels slowly: size, analyst coverage, and the profitability of momentum strategies", *Journal of Finance* 55:265–295.
- Huberman, G. (2001), "Familiarity breeds investment", *Review of Financial Studies* 14:659–680.
- Hvidkjaer, S. (2001), "A trade-based analysis of momentum", Working Paper (University of Maryland, College Park, MD).
- Ikenberry, D., J. Lakonishok and T. Vermaelen (1995), "Market underreaction to open market share repurchases", *Journal of Financial Economics* 39:181–208.
- Jegadeesh, N., and S. Titman (1993), "Returns to buying winners and selling losers: implications for stock market efficiency", *Journal of Finance* 48:65–91.
- Jones, C., and O. Lamont (2002), "Short-sale constraints and stock returns", *Journal of Financial Economics* 66:207–239.
- Jung, K., Y. Kim and R. Stulz (1996), "Timing, investment opportunities, managerial discretion, and the security issue decision", *Journal of Financial Economics* 42:159–185.
- Kahneman, D., and A. Tversky (1974), "Judgment under uncertainty: heuristics and biases", *Science* 185:1124–1131.
- Kahneman, D., and A. Tversky (1979), "Prospect theory: an analysis of decision under risk", *Econometrica* 47:263–291.
- Kahneman, D., and A. Tversky, eds (2000), *Choices, Values and Frames* (Cambridge University Press).
- Kahneman, D., P. Slovic and A. Tversky, eds (1982), *Judgment Under Uncertainty: Heuristics and Biases* (Cambridge University Press).
- Kaul, A., V. Mehrotra and R. Morck (2000), "Demand curves for stocks do slope down: new evidence from an index weights adjustment", *Journal of Finance* 55:893–912.
- Knight, F. (1921), *Risk, Uncertainty and Profit* (Houghton Mifflin, Boston, New York).
- Korajczyk, R., D. Lucas and R. McDonald (1991), "The effects of information releases on the pricing and timing of equity issues", *Review of Financial Studies* 4:685–708.
- La Porta, R., J. Lakonishok, A. Shleifer and R. Vishny (1997), "Good news for value stocks: further evidence on market efficiency", *Journal of Finance* 49:1541–1578.
- Lakonishok, J., A. Shleifer and R. Vishny (1994), "Contrarian investment, extrapolation and risk", *Journal of Finance* 49:1541–1578.

