


Do Cross-Sectional Predictors Contain Systematic Information?

Joseph Engelberg 
University of California, San Diego Rady School of Management
jengelberg@ucsd.edu

R. David McLean
Georgetown University, McDonough School of Business
dm1448@georgetown.edu

Jeffrey Pontiff
Boston College, Carroll School of Management
pontiff@bc.edu (corresponding author)

Matthew C. Ringgenberg 
University of Utah, David Eccles School of Business
matthew.ringgenberg@eccles.utah.edu

Abstract

Firm-level variables that predict cross-sectional stock returns, such as price-to-earnings and short interest, are often averaged and used to predict market returns. Using various samples of cross-sectional predictors and accounting for the number of predictors and their interdependence, we find only weak evidence that cross-sectional predictors make good time-series predictors, especially out-of-sample. The results suggest that cross-sectional predictors do not generally contain systematic information.

I. Introduction

Is the market risk premium predictable? Financial research has strived to answer this question going back to at least to Dow (1920), and a large number of variables have been used to predict the time-series of stock returns. A growing number of papers with time-series predictors begin with a cross-sectional predictor and ask whether that cross-sectional variable can be aggregated and then used to predict market returns. Such papers presume that cross-sectional predictors contain systematic information. For example, Baker and Wurgler (2000) cite the extensive market-timing literature, which relates equity issuance to returns in the cross-section, and suggest “issuers try to time *both* their idiosyncratic return and the

The authors thank an anonymous referee, Hendrik Bessembinder (the editor), John Campbell, Mike Cooper, Amit Goyal, Robin Greenwood, Campbell Harvey, Travis Johnson, Bryan Kelly, Owen Lamont, Yan Liu, Seth Pruitt, Allan Timmermann, and Michael Wolf, and conference and seminar participants at the 2018 Society for Financial Studies Cavalcade, the 2019 American Finance Association, MIT (Accounting), TCU, UC-Berkeley, University of Kentucky, University of Michigan, UNLV, University of Utah, UC Riverside, University of Virginia, Washington University in St. Louis, University of Oxford, and Warwick Business School. All errors are our own.

market return.” They find that an aggregated capital raising variable predicts market returns. Similarly, Lewellen (2004) finds that aggregated price-to-fundamental ratios based on dividends, book-to-market, and earnings predict market returns. Robert Shiller’s famous CAPE ratio is based on the idea that aggregated price-to-earnings ratios predict market returns. Hirshleifer, Hou, and Teoh (2009) aggregate Sloan’s (1996) accrual anomaly and use it to predict the market. Baker and Wurgler (2007) incorporate several firm-level variables to create an investor sentiment index and show that it predicts market returns.

Drawing on this large literature, we examine market return predictability from cross-sectional variables through the lens of multiple hypothesis testing. After all, there are many cross-sectional predictors that *could be chosen* in order to predict market returns. Published papers present individual hypothesis tests for predictors that *were chosen*. This article asks: do cross-sectional variables generally aggregate to make good time-series predictors? Using various samples of cross-sectional predictors and incorporating a multiple hypothesis testing framework, we find that the answer is “no.”

We begin by creating a sample of time-series predictors constructed from the population of cross-sectional predictors documented in the literature. We use a sample of 140 cross-sectional predictors that is essentially the population of firm-level characteristics that have demonstrated cross-sectional predictability in the academic literature.¹ Of these 140 predictors, 26 have already been aggregated and used to predict market returns in a published study.² These papers are well-cited, having received over 12,000 Google Scholar citations in total.

It is impossible to know the true set of variables examined by researchers. Statistically significant findings are more likely to be published, so the 26 variables in the aforementioned papers represent a lower bound on the set of predictors considered by researchers. We begin our analysis by examining this sample; we then extend our analyses to examine other subsets motivated by economic theory, as well as the entire sample of 140 predictors that *could* have been considered. Our analyses thus examine the minimum and maximum number of cross-sectional variables that could have been used to generate time-series variables.

We begin by taking each of the 140 firm-level predictors and calculating monthly cross-sectional averages to get a single, monthly value. For each predictor, we construct equal-weighted and value-weighted averages; the resulting database has 280 different predictive variables (140 equal-weighted and 140 value-weighted).³ To examine the market-level predictability of these variables, we perform both in-sample and out-of-sample tests. Like other papers in the time-series literature, our in-sample tests use the entire sample of data and estimate a single parameter estimate from a time-series regression of the market risk premium on the predictor. Our out-of-sample tests consist of expanding, rolling-window regressions

¹This builds on the 97-predictor list used in McLean and Pontiff (2016) and Engelberg, McLean, and Pontiff (2018).

²Table IA.VI in the Supplementary Material lists these 26 cross-sectional predictors along with the 23 time-series papers that each was featured in.

³We follow existing practice when constructing aggregate predictors from cross-sectional data: although some researchers consider equal-weights (e.g., Goyal and Santa-Clara (2003), Hirshleifer et al. (2009)), value-weights are most common. For completeness, we present results with both types of weighting.

that use only information available at each point in time to examine whether a predictor is useful for forecasting the market's risk premium.

At first glance, it appears that many cross-sectional predictors are good market-level predictors in-sample. When we examine the predictors already studied by the existing literature, we find that 20% of them predict 1-year market returns in an ordinary least squares (OLS) regression with coefficients that are significant at the 10% level or better, and 8% of them are significant at the 1% level.⁴ The strength of this result is strongly related to the horizon of predictability: when considering 1-month market returns, the number of significant variables falls to 12% at the 10% level and 2% at the 1% level.

We then consider two predictor subgroups that can clearly be motivated by economic theory. The first is a subsample of VALUATION predictors, which are variables that are based on ratios of market prices to fundamentals. Many VALUATION predictors have received attention in the market risk premium literature, including the dividend-to-price and earnings-to-price ratios. Moreover, since virtually all VALUATION ratios should be a function of discount rates, theory suggests that they should all work in a time-series setting (Lewellen (2004), Kelly and Pruitt (2013)). We also form a subsample of OPINION predictors (e.g., institutional trading, analyst upgrades) which can be motivated with the sentiment explanation of Baker and Wurgler (2006) or with the information explanation of Seyhun (1988). However, despite the economic motivation for these subsamples, we find only weak in-sample predictability for the VALUATION and OPINION subcategories. In fact, the VALUATION and OPINION subcategories perform worse than the sample of predictors from the existing literature, suggesting the predictors in the existing literature are a special subset of all cross-sectional variables.

We also examine a third subsample, BEST_CROSS-SECTIONAL, which consists of the 10 predictors with the highest cross-sectional *t*-statistics. A number of papers have found evidence that cross-sectional return predictability is smaller in recent periods (e.g., McLean and Pontiff (2016), Green, Hand, and Zhang (2017)), so it is possible that the weak performance of some of our predictors (like those in the VALUATION and OPINION subsamples) results from weakness of the underlying cross-sectional predictors. To account for this possibility, we examine a subset of predictors formed from the cross-sectional variables with the best performance. Yet, once again, for this subsample, we find only weak evidence of return predictability.

The 51 predictors from the existing literature that we discuss above represent a lower bound on the set of variables actually considered. Harvey, Liu, and Zhu (2016) note that there is a publication bias in Finance, because journals are less likely to publish results that are not statistically significant. Chen (2021) notes that selective reporting of results makes it difficult to know the true set of variables considered, and this has implications for statistical inference. We therefore examine all 280 possible predictors constructed from the 140 cross-sectional variables in the

⁴There are 51 variables since we create both equal-weighted and value-weighted predictors for each of the 26 variables in the existing literature and one of the resulting 52 possible variables is dropped because it is nonstationary. While we do not find that 51 out of 51 are significant, it is important to note these are reanalyses, not replications, in the language of Welch (2019) because we do not use the same sample or code used in the original papers. Moreover, we examine two versions of each predictor (equal-weighted and value-weighted) while many of the original papers only examine one of these.

literature – this set represents an upper bound on the cross-sectional predictors that could have been used to make time-series predictors. Out of the 280 predictive variables, 253 of them are stationary and exhibit sufficient time-series variation.⁵ Of these 253 predictors, 43 (17%) predict 1-year market returns in an OLS regression with coefficients that are significant at the 10% level or better, and 14 (6%) are significant at the 1% level. Again, the strength of this result is strongly related to the horizon of predictability: when considering 1-month market returns, only 27 of the 253 predictors (11%) are significant at the 10% level and 7 (3%) are significant at the 1% level. We also find that several cross-sectional predictors (such as ASSET_TURNOVER and Z_SCORE) are among the best performers for market predictability, but have yet to be documented in this literature. For example, value-weighted ASSET_TURNOVER, which has not been previously proposed as a predictor of market risk premia, predicts the market risk premium with an R^2 of 17.8% at the 1-year horizon.