- Lamont, O., and R. Thaler (2003), "Can the market add and subtract? Mispricing in tech stock carve-outs", *Journal of Political Economy* 111:227–268.
- Lee, C., A. Shleifer and R. Thaler (1991), "Investor sentiment and the closed-end fund puzzle", *Journal of Finance* 46:75–110.
- LeRoy, S., and R. Porter (1981), "The present-value relation: tests based on implied variance bounds", *Econometrica* 49:97–113.
- Lewis, K. (1999), "Trying to explain home bias in equities and consumption", *Journal of Economic Literature* 37:571–608.
- Lintner, J. (1956), "Distribution of incomes of corporations among dividends, retained earnings and taxes", *American Economic Review* 46:97–113.
- Loomes, G., and R. Sugden (1982), "Regret theory: an alternative theory of rational choice under uncertainty", *The Economic Journal* 92:805–824.
- Lord, C., L. Ross and M. Lepper (1979), "Biased assimilation and attitude polarization: the effects of prior theories on subsequently considered evidence", *Journal of Personality and Social Psychology* 37:2098–2109.
- Loughran, T., and J. Ritter (1995), "The new issues puzzle", *Journal of Finance* 50:23–50.
- Loughran, T., and J. Ritter (2000), "Uniformly least powerful tests of market efficiency", *Journal of Financial Economics* 55:361–389.
- Loughran, T., and J. Ritter (2002), "Why don't issuers get upset about leaving money on the table?", *Review of Financial Studies* 15:413–443.
- Loughran, T., J. Ritter and K. Rydgqvist (1994), "Initial public offerings: international insights", *Pacific Basin Finance Journal* 2:165–199.
- Lyon, J., B. Barber and C. Tsai (1999), "Improved methods for tests of long-run abnormal stock returns", *Journal of Finance* 54:165–201.
- Maenhout, P. (1999), "Robust portfolio rules and asset pricing", Working Paper (INSEAD, Paris).
- Malmendier, U., and G. Tate (2001), "CEO overconfidence and corporate investment", Working Paper (Harvard University).
- Mankiw, N.G., and S. Zeldes (1991), "The consumption of stockholders and non-stockholders", *Journal of Financial Economics* 29:97–112.
- Markowitz, H. (1952), "The utility of wealth", *Journal of Political Economy* 60:151–158.
- Mehra, R., and E. Prescott (1985), "The equity premium: a puzzle", *Journal of Monetary Economics* 15:145–161.
- Merton, R. (1987), "A simple model of capital market equilibrium with incomplete information", *Journal of Finance* 42:483–510.
- Michaely, R., R. Thaler and K. Womack (1995), "Price reactions to dividend initiations and omissions", *Journal of Finance* 50:573–608.
- Miller, E. (1977), "Risk, uncertainty and divergence of opinion", *Journal of Finance* 32:1151–1168.
- Mitchell, M., and E. Stafford (2001), "Managerial decisions and long-term stock price performance", *Journal of Business* 73:287–329.
- Mitchell, M., T. Pulvino and E. Stafford (2002), "Limited arbitrage in equity markets", *Journal of Finance* 57:551–584.
- Modigliani, F., and R. Cohn (1979), "Inflation and the stock market", *Financial Analysts Journal* 35:24–44.
- Morck, R., A. Shleifer and R. Vishny (1990), "The stock market and investment: is the market a sideshow?", *Brookings Papers on Economic Activity* 0:157–202.
- Mullainathan, S. (2001), "Thinking through categories", Working Paper (MIT, Cambridge, MA).
- Odean, T. (1998), "Are investors reluctant to realize their losses?", *Journal of Finance* 53:1775–1798.
- Odean, T. (1999), "Do investors trade too much?", *American Economic Review* 89:1279–1298.
- Ofek, E., and M. Richardson (2003), "Dot-com mania: market inefficiency in the internet sector", *Journal of Finance* 58:1113–1137.

- Pagano, M., F. Panetta and L. Zingales (1998), "Why do companies go public? An empirical analysis", *Journal of Finance* 53:27–64.
- Polk, C., and P. Sapienza (2001), "The real effects of investor sentiment", Working Paper (Northwestern University, Evanston, IL).
- Poteshman, A. (2001), "Underreaction, overreaction and increasing misreaction to information in the options market", *Journal of Finance* 56:851–876.
- Quiggin, J. (1982), "A theory of anticipated utility", *Journal of Economic Behavior and Organization* 3:323–343.
- Rabin, M. (1998), "Psychology and economics", *Journal of Economic Literature* 36:11–46.
- Rabin, M. (2000), "Risk aversion and expected utility theory: a calibration theorem", *Econometrica* 68:1281–1292.
- Rabin, M. (2002), "Inference by believers in the law of small numbers", *Quarterly Journal of Economics* 117:775–816.
- Redelmeier, D., and A. Tversky (1992), "On the framing of multiple prospects", *Psychological Science* 3:191–193.
- Ritter, J., and R. Warr (2002), "The decline of inflation and the bull market of 1982 to 1997", *Journal of Financial and Quantitative Analysis* 37:29–61.
- Roll, R. (1977), "A critique of the asset pricing theory's tests: part I", *Journal of Financial Economics* 4:129–174.
- Roll, R. (1983), "Vas ist das?", *Journal of Portfolio Management* 9:18–28.
- Roll, R. (1986), "The hubris hypothesis of corporate takeovers", *Journal of Business* 59:197–216.
- Rosenberg, B., K. Reid and R. Lanstein (1985), "Persuasive evidence of market inefficiency", *Journal of Portfolio Management* 11:9–17.
- Ross, S. (2001), *Lectures Notes on Market Efficiency* (MIT, Cambridge, MA).
- Rouwenhorst, G. (1998), "International momentum strategies", *Journal of Finance* 53:267–284.
- Rubinstein, M. (2001), "Rational markets: yes or no? The affirmative case", *Financial Analysts Journal* (May-June), pp. 15–29.
- Santos, M., and M. Woodford (1997), "Rational asset pricing bubbles", *Econometrica* 65:19–58.
- Sargent, T. (1993), *Bounded Rationality in Macroeconomics* (Oxford University Press).
- Savage, L. (1964), *The Foundations of Statistics* (Wiley, New York).
- Scheinkman, J., and W. Xiong (2003), "Overconfidence and speculative bubbles", *Journal of Political Economy*, forthcoming.
- Segal, U. (1987), "Some remarks on Quiggin's anticipated utility", *Journal of Economic Behavior and Organization* 8:145–154.
- Segal, U. (1989), "Anticipated utility: a measure representation approach", *Annals of Operations Research* 19:359–373.
- Shafir, E., P. Diamond and A. Tversky (1997), "Money illusion", *Quarterly Journal of Economics* 112:341–374.
- Shefrin, H., and M. Statman (1984), "Explaining investor preference for cash dividends", *Journal of Financial Economics* 13:253–282.
- Shefrin, H., and M. Statman (1985), "The disposition to sell winners too early and ride losers too long", *Journal of Finance* 40:777–790.
- Shiller, R. (1981), "Do stock prices move too much to be justified by subsequent changes in dividends?", *American Economic Review* 71:421–436.
- Shiller, R. (1984), "Stock prices and social dynamics", *Brookings Papers on Economic Activity* 2: 457–498.
- Shleifer, A. (1986), "Do demand curves for stocks slope down?", *Journal of Finance* 41:579–90.
- Shleifer, A. (2000), *Inefficient Markets: An Introduction to Behavioral Finance* (Oxford University Press).
- Shleifer, A., and L. Summers (1990), "The noise trader approach to finance", *Journal of Economic Perspectives* 4:19–33.

- Shleifer, A., and R. Vishny (1997), "The limits of arbitrage", *Journal of Finance* 52:35–55.
- Shleifer, A., and R. Vishny (2004), "Stock market driven acquisitions", *Journal of Financial Economics*, forthcoming.
- Stein, J. (1996), "Rational capital budgeting in an irrational world", *Journal of Business* 69:429–455.
- Summers, L. (1986), "Does the stock market rationally reflect fundamental values?", *Journal of Finance* 41:591–601.
- Thaler, R. (2000), "Mental accounting matters", in: D. Kahneman and A. Tversky, eds., *Choice, Values and Frames* (Cambridge University Press, Cambridge, UK) pp. 241–268.
- Thaler, R., and E. Johnson (1990), "Gambling with the house money and trying to break even: the effects of prior outcomes on risky choice", *Management Science* 36:643–660.
- Thaler, R., A. Tversky, D. Kahneman and A. Schwartz (1997), "The effect of myopia and loss aversion on risk-taking: an experimental test", *Quarterly Journal of Economics* 112:647–661.
- Tversky, A., and D. Kahneman (1986), "Rational choice and the framing of decisions", *Journal of Business* 59:251–278.
- Tversky, A., and D. Kahneman (1992), "Advances in prospect theory: cumulative representation of uncertainty", *Journal of Risk and Uncertainty* 5:297–323.
- Veronesi, P. (1999), "Stock market overreaction to bad news in good times: a rational expectations equilibrium model", *Review of Financial Studies* 12:975–1007.
- Vijh, A. (1994), "S&P 500 trading strategies and stock betas", *Review of Financial Studies* 7:215–251.
- von Neumann, J., and O. Morgenstern (1944), *Theory of Games and Economic Behavior* (Princeton University Press).
- Vuolteenaho, T. (2002), "What drives firm-level stock returns?", *Journal of Finance* 57:233–264.
- Weil, P. (1989), "The equity premium puzzle and the risk-free rate puzzle", *Journal of Monetary Economics* 24:401–421.
- Weinstein, N. (1980), "Unrealistic optimism about future life events", *Journal of Personality and Social Psychology* 39:806–820.
- Wurgler, J., and K. Zhuravskaya (2002), "Does arbitrage flatten demand curves for stocks?", *Journal of Business* 75:583–608.
- Yaari, M. (1987), "The dual theory of choice under risk", *Econometrica* 55:95–115.