Since we examine a large number of predictors, we expect that some variables will appear significant by chance. Moreover, many of these variables are related to each other so the tests we conduct are not independent. To address both the number of tests we conduct as well as the dependency among the tests, we perform the Romano and Wolf (2016) resampling-based stepdown procedure to compute adjusted p -values that control the family-wise error rate while accounting for the number of tests and dependence. Among methods that control the family-wise error rate, the Romano and Wolf (2016) procedure is preferable to the well-known Bonferroni (Dunn (1961)) or Holm (1979) tests because it considers the dependence structure across multiple tests and thus has more power (Romano and Wolf (2005), (2016)).⁶ Our paper is the first to apply this procedure to predicting the equity market risk premium. In robustness tests, we also estimate adjusted p -values while controlling the false discovery rate (Benjamini and Yekutieli (2001)). This approach is less conservative than controlling the family-wise error rate (Romano et al. (2008)).⁷ We reach the same conclusion with both approaches.

When we apply the Romano and Wolf (2016) stepdown procedure to the 51 predictors from the existing literature, we find weaker evidence that cross-sectional predictors contain systematic information. While 10 are significant at the 1-year horizon using single hypothesis testing, only 3 are significant at the 1-year horizon using the Romano and Wolf procedure. The results suggest the predictors in

⁵We test each of the 280 predictors for a unit root using an Augmented Dickey–Fuller (1979) test. If we fail to reject the null that a variable is nonstationary, we then calculate the first-difference. If we fail to reject the null that the first-differenced variable is nonstationary, we drop the variable. We also drop variables that aggregate to form a variable that does not exhibit time-series variation and we filter variables that should not aggregate according to economic logic. See Section II for details.

⁶Romano, Shaikh, and Wolf (2008) examine simulation evidence for a variety of multiple testing techniques and find that the Romano and Wolf (2005) procedure has good power, especially relative to methods that do not account for the dependence of the individual test statistics. While some other methods that control generalized error rates have even better power, this comes at the cost of having higher false rejection rates. They state, “It appears that when the number of strategies is in the thousands, controlling the [family-wise error rate] becomes too stringent.” Since we examine a maximum of 269 strategies, we focus on controlling the family-wise error rate.

⁷Romano et al. (2008) show that false discovery rate methods generally exhibit fewer false negatives at the cost of allowing more false positives.

the existing literature may be a result of selective testing and reporting. Of course, it is not possible to know the true set of cross-sectional variables considered, and Chen (2021) shows that accounting for this could raise or lower the bar for statistical significance. To account for this, we then examine our other samples using the Romano and Wolf procedure. For the VALUATION and OPINION subsamples, as well as the BEST_CROSS-SECTIONAL subsample, we find weak evidence of return predictability using multiple testing methods.

We then turn to our full sample of predictors. Using in-sample regressions, when we examine the full set of 253 predictors and perform the Romano and Wolf (2016) stepdown procedure we are unable to reject the null of no-predictability at the 1% level for any predictor at any horizon. At the 1-month and 3-month horizons, no variable is significant even at the 10% level. At the 12-month horizon only two variables, Z_SCORE and ASSET_TURNOVER, have Romano and Wolf p -values that are less than 5%. When we examine more permissive p -values based on controlling the false discovery rate, there are only a few more significant predictors. We find that one predictor is significant at the 5% level at the 1-month horizon and 8 (out of 253) are significant at the 12-month horizon. In short, we find most of the cross-sectional variables that were statistically significant when examined in isolation are no longer significant when examined in the context of all cross-sectional predictors. However, a few variables, like Z_SCORE and ASSET_TURNOVER, do still show evidence of return predictability.

Goyal and Welch (2008) argue that out-of-sample regressions serve as a useful diagnostic. Accordingly, we turn to out-of-sample forecasting regressions. Again, we start with the predictors already examined in the literature, and then examine other samples of predictors including the subgroups motivated by economic theory and the set of all 140 cross-sectional variables that could have been considered. In all samples, we find that things look even bleaker for cross-sectional predictors when we consider out-of-sample tests. While 11 out of the 51 predictors from the existing literature are significant at the 12-month horizon, none are significant at the 1-month horizon. Similarly, among the 253 stationary predictors, we examine out-of-sample regressions. We find that 3% significantly predict market returns at the 1-month horizon and 17% predict market returns at the 12-month horizon (compared with 11% and 17% in-sample). Moreover, once we adjust the individual p -values using the Romano and Wolf (2016) stepdown procedure, we no longer find evidence against the null at the 10% level for any predictor at any horizon. For example, in our out-of-sample tests, ASSET_TURNOVER again appears to be a strong predictor. It exhibits positive out-of-sample R^2 values at every forecasting horizon, with an impressive R^2 of 17% at the 1-year horizon. However, the corresponding Romano and Wolf (2016) adjusted p -value is not statistically significant at any of the usual levels. We find similar results when controlling the false discovery rate. Taking all of the results together, our findings suggest that cross-sectional predictors do not contain systematic information.

Our paper makes a number of contributions. First, we provide context for existing papers that propose a particular cross-sectional predictor should be transformed into a time-series one. Existing results in the literature that have economically large coefficients and impressive t -statistics could be the result of chance if enough predictors are considered. Our findings show that cross-sectional predictors

as a group do not make good time-series predictors. The literature's elevation of individual predictors should consider the economic motivation behind each predictor, and the results should be interpreted in light of our data mining critique.

Second, our results provide new insight into the nature of return predictability, both in the cross-section and the time-series. By aggregating cross-sectional predictors into time-series variables, we are able to understand whether cross-sectional variables contain information about the systematic components of returns. Our results suggest they do not. This is surprising, as Wen (2019) finds that asset growth can be aggregated to predict market returns, and several studies (e.g., Hou, Xue, and Zhang (2015)) find that factor models with an investment factor can explain many cross-sectional anomalies. Taken together, these papers suggest that most anomalies should aggregate to predict market returns, yet we find little evidence that they do.

Finally, we contribute to the extensive equity return premium literature. While Goyal and Welch (2008) show that 14 popular time-series variables do not significantly predict returns in out-of-sample tests, subsequent papers have documented evidence of return predictability using firm-level variables aggregated across stocks (e.g., Hirshleifer et al. (2009), Rapach, Ringgenberg, and Zhou (2016)). Our results extend these findings by showing that several other cross-sectional predictors can be aggregated to form good time-series predictors. However, we find that many of these predictors are no longer significant after adjusting for multiple hypothesis testing. Our findings emphasize the importance of considering the impact of data snooping bias when examining time-series return predictability.

The remainder of this article proceeds as follows: [Section II](#) briefly describes the existing literature and outlines the theoretical relation between cross-sectional predictive variables and time-series return predictability. [Section III](#) describes the data used in this study. [Section IV](#) characterizes our findings and [Section V](#) concludes.

II. Background

Financial researchers have examined the predictability of stock returns for over a century (e.g., Gibson (1906)) and a large literature has documented evidence of predictability in the cross-section of stock returns. A separate literature has examined the predictability of the equity risk-premium using time-series predictive variables. To date, these two literatures have evolved relatively independently. We connect these two literatures.

A. Time-Series Return Predictability

A number of papers find in-sample evidence of time-series return predictability, but out-of-sample evidence is rare, suggesting that many predictors are the result of data snooping (i.e., overfitting). For example, Bossaerts and Hillion (1999) use model selection criteria from the statistics literature to choose candidate predictors, which allows them to partially avoid data snooping biases, yet they find that the resulting predictors are unable to forecast out-of-sample returns. Similarly, Goyal and Welch (2008) examine 14 popular predictors from the existing literature and

find that they fail to forecast the equity risk premium in out-of-sample tests. Cooper and Gulen (2006) note that researchers have many different choices regarding the specification of predictability tests, including the predictor variables, the estimation periods, and the assets being forecasted. They perform specification searches across these parameters and find that return predictability is highly sensitive to these parameter choices. More recently, Bartsch, Dichtl, Drobetz, and Neuhierl (2021) examine a wide variety of possible permutations of the predictors in Goyal and Welch (2008) and the technical predictors in Neely, Rapach, Tu, and Zhou (2014) and estimate that most out-of-sample performance for these variables is from data snooping.

In light of the poor performance of many predictive variables in out-of-sample tests, researchers have focused on developing methodologies that are robust to data snooping concerns (see Harvey, Liu, and Saretto (2020)). Foster, Smith, and Whaley (1997) develop a procedure to account for data snooping biases when evaluating the fit of predictive regressions. White (2000) develops a reality check bootstrap (RCB) to account for data snooping biases that result from specification searches, and Sullivan, Timmermann, and White (1999) apply the RCB procedure to a set of technical trading rules. While the White RCB procedure determines whether the best predictor among a group is statistically significant after adjusting for data snooping biases, Romano and Wolf (2005), (2016) show how to adjust the p -value for each individual predictor to account for data snooping biases. We are the first to apply the Romano and Wolf stepdown procedure to a large set of predictive variables derived from existing academic studies. In robustness checks, we also examine adjusted p -values that control the false discovery rate (Benjamini and Yekutieli (2001)) which is less conservative than the Romano and Wolf procedure.

B. Cross-Sectional Return Predictability

While the literature on time-series return predictability has generally found that most predictors fail to perform in out-of-sample tests, a large literature finds evidence of return predictability in the cross-section of stocks. More recently, a number of papers reaffirm earlier studies that find return predictability in the cross-section of stocks. For example, using quintile portfolios, McLean and Pontiff (2016) find that 88% of firm-level predictors generate t -statistics greater than 1.5 in the sample period used by the original study. Chen and Zimmerman (2020) replicate 319 firm-level predictors and find that only 3 fail to reproduce the statistical significance of the original study. Some studies question the breadth of predictability and whether implementable portfolio strategies are possible. Hou et al. (2015) use a sample of 452 anomalies and conclude that most fail to replicate. Chen and Zimmerman reconcile their results with Hou et al. (2015) by noting that Hou et al. eliminate microcap stocks from their sample (60% of CRSP), that many of Hou et al.'s anomalies are not unique and are based on the same firm-level variable measured over different horizons or intervals, and that many of Hou et al.'s anomalies are unique to their paper and were never significant to begin with. The fact that Hou et al. find that eliminating microcaps weakens return-predictability is consistent with many studies (see Barberis and Thaler (2003) or Pontiff (2006) for a

thorough review) which find that anomalies are weaker in large stocks and strongest in small stocks that are more costly to hold and trade (Pontiff (1996)).

There is though, an evolving discussion centered on the extent to which cross-sectional predictability is the outcome of data mining and whether the literature has overlooked multiple testing issues. Harvey et al. (2016) conclude that the majority of findings in financial economics papers are likely false, and that *t*-statistic hurdles should be raised to 3. McLean and Pontiff (2016) examine 97 anomalies and find that at most 25% of anomaly returns can be explained by data mining. Linnainmaa and Roberts (2018) conclude that the majority of accounting-based anomalies are the outcome of data mining.

Yet other recent studies find that data mining is likely not the first-order cause of anomaly returns. McLean and Pontiff (2016) find that anomaly returns are 50% lower post-publication, an effect they attribute to both data mining and arbitrage informed by the original publication. Importantly, they reject the null that anomaly returns are zero post-publication. Jacobs and Muller (2020), find that anomaly returns are, on average, significant outside of the US, and do not decline post-publication as they do in the USA. Chen and Zimmerman (2020) find that publication bias only accounts for about 12% of anomaly returns. Chen (2021) argues that it would take 10,000 researchers hundreds of years to produce the large *t*-statistics reported in the anomaly literature.

Overall, the existence, and nature, of cross-sectional return predictability is an important debate that has yet to be resolved. However, none of the existing studies examine the relation between cross-sectional predictors and aggregate market returns, which is the focus of our paper.

C. The Information in Cross-Sectional Variables

While there are extensive literatures on both time-series predictability and cross-sectional predictability, with a few notable exceptions they have evolved independently. Several papers show that firm-level anomalies aggregate to market-wide predictors. As we mentioned earlier, to the best of our knowledge 26 predictors have been aggregated and used to predict market returns in a published study. Table 1 lists these 26 cross-sectional predictors along with the 23 time-series papers that each was featured in. Some examples include Pontiff and Schall (1998) with book-to-market ratios, Campbell and Shiller (1988) with P/E ratios, and Chordia, Roll, and Subrahmanyam (2002) with insider trading. More recently, Hirshleifer et al. (2009) find that firm-level accruals and cash flow, when aggregated across stocks, contain information about market returns, and Wen (2019) shows that aggregate asset growth predicts market returns. Finally, Rapach et al. (2016) show that firm-level short interest aggregates to form one of the strongest known predictors of market returns.

D. The Information in Nonsystematic Predictors

In this article, our goal is to understand the relation between cross-sectional return predictability and time-series return predictability. While it may seem natural

TABLE 1
Summary Statistics

Table 1 displays summary statistics across all of the cross-sectional variables from which we construct time-series predictors. To construct time-series predictors out of cross-sectional variables, we calculate the value-weighted and equal-weighted mean across all stocks on each date. The first row displays statistics for variables already examined in the existing literature on time-series return predictability (PREDICTORS_FROM_EXISTING_PAPERS), the second row shows statistics for all possible variables derived from the cross-sectional literature (ALL_POSSIBLE), and the remaining rows examine subsamples formed on the 10 most statistically significant cross-sectional predictors (BEST_CROSS-SECTIONAL) and two different groupings from the categories in McLean and Pontiff (2016): i) OPINION and ii) VALUATION. We display the mean and median values of cross-sectional t -statistics as well as the number of citations for each predictor from Google Scholar as of 2018.

Type	N	Cross-Sectional t -Statistic		Citations	
		Mean	Median	Mean	Median
PREDICTORS_FROM_EXISTING_PAPERS	26	4.26	3.26	3,312	1,531
ALL_POSSIBLE	140	3.03	2.71	1,432	517
BEST_CROSS-SECTIONAL	10	8.38	7.96	854	740
OPINION	22	2.97	2.48	798	689
VALUATION	15	3.14	2.88	2,262	273

that cross-sectional return predictors should aggregate to generate time-series return predictors, it is possible to have one without the other.⁸ Specifically, cross-sectional variables could predict returns because they forecast the systematic component of returns or the nonsystematic component of returns. As such, cross-sectional return predictors do not necessarily aggregate to form good time-series predictors. To see this, define a variable $X_{i,t}^{\text{Nonsyst}}$ as a nonsystematic predictor if it forecasts the *nonsystematic* portion of stock returns for asset i on date $t + 1$. Without loss of generalization,⁹ we can express the return on asset i using the market model:

$$(1) \quad R_{i,t} = R_f + \beta_i (R_{m,t} - R_f) + \varepsilon_{i,t},$$

where $R_{i,t}$ is the return on stock i on date t , R_f is the risk-free rate, $R_{m,t}$ is the market return, and $\varepsilon_{i,t}$ is the portion of stock i 's return that is orthogonal to the market's return. We define a nonsystematic predictor as a variable $X_{i,t-1}^{\text{Nonsyst}}$ that satisfies $\gamma_1 \neq 0$ in a linear regression of the form¹⁰:

$$(2) \quad \widehat{\varepsilon}_{i,t} = \gamma_0 + \gamma_1 X_{i,t-1}^{\text{Nonsyst}} + \omega_{i,t},$$

where $\widehat{\varepsilon}_{i,t}$ is the abnormal return from the market model (Sharpe (1964), Lintner (1965)). In other words, a nonsystematic predictor, by definition, forecasts the portion of asset i 's return that is *not* explained by aggregate market movements. However, while $X_{i,t-1}^{\text{Nonsyst}}$ contains information about individual stock returns, it

⁸Indeed, several existing papers find that firm-level relations do not hold at the aggregate-level. Kothari, Lewellen, and Warner (2006) document a negative relation between returns and earnings surprise at the aggregate-level, in contrast to the positive relation documented at the firm-level. Similarly, Hirshleifer et al. (2009) find that the relation between accruals and returns changes sign between firm-level and aggregate-level analyses.