FINANCE, OPTIMIZATION, AND THE IRREDUCIBLY IRRATIONAL COMPONENT OF HUMAN BEHAVIOR

ROBERT J. SHILLER

Yale University

Financial theory has been least successful in systematizing our knowledge about the *sources* of volatility in financial markets. This problem with financial theory has influenced my own changing research direction, and that of many others, in recent years.

The most successful applications of financial theory have been in areas where we do not need to know all the sources of volatility. For example, the theory of derivative pricing, covered here by Robert Whaley (Chapter 19), has led to an explosion of new risk management instruments. Market microstructure theory, covered here by Hans Stoll (Chapter 9) and by David Easley and Maureen O'Hara (Chapter 17), has led to significant changes in the design and regulation of trading markets. Agency theory, covered here by Jeremy Stein (Chapter 2), has led to a revolution in managerial compensation that, despite some transitory glitches, promises to make our economy much more efficient. But none of these theories depends on a systematic understanding of the ultimate sources of market volatility.

The recent international stock market boom, peaking in early 2000, and then falling in half or even further in many countries, is a stark example of this volatility. This example does not *prove* anything about anyone's model; it is only one observation. But, the reasons for this boom and crash, which has occupied the attention and anxieties of hundreds of millions of people, are certainly not elucidated by any well-established financial theory.

In looking at circumstantial evidence about this event, it seems to suggest some phenomena that are hardly discussed systematically today. The boom and crash in the stock market correspond somewhat to a boom and crash in measured earnings. But, these earnings movements do not really *explain* the market price changes, for there is no clear evidence of an exogenous cause for the earnings change: the earnings changes might just as well be caused by the stock price changes affecting individual behavior, or by human factors that caused the stock market changes, through various feedback effects.

It seems as if many of the "states of nature" that must be priced in financial markets include states of our own collective psychology. Casual observers of the stock market boom and crash note an array of cultural changes that accompanied these events.

People seemed to become more enthusiastic about investing in the stock market as the 1990s progressed, and survey evidence supports this notion. People began increasingly to believe in business, willing to accept unquestioningly their earnings statements and their pronouncements, and young people changed their aspirations in response to this market. A dramatic “new era” theory of our economy entered public perceptions, and began to be regarded as an established fact. The same theory had resonance over much of the world, and threw the stock markets of the world into unusual synchrony. The news media have a name for this phenomenon, “irrational exuberance”, but this subject has had little academic scrutiny.

There are many different approaches to understanding such events as a stock market boom and crash, and different departments in the university have their own special strengths in doing so. The psychology department offers insights into such human factors as overconfidence, self-esteem, social behavior, and attention anomalies. The sociology department offers us insights into patterns of collective thinking and belief in authority. The history department will have a different kind of insight, perhaps a more comprehensive, more inductive approach to understanding financial phenomena. We must not forget that when we study financial data we are studying history.

But most financial theory derives formally from an abstract paradigm of individual optimization. Optimizing models have been the cornerstone of economic and financial theory. Such models are indeed the dominant tool of papers in this book.