⁹In Section B of the Supplementary Material, we show the logic in this section generalizes to any model with nonsystematic and systematic components including the Fama and French (1992) 3-factor model and the Fama and French (2015) 5-factor model.

¹⁰For simplicity, we ignore the sign of the abnormal return and define an anomaly as any variable that predicts abnormal returns in either direction.

will not aggregate to generate market return predictability. To see this, average equation (2) across assets all N stocks in the economy and multiply both sides by $\frac{\text{mc}_{i,t}}{\sum_i \text{mc}_{i,t}}$, where $\text{mc}_{i,t}$ is the market capitalization of stock i on date t :

$$(3) \quad \frac{\text{mc}_{i,t}}{\sum_i \text{mc}_{i,t}} \sum_{i=1}^N [R_{i,t} - R_f - \beta_i (R_{m,t} - R_f)] = \bar{\gamma}_0 + \bar{\gamma}_1 \overline{X_{t-1}^{\text{Nonsyst}}},$$

where the bar above a variable denotes the value-weighted mean. It is simple to show that the left-hand side of equation (3) is equal to zero. Thus, the value-weighted variable $\overline{X_{t-1}^{\text{Nonsyst}}}$ contains *no* information about aggregate market returns.

E. The Information in Systematic Predictors

While nonsystematic predictors contain no information about aggregate market returns, it is possible to have a variable that predicts information in the cross-section that contains information about the aggregate risk-premium. Define a variable $X_{i,t-1}^{\text{systr}}$ as a systematic predictor if it forecasts the *systematic* portion of stock returns for asset i on date $t + 1$. Thus, define a systematic predictor as a variable $X_{i,t-1}^{\text{systr}}$ that satisfies $\gamma_1 \neq 0$ in a linear regression of the form:

$$(4) \quad \beta_i (R_{m,t} - R_f) = \gamma_0 + \gamma_1 X_{i,t-1}^{\text{systr}} + \omega_{i,t}.$$

Because the market beta is 1, it is easy to show that the left-hand side of equation (4) implies a direct linear relation between the predictor variable and the market risk-premium. Notice also that this relation goes in both directions: if a time-series predictor is constructed from individual assets, it *must* contain information about the systematic portion of individual asset returns.^{11,12} This implication provides additional economic information to test the validity of proposed predictors. In other words, when evaluating predictors constructed from individual characteristics, we should focus on the subsample of individual characteristics that contain information about individual asset returns. Accordingly, in the rest of the article, we examine the aggregate information in a set of 140 predictors that have been previously shown to contain information about individual asset returns (McLean and Pontiff (2016)).

¹¹In the market model (Sharpe (1964), Lintner (1965)), the systematic portion of stock returns reflects compensation for bearing systematic risk. However, outside of the market model, a common component of returns could exist that is not related to systematic risk (e.g., consumer sentiment). Similar to the systematic portion of returns in the market model, such a variable could be related to the cross-section of returns and, since it has a common component, it could aggregate to contain information about the equity risk premium.

¹²Theoretically, it is possible that investors switch from gathering systematic information to gathering idiosyncratic information depending on economic conditions as in Kacperczyk, Van Nieuwerburgh, and Veldkamp (2016). As a result, some variables could contain idiosyncratic information at times, and systematic information at other times.

III. Data

To examine the relation between cross-sectional anomaly variables and the equity risk premium, we combine daily data from the Center for Research in Security Prices (CRSP) and Compustat over the period 1926 through 2017.

We calculate the equity risk premium as the log return on the S&P 500 index minus the log return on a 1-month Treasury bill as in Goyal and Welch (2008).¹³ We construct aggregate time-series variables for each of the 97 cross-sectional anomalies in McLean and Pontiff (2016) and we supplement this data set with 43 additional variables from the extant literature to arrive at 140 candidate predictors.¹⁴ As we explain in Section I, we explore various subgroups within the 140. Table 1 contains summary statistics for the cross-sectional variables we use to form time-series predictors. Of the 140 variables, 15 are classified as VALUATION predictors and 22 are classified as OPINION predictors. BEST_CROSS-SECTIONAL. This subcategory consists of the 10 predictors with the greatest cross-sectional *t*-statistic, using the sample period in the original paper. By definition, the BEST_CROSS-SECTIONAL predictors exhibit strong cross-sectional predictive power; indeed, the mean cross-sectional *t*-statistic for these predictors is approximately 9. Predictors from existing papers consist of the 26 cross-sectional variables that have been used to predict market returns in a published study.

We construct two time-series predictors from each cross-sectional variable based on the equal-weighted average and value-weighted average. Unlike with cross-sectional estimation, consideration of stationarity is crucial in estimation of the market risk premium (e.g., Campbell (1991), Hodrick (1992)). We test each time-series predictor for a unit root using an Augmented Dickey–Fuller (1979) test. Because some of the resulting time-series variables are nonstationary, we proceed as follows: if we reject the null that the raw (untransformed) variable is nonstationary, we use it in our tests. If not, we calculate the first difference for each variable; if we reject the null that it is nonstationary then we use it, otherwise the variable is excluded.

Some of our cross-sectional predictors should (theoretically) aggregate to form a variable that is constant across time (e.g., CAPM beta). Accordingly, after constructing time-series versions of each variable, we apply a manual filter to drop variables that should not aggregate to form an economically meaningful variable and we also examine each variable's time-series standard deviation and exclude predictors that exhibit little to no time-series variation.¹⁵ Of our original 280 different predictive variables (140 equal-weighted and 140 value-weighted variables), 253 survive the screening procedures discussed above. We use these 253 variables in our main regressions.

To verify that the Romano and Wolf (2016) adjusted *p*-values are not sensitive to the set of variables we consider, we also examine three alternate sets of candidate predictors: one that has more variables (269 variables) and two that have fewer

¹³We download this data from Amit Goyal's website (<http://www.hec.unil.ch/agoyal/>).

¹⁴See the Supplementary Material for an overview of the construction of these 140 variables.

¹⁵We calculate a measure of variation for each variable as the ratio of its time-series standard deviation to the absolute value of its time-series mean, and we drop variables with a ratio less than 0.06. See the Supplementary Material for a list of variables dropped using the manual filter.

(137 and 51 variables, respectively). The first set expands the list of predictors to 269 variables by omitting the manual filter discussed above and adding stochastic detrending as a way to avoid nonstationarity. Specifically, if we reject the null that the raw (untransformed) variable is nonstationary, we use it in our tests. If not, we calculate deviations from a linear trend model for each variable; if we reject the null that it is nonstationary then we use it. If not, we calculate the first difference for each variable; if we reject the null that it is nonstationary then we use it, otherwise we drop the variable. For the linear trend, we estimate a model of the form:

$$(5) \quad x_t = a + bt + u_t \text{ for } t = 1, \dots, T,$$

for each predictor variable x_t and time period t . We take the fitted residual, \hat{u}_t , as our detrended measure. By construction, \hat{u}_t has a mean of zero and we normalize it to have a standard deviation of 1.¹⁶ This set represents the maximum number of candidate variables possible; we use all of the candidate predictors that we can with minimal assumptions (i.e., they must have time-series variation and pass the unit root test).

One potential criticism is that there may be only a few potent predictor variables. In this view of the world, a large sample such as ours might be padded with obviously unlikely candidates, and a multiple hypothesis test would have low statistical power. To address this potential criticism, we consider two smaller sets of predictors. One set contains 137 predictors and uses only the raw version of each variable (i.e., it does not use first-differencing or stochastic detrending). This selection reflects the notion that authors, referees, and editors avoid time-series variables that require extra manipulation. For this set, we again apply the manual filter to drop predictors that exhibit little to no time-series variation. The other set represents the minimum number of candidate variables possible; we use only the predictors that have *already* been examined in the existing time-series literature. This set contains 51 variables and, by definition, contains variables that the literature views as reasonable candidate variables.

IV. Results

In this section, we examine whether cross-sectional predictors, in general, contain information about the equity risk premium. We start by examining in-sample tests that use the entire sample of data and estimate a single parameter estimate from a time-series regression of the market risk premium on the predictor. We then examine out-of-sample tests that use rolling regressions to test whether a variable is useful for predicting the future equity risk premium, using only information available at each date.

A. In-Sample Tests

As noted in Goyal and Welch (2008), “It is unreasonable to propose a model if the [in-sample] performance is insignificant, regardless of its [out-of-sample]

¹⁶For our in-sample analyses, we estimate the linear trend model using all available data. For our out-of-sample analyses, we estimate the trend model only using data available at each point in time to avoid a look-ahead bias.

performance.” As discussed in Section II, we start with 140 variables from the existing literature and, using these, we form 268 candidate predictor variables. The sample length for each predictive variable depends on data availability. Some predictors have data available as far back as 1926, while other variables have samples that start more recently. In our in-sample tests, the length of the time series varies depending on data availability.

For each variable, we run predictive regression models of the form:

$$(6) \quad r_{t:t+h} = \alpha + \beta x_t + \varepsilon_{t:t+h} \text{ for } t = 1, \dots, T - h,$$

where $r_{t:t+h} = (1/h)(r_{t+1} + \dots + r_{t+h})$, r_t is the continuously compounded S&P 500 return for month t from CRSP including dividends and excess of the monthly risk-free rate from Goyal and Welch (2008), x_t is the predictor variable, and h denotes the forecast horizon. We examine four different forecast horizons: 1 month ahead, 1 quarter ahead, one-half year ahead, and 1 year ahead (i.e., $h = 1, 3, 6,$ or 12). For each predictor at each forecast horizon, the regression is estimated using all available data, leading to a single parameter estimate (β) that measures the predictive ability of the candidate variable at that horizon.¹⁷

The results are shown in Figure 1 and Tables 2 and 3. In Figure 1, we provide a visual display of the relation between cross-sectional predictors and their time-series counterparts. For each predictor, the horizontal axis plots the cross-sectional t -statistic, while the vertical axis plots the time-series t -statistic. Graph A displays results for 1-month time-series regressions, while Graphs B–D display results for 3-month, 6-month, and 12-month regressions.¹⁸ If there is a direct mapping from cross-sectional predictability to time-series predictability, then at a minimum, the signs of the two results should be the same. This would imply that all of the results should either be in quadrant I (the top right of each graph) or quadrant III (the bottom left). While a number of results are in quadrants I and III, they represent about 58% of all observations at the 12-month horizon. In other words, there are a number of results in quadrant II (top left) and IV (bottom right) suggesting that many time-series predictors have the opposite sign of their cross-sectional counterpart. In each graph, we also plot a linear trendline; if time-series predictors have the same t -statistics as their cross-sectional counterpart, the trendline should have a 45-degree slope increasing from the left of the figure to the right. In all graphs, the line does slope up from left to right, indicating a relation between time-series and cross-sectional predictors, but the line is flatter than 45°. However, the line does become steeper as the forecasting horizon increases. At the 1-month horizon, the Pearson correlation coefficient between the time-series and cross-sectional predictors is 0.175 while at the 12-month horizon the correlation rises to 0.289.

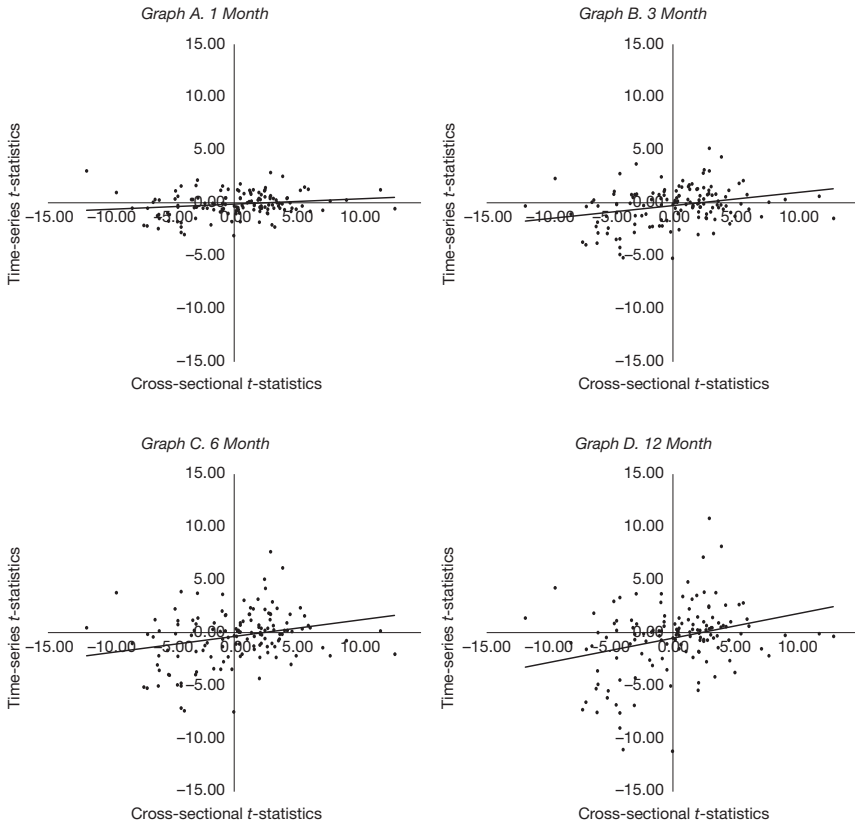
To formally examine the time-series information in cross-sectional predictors, we next turn to the regression results. Table 2 provides a summary of the performance of the candidate predictors, broken out by various subcategories. We report the fraction of predictors that are statistically significant at the 10% level or better

¹⁷In untabulated results, we estimate equation (6) using the weighted least squares method of Johnson (2019). This estimator does not affect our conclusions.

¹⁸We only display results for value-weighted time-series predictors in Figure 1, however, in unreported results the results are similar for equal-weighted predictors.

FIGURE 1
In-Sample Relation Between Cross-Sectional and Time-Series Predictors

Figure 1 displays a plot of the relation between cross-sectional and time-series predictors. For each predictor, we plot the t -statistic from a cross-sectional regression (shown on the horizontal axis) and the t -statistic from value-weighted in-sample time-series regressions (shown on the vertical axis). Graph A examines 1-month time-series regressions, whereas Graphs B–D examine 3-month, 6-month, and 12-month regressions, respectively. In each graph, the diagonal black line denotes a linear trendline.



using 2-sided t -statistics computed using a stationary block bootstrap (Politis and Romano (1994)) with 1,000 draws to account for the Stambaugh (1999) bias and the fact that the model uses overlapping observations when $h > 1$ (Hodrick (1992), Goetzmann and Jorion (1993), and Nelson and Kim (1993)).¹⁹

Panel A of Table 2 displays the results from the four different sets of predictors: i) PREDICTORS_FROM_EXISTING_PAPERS examines only

¹⁹Specifically, we resample the original data using a stationary block bootstrap with a mean block size of 5, however our conclusions are robust to alternate block sizes of 10, 25, and 50. To avoid the overlapping dependent variable issue, we aggregate the independent variable rather than the dependent variable as suggested by Cochrane (1991) and Jegadeesh (1991); our general conclusions – about the lack of statistically significant predictability after accounting for the number of tests run – do not change under any of these alternate approaches.

TABLE 2
Summary of In-Sample Performance Using Unadjusted p -Values

Table 2 displays a count of the number of predictive variables that are statistically significant at the 10% level or better, as a fraction of the total number of variables examined. We calculate statistical significance using bootstrap p -values. For each anomaly, we estimate an in-sample predictive regression of the form:

$$r_{t:t+h} = \alpha + \beta x_t + \varepsilon_{t:t+h} \text{ for } t = 1, \dots, T - h,$$

where $r_{t:t+h} = (1/h)(r_{t+1} + \dots + r_{t+h})$, r_t is the continuously compounded S&P 500 return for month t from CRSP including dividends and excess of the monthly risk-free rate from Goyal and Welch (2008), h indicates the forecast horizon in months, and x_t is one of the 140 predictor variables. To construct time-series predictors out of cross-sectional predictors, we calculate the value-weighted and equal-weighted mean across all stocks on each date resulting in 280 possible predictors. In Panel A, we consider four different definitions: i) PREDICTORS_FROM_EXISTING_PAPERS uses only those variables that are used in the existing literature on time-series return predictability. ii) RAW_PREDICTORS examines every possible variable for which we reject the null that the raw variable is nonstationary. iii) RAW_PREDICTORS + FIRST_DIFF examines every possible variable however if a variable is not stationary in raw form, we then examine whether it is nonstationary in first-differenced form. If we fail to reject the null that the first differenced variable is nonstationary, we drop the variable. iv) ALL_POSSIBLE_PREDICTORS examines every possible variable. If we fail to reject the null that a variable is nonstationary, we calculate deviations from a linear trend model. If we fail to reject the null that the linearly detrended variable is nonstationary, we calculate the first difference. If we fail to reject the null that the first differenced variable is nonstationary, we drop the variable. In Panel B, we examine subsamples of the variables in RAW_PREDICTORS + FIRST_DIFF formed on the 10 most statistically significant cross-sectional predictors (BEST_CROSS-SECTIONAL), and two different groupings based on the categories in McLean and Pontiff (2016): i) OPINION and ii) VALUATION. In Panel C, we examine EQUAL-WEIGHTED vs. VALUE-WEIGHTED_PREDICTORS for the variables in RAW_PREDICTORS + FIRST_DIFF.

Predictive Variable	Return Horizon (h)			
	$h = 1$	$h = 3$	$h = 6$	$h = 12$
<i>Panel A. Candidate Predictors (Number Significant/Total Examined)</i>				
PREDICTORS_FROM_EXISTING_PAPERS	6/51	7/51	6/51	10/51
RAW_PREDICTORS	19/137	21/137	30/137	37/137
RAW_PREDICTORS + FIRST_DIFF	27/253	27/253	36/253	43/253
ALL_POSSIBLE_PREDICTORS	34/269	43/269	45/269	54/269
<i>Panel B. By Subcategory (Number Significant/Total Examined)</i>				
BEST_CROSS-SECTIONAL	1/20	1/20	1/20	2/20
OPINION	4/38	3/38	2/38	3/38
VALUATION	0/24	0/24	2/24	6/24
<i>Panel C. By Aggregation Method (Number Significant/Total Examined)</i>				
EQUAL-WEIGHTED_PREDICTORS	10/125	8/125	14/125	17/125
VALUE-WEIGHTED_PREDICTORS	17/128	19/128	22/128	26/128

predictors that have *already* been examined in the time-series literature; ii) RAW_PREDICTORS expands the list to include *all* predictors that are stationary in raw form (i.e., before applying transformations like the first difference) and it imposes filters to remove variables that should not aggregate and/or do not exhibit time-series variation; iii) RAW_PREDICTORS + FIRST_DIFF is similar to ii) except it adds the first-difference transformation (i.e., if we fail to reject the null that the raw variable is nonstationary, we then calculate the first difference and include it if it passes the unit root test); iv) ALL_POSSIBLE_PREDICTORS examines as many predictors as possible (i.e., we only require that they have time-series variation and pass the unit root test). Across all four sets, we find evidence of return predictability. In our main specification (RAW_PREDICTORS + FIRST_DIFF), on average 13% are statistically significant at the 10% level or better across the various forecasting horizons. This number generally increases as the forecast horizon increases from 27 (11%) at the 1-month horizon to 43 (17%) at the 12-month horizon. The numbers are similar for the other sets of predictors. These results run counter to explanations that the results are hinged on the initial set of variables being over- or under-aggressive.

TABLE 3
Best In-Sample Predictive Regression Results Using Romano and Wolf p -Values

Table 3 reports the ordinary least squares estimate of β , p -values, and the R^2 statistic from in-sample predictive regression models of the form:

$$r_{t:t+h} = \alpha + \beta x_t + e_{t:t+h} \text{ for } t = 1, \dots, T - h,$$

where $r_{t:t+h} = (1/h)(r_{t+1} + \dots + r_{t+h})$, r_t is the continuously compounded S&P 500 return for month t from CRSP including dividends and excess of the monthly risk-free rate from Goyal and Welch (2008), h indicates the forecast horizon in months, and x_t is the predictor variable shown in columns 2 and 9. For each horizon, we run 253 regressions using the variables in the RAW_PREDICTORS + FIRST_DIFF set of predictors. Panel A displays results for the 1-month horizon, Panel B shows the 3-month horizon, Panel C shows the 6-month horizon, and Panel D shows the 12-month horizon. Within each panel, predictors are sorted by their Romano and Wolf p -value and then their unadjusted p -value. We report all predictors that have unadjusted p -values less than 10% for a given horizon. Unadjusted p -values are shown in columns 6 and 13 and Romano and Wolf (2016) adjusted p -values are shown in columns 7 and 14.

1	2	3	4	5	6	7	8	9	10	11	12	13	14
Rank	Predictor	EW or VW	$\hat{\beta}$	(%) R^2	p -Value	Raw RW	Rank	Predictor	EW or VW	$\hat{\beta}$	(%) R^2	p -Value	Raw RW
<i>Panel A. 1 Month Horizon</i>													
1	Z_SCORE	VW	-0.0056	1.6	0.00	0.83	15	AMIHUD'S_MEASURE	EW	0.0026	0.4	0.03	1.00
2	ASSET_TURNOVER	VW	0.0053	1.5	0.00	0.89	16	CASH_FLOW_VARIANCE	EW	0.0034	0.6	0.07	1.00
3	Z_SCORE	EW	-0.0052	1.4	0.00	0.91	17	ΔREC. + ACCRUAL	EW	-0.0046	1.1	0.03	1.00
4	LT_REVERSAL	VW	-0.0046	0.7	0.01	0.92	18	SPREADS	EW	0.0026	0.4	0.06	1.00
5	ΔNC_OP_ASSETS	VW	-0.0036	0.7	0.03	0.98	19	CAPEX_GROWTH	VW	-0.0032	0.5	0.07	1.00
6	STOCK_SPLIT	EW	-0.0026	0.4	0.10	1.00	20	SEO	VW	-0.0042	0.9	0.01	1.00
7	ΔTAX_TO_ASSETS	EW	0.0055	1.6	0.00	1.00	21	ASSET_GROWTH	VW	-0.0035	0.7	0.02	1.00
8	VOLUME/MV	EW	-0.0044	1.0	0.05	1.00	22	MOMENT_LT_REVERSAL	VW	0.0025	0.2	0.06	1.00
9	COSKEWNESS	EW	0.0036	0.7	0.04	1.00	23	SUSTAINABLE_GROWTH	VW	-0.0031	0.5	0.06	1.00
10	SECURE/TOTAL_DEBT	VW	0.0040	0.9	0.03	1.00	24	COSKEWNESS	VW	0.0037	0.8	0.07	1.00
11	SPINOFFS	VW	-0.0040	0.9	0.05	1.00	25	%OPERATING_ACCRUAL	VW	-0.0038	0.9	0.10	1.00
12	PENSION_FUNDING	VW	-0.0035	0.7	0.01	1.00	26	FORECAST_DISPERSION	VW	0.0032	0.6	0.04	1.00
13	IPOS	VW	0.0035	0.7	0.05	1.00	27	ΔREC. + ACCRUAL	VW	-0.0049	1.3	0.05	1.00
14	TARGET_PRICE	EW	-0.0054	1.7	0.05	1.00							
<i>Panel B. 3 Month Horizon</i>													
1	Z_SCORE	VW	-0.0055	4.6	0.00	0.33	15	SPREADS	VW	0.0028	1.2	0.06	0.99
2	ASSET_TURNOVER	VW	0.0056	4.7	0.00	0.33	16	ABN_ANALYST_INTENSE	EW	-0.0026	1.0	0.01	1.00
3	LT_REVERSAL	VW	-0.0048	2.2	0.02	0.41	17	REVERSE_SPLIT	EW	-0.0020	0.6	0.08	1.00
4	Z_SCORE	EW	-0.0046	3.1	0.00	0.69	18	ΔTAX_TO_ASSETS	EW	0.0044	2.9	0.00	1.00
5	SEO	VW	-0.0045	3.1	0.00	0.72	19	VOLUME	VW	-0.0020	0.4	0.09	1.00
6	TARGET_PRICE	EW	-0.0062	6.0	0.01	0.84	20	INVENTORY_GROWTH	VW	-0.0023	0.8	0.08	1.00
7	ΔNC_OP_ASSETS	VW	-0.0033	1.8	0.04	0.84	21	SECURE/TOTAL_DEBT	VW	0.0021	0.7	0.08	1.00
8	ASSET_GROWTH	VW	-0.0036	2.1	0.01	0.85	22	DIVIDEND_OMISSION	EW	0.0026	1.0	0.07	1.00
9	COSKEWNESS	VW	0.0036	2.1	0.04	0.85	23	SPREADS	EW	0.0024	0.9	0.05	1.00
10	%OPERATING_ACCRUAL	VW	-0.0045	3.4	0.05	0.93	24	REVERSE_SPLIT	VW	0.0028	1.3	0.08	1.00
11	SUSTAINABLE_GROWTH	VW	-0.0033	1.7	0.03	0.95	25	ACCRUALS	VW	0.0029	1.3	0.07	1.00
12	CASH_FLOW_VARIANCE	EW	0.0033	1.7	0.06	0.99	26	SHARE_ISSUES_PW	VW	-0.0031	1.4	0.09	1.00
13	CAPEX_GROWTH	VW	-0.0032	1.6	0.06	0.99	27	MOMENT_LT_REVERSAL	VW	0.0022	0.4	0.08	1.00
14	ΔREC. + ACCRUAL	VW	-0.0037	2.1	0.06	0.99							
<i>Panel C. 6 Month Horizon</i>													
1	ASSET_TURNOVER	VW	0.0059	9.7	0.00	0.07	19	VOLUME_TREND	EW	0.0032	3.0	0.08	0.94
2	Z_SCORE	VW	-0.0057	9.1	0.00	0.10	20	SPREADS	EW	0.0026	2.0	0.01	0.94
3	LT_REVERSAL	VW	-0.0049	4.6	0.00	0.13	21	REVERSE_SPLIT	VW	0.0029	2.6	0.01	0.94
4	Z_SCORE	EW	-0.0044	5.2	0.00	0.46	22	MOMENT_REVERSE	VW	-0.0026	1.2	0.07	0.94
5	ASSET_GROWTH	VW	-0.0037	3.9	0.00	0.60	23	ΔSALES_ΔINVENTORY	EW	-0.0026	2.0	0.04	0.95
6	SPREADS	VW	0.0035	3.7	0.00	0.61	24	SHARE_ISSUES_PW	VW	-0.0028	2.1	0.07	0.95
7	ΔNC_OP_ASSETS	VW	-0.0032	3.1	0.03	0.64	25	R&D/MV	EW	0.0026	2.0	0.07	0.96
8	SUSTAINABLE_GROWTH	VW	-0.0036	3.8	0.01	0.64	26	CASH_FLOW_VARIANCE	EW	0.0029	2.6	0.05	0.96
9	%OPERATING_ACCRUAL	VW	-0.0047	6.7	0.02	0.66	27	DIVIDEND_OMISSION	EW	0.0021	1.3	0.09	0.99
10	SEO	VW	-0.0038	4.0	0.01	0.73	28	ORG_CAPITAL	EW	0.0023	1.5	0.07	0.99
11	TARGET_PRICE	EW	-0.0052	7.0	0.01	0.90	29	PROFIT_MARGIN	VW	0.0025	1.7	0.09	0.99
12	ΔTAX_TO_ASSETS	EW	0.0048	6.0	0.00	0.92	30	ABN_ANALYST_INTENSE	EW	-0.0012	0.4	0.05	1.00
13	ΔCAPEX_ΔIND_CAPEX	EW	-0.0029	2.4	0.00	0.92	31	SHARE_ISSUES_DT	EW	-0.0021	1.2	0.05	1.00
14	CAPEX_GROWTH	VW	-0.0031	3.0	0.06	0.92	32	CASH_TO_ASSETS	VW	0.0014	0.6	0.04	1.00
15	ACCRUALS	VW	0.0029	2.4	0.03	0.92	33	ΔTAX_TO_ASSETS	VW	0.0028	2.0	0.09	1.00
16	ΔCAPEX_ΔIND_CAPEX	VW	-0.0029	2.5	0.02	0.92	34	SHARE_ISSUES_DT	VW	-0.0017	0.8	0.06	1.00
17	PRICE	VW	-0.0028	1.4	0.09	0.92	35	ΔSALES_ΔSG&A	EW	0.0019	1.0	0.09	1.00
18	COSKEWNESS	VW	0.0027	2.2	0.06	0.93	36	CASH_FLOW_VARIANCE	VW	0.0020	1.2	0.10	1.00

(continued on next page)

TABLE 3 (continued)
Best In-Sample Predictive Regression Results Using Romano and Wolf p -Values

1	2	3	4	5	6	7	8	9	10	11	12	13	14
Rank	Predictor	EW or VW	$\hat{\beta}$	Raw		Rank	Predictor	EW or VW	$\hat{\beta}$	($\%$) F^2	Raw		Rank
				($\%$) F^2	p -Value						p -Value	p -Value	
<i>Panel D. 12 Month Horizon</i>													
1	Z_SCORE	VW	-0.0058	18.5	0.00	0.02	23	Δ SALES- Δ INVENTORY	EW	-0.0022	2.8	0.03	0.95
2	ASSET_TURNOVER	VW	0.0058	17.8	0.00	0.03	24	MOMENT-REVERSE	VW	-0.0022	1.7	0.06	0.95
3	LT_REVERSAL	VW	-0.0044	7.3	0.00	0.13	25	R&D/MV	EW	0.0022	2.5	0.08	0.97
4	Z_SCORE	EW	-0.0043	9.8	0.00	0.28	26	CASH_FLOW_ VARIANCE	EW	0.0023	3.0	0.05	0.97
5	SUSTAINABLE_ GROWTH	VW	-0.0039	8.3	0.00	0.30	27	REVERSE_SPLIT	VW	0.0022	2.8	0.03	0.97
6	Δ NC_OP_ASSETS	VW	-0.0035	6.8	0.01	0.31	28	EXCHANGE_SWITCH	VW	-0.0019	2.1	0.04	0.97
7	ASSET_GROWTH	VW	-0.0037	7.7	0.00	0.36	29	COSKEWNESS	VW	0.0019	2.1	0.06	0.97
8	SPREADS	VW	0.0035	7.2	0.00	0.37	30	SHARE_ISSUES_DT	EW	-0.0018	1.8	0.03	0.99
9	PRICE	VW	-0.0036	4.1	0.02	0.45	31	E/P	EW	-0.0017	1.7	0.08	0.99
10	%OPERATING_ACCRUAL	VW	-0.0046	12.2	0.01	0.47	32	AMIHUD'S_MEASURE	EW	0.0018	1.8	0.09	0.99
11	CAPEX_GROWTH	VW	-0.0034	7.1	0.02	0.63	33	HYBRID_COVAR_ _RISK	VW	0.0018	2.0	0.10	0.99
12	Δ CAPEX- Δ IND_CAPEX	EW	-0.0030	5.2	0.00	0.70	34	M/B_AND_ACCRUALS	VW	-0.0020	2.3	0.09	0.99
13	SPREADS	EW	0.0027	4.2	0.00	0.77	35	STOCK_SPLIT	EW	-0.0017	1.8	0.08	1.00
14	Δ CAPEX- Δ IND_CAPEX	EW	-0.0028	4.6	0.00	0.77	36	Δ TAX_TO_ASSETS	EW	0.0029	4.0	0.01	1.00
15	ORG_CAPITAL	EW	0.0030	4.9	0.01	0.79	37	COSKEWNESS	EW	0.0010	0.6	0.03	1.00
16	EMPLOYEE_GROWTH	VW	-0.0027	4.3	0.04	0.80	38	E/P	VW	0.0008	0.4	0.04	1.00
17	VOLUME_TREND	EW	0.0030	5.0	0.03	0.86	39	CF/MV	VW	0.0007	0.3	0.06	1.00
18	LT_REVERSAL	EW	-0.0025	2.3	0.05	0.86	40	Δ SALES- Δ SG&A	EW	0.0016	1.5	0.09	1.00
19	SHARE_ISSUES_PW	VW	-0.0027	4.0	0.06	0.87	41	Δ #INSTITUT. _OWNERS	VW	-0.0015	1.3	0.09	1.00
20	Δ ASSET_TURN	VW	0.0028	4.2	0.02	0.87	42	Δ PRICE_FORECAST	VW	-0.0011	0.5	0.10	1.00
21	TARGET_PRICE	EW	-0.0043	8.7	0.03	0.92	43	MOMENT_LT_ REVERSAL	VW	0.0015	0.8	0.03	1.00
22	SEO	VW	-0.0026	3.6	0.08	0.93							

In Panel B of Table 2, we examine three subcategories of the 253 variables in our main specification. The first two subcategories, VALUATION and OPINION, are motivated by economic theory. The VALUATION subcategory is composed of variables that are a function of discount rates, so theory suggests that they should be related to returns in a time-series setting (Lewellen (2004), Kelly and Pruitt (2013)). The OPINION subcategory of predictors consists of variables like institutional trading and analyst upgrades, which can be motivated with the sentiment explanation of Baker and Wurgler (2006) or with the information explanation of Seyhun (1988). Finally, we examine a third subcategory, BEST_CROSS-SECTIONAL, defined as the 10 most statistically significant anomalies in the cross-sectional literature. Several papers find a reduction of cross-sectional return predictability in recent periods (e.g., McLean and Pontiff (2016), Green et al. (2017)). If the cross-sectional predictors we start with are only weakly related to returns in the cross-section, then they may lead to weak performance in our time-series tests. Accordingly, we examine a subcategory that focuses on the best performers in the cross-sectional literature.

When we examine the three subcategories defined above, the results are similar. Interestingly, the results are generally weaker for the VALUATION and OPINION categories than for the entire set of predictors. Since valuation ratios may be a function of discount rates, they are arguably the most likely cross-sectional predictors to work in a time-series setting (Kelly and Pruitt (2013)). Indeed, several different VALUATION predictors have been studied in the existing market risk premium literature, notably the dividend-to-price and earnings-to-price ratios.

Our results suggest these predictors are an exception, rather than the norm. Across all horizons, only 8% of VALUATION predictors and 8% of OPINION predictors are statistically significant, vs. 13% across all predictors.

When we examine the results for the “best” cross-sectional variables, there is some evidence of predictability, but again, the results are weaker than when we examine all possible predictors. For the BEST_CROSS-SECTIONAL predictors (those with the top 10 highest cross-sectional t -statistics), 10% are statistically significant at the 1-year horizon. The results are generally consistent with the idea that the best cross-sectional variables forecast returns because they contain information about the systematic component of returns. As a result, these variables also aggregate to form good time-series predictors. However, the fact that the subset of BEST_CROSS-SECTIONAL predictors is not stronger than the entire set suggests that good cross-sectional predictors are not necessarily good time-series predictors.

Finally, we examine the in-sample results broken out by the different methodologies used to construct the aggregate predictor variables. Specifically, in Panel C of Table 2, we present results for value- or equal-weighted average predictors. The findings are largely consistent across the two methodologies: across all predictors, 13% are statistically significant at the 1-month horizon when value-weighted and 8% are statistically significant at the 1-month horizon when equal-weighted. As the horizon extends, the value-weighted predictors appear to perform slightly better, but the results are generally similar across the two groups. For example, across all predictors, 20% are statistically significant at the annual horizon when value-weighted vs. 14% when equal-weighted. On the surface, before a deeper consideration of multiple hypothesis testing is considered, Table 2 suggests that some cross-sectional anomaly variables predict the market risk premium.

B. Multiple Hypothesis Testing

The results in Table 2 examine as many as 269 different regression models, so the results are subject to concerns about data snooping. Put differently, with 269 different regression models, some tests will likely be statistically significant due to type I errors. Fortunately, a growing literature shows how to adjust p -values to account for the number of models considered.

White (2000) develops a RCB to correct for data-snooping. While a number of approaches exist to adjust p -values for the bias that results from multiple hypothesis testing, the White approach has several desirable properties. First, it uses a bootstrap procedure to estimate the dependence structure of the p -values across all considered models. In contrast, the Bonferroni, Dunn (1961) and Holm (1979) approaches assume the worst-case dependence structure. This causes them to be overly conservative in that they do not reject the null hypothesis enough when the null is false. Because it estimates the actual dependence structure, the White RCB has greater power than the Bonferroni and Holm methods. Second, the White procedure is particularly suited to the application of return predictability regressions because the procedure uses a loss function that compares the performance of a predictor to a benchmark model. Economic theory suggests the equity risk premium should be

positive; as a result, a strategy that simply predicts positive returns all the time would frequently be correct. Accordingly, the return predictability literature often compares the forecast accuracy of a predictive variable to the so-called “prevailing mean” model that uses the prevailing mean return as the forecast of next period’s equity risk premium. The White RCB then uses a loss function that compares the mean squared forecast error (MSFE) for each predictor to the MSFE from a benchmark model that uses the prevailing mean return.

Despite these advantages, the White procedure does have a drawback: if the null is rejected, it indicates that the *best* predictor examined is better than the benchmark, but it does not reveal whether the other predictors are better. Put differently, the White procedure does not test whether the second-best predictor, or the k th best predictor, is better than the benchmark. Accordingly, Romano and Wolf (2005), (2016) develop a step-down procedure that extends the White procedure to calculate whether individual predictors are better than the benchmark.²⁰ The resulting procedure controls the familywise error rate and provides p -values for each individual predictor.

This is the first paper to use this procedure to evaluate the predictability of the market risk premium. The closest related papers are Sullivan, Timmerman, and White (1999) and Chordia, Goyal, and Saretto (2020). Sullivan, Timmerman, and White (1999) apply the White procedure to examine the performance of technical trading strategies in predicting the equity risk premium. Sullivan et al. find no evidence that technical trading strategies generate portfolio performance that outperforms a benchmark. More recently, Chordia et al. (2020) examine the performance of trading strategies in the cross-section and they use several methods to correct for multiple hypothesis testing bias, including the Romano and Wolf