It is common to see behavioral finance as quite a different thing, based on a “non-optimizing” theory. But what does it mean to be optimizing or not? Presumably, one could describe any human behavior as the solution to *some* optimization problem. We see described in this handbook, under the rubric of optimizing models, many variations on traditional expected utility theory. We see discussions, in Chapter 13 by John Campbell and Chapter 14 by Rajnish Mehra and Edward Prescott, of Epstein-Zin utility, which is not technically an expected utility maximization theory. We see discussions as well of habit formation, or keeping up with the Jones’s. These are pushing out the borders of rational optimization theory even if they can be modeled in terms of some form of expected utility. Even prospect theory, the most influential construct of behavioral economics, was couched by its framers, Daniel Kahneman and Amos Tversky, as an optimization problem.

What sets the behavioral theories apart from mainstream optimizing models in most cases is their insistence on the inconsistencies of human judgment. People’s responses to any given situation are affected by framing, by salience, and by the internal dynamics of their own attention. An insistence that people are inconsistent in their behavior will never be the basis of an *elegant* financial theory.

If one is to have a theory of stock prices, of derivatives prices, of corporate investment, of banking and other intermediaries, all subjects of this volume, one is naturally drawn to the idea that people are consistently rational optimizers of some well-defined and sensible objective. For all these phenomena are related to financial instruments: securities, options and corporate loans are tools that people use with purpose, and so a theory of these would naturally begin with an understanding of this

real purpose. Just as to understand a hammer or drill, one must reflect on what their use is, one must understand the complex uses of financial instruments and institutions. We cannot proceed on behavioral finance alone.

Dybvig and Ross, in Chapter 10 in this volume, get at what they view as a big problem with behavioral finance. They argue that a problem with the psychological theories is that "... they tend to be isolated stories rather than general specifications and are often hard to generalize. For example, prospect theory says that agents put extra weight on very unlikely outcomes, but it is not clear what this means in a model with a continuum of states" (p. 609). They are indirectly referring to the issue of psychological framing, which is discussed in detail in this Handbook by Nicholas Barberis and Richard Thaler (Chapter 18). Framing refers to people's tendency to act inconsistently, so that their behavior depends on a *suggested* or *currently convenient* frame of reference, and hence can be inconsistent from time to time. Prospect theory says that people tend often to overestimate small probabilities that are salient to them from their frame of reference, but of course it cannot mean that people overestimate all small probabilities. Prospect theory is ultimately a theory of people's faulty attention mechanism and cannot be the basis for an overarching financial theory. But, it should be part of the adjustments made to the theory before it is applied.

There is a tendency for some observers to dismiss behavioral finance because of the transience of financial anomalies. After an anomaly, such as the January effect or the small-firm effect, is discovered and the discovery is given news-media attention, the anomaly tends to disappear. William Schwert, in his Chapter 15, asserts that "All these findings raise the possibility that the anomalies are more apparent than real" (p. 941). But, these anomalies are very real even though they are transient, since they ultimately account for a significant part of market volatility. If our paradigm emphasizes the inconsistency of individual behavior, then it accords with changing anomalies.

I wish to argue here that both approaches to finance, the behavioral approach, and the rational optimizing approach, have their own contributions to make, and that much work remains to be done on integrating them. There are not enough people who take an active and constructive interest in both approaches. When there is a conflict of paradigms, as appears to be the case here, it is often most fruitful for research to be conducted at the point of conflict between the paradigms. There is a definite schism in much of this volume, with the references covered in the one chapter on behavioral finance figuring at most peripherally in the other essays.

I find it odd that there should be a "field" called behavioral finance. Today, we have graduate students naming this as a field for their oral examinations. But all branches of finance should take account of the various social sciences, and ideally we would neither be giving oral examinations in behavioral finance nor be corralling all of the resulting insights into a single chapter of the handbook. I take the isolation of the single behavioral finance chapter in this handbook as a sign that acceptance and understanding of insights from other social sciences has only just begun to permeate the finance profession. When the next major edition of this handbook appears, let us suppose in a decade or more, I would not expect to see a chapter on behavioral finance. The material

should be thoroughly dispersed among most of the chapters. One may hope that by then some way of integrating it in a productive way with the insights that we have from our optimizing models will be generally understood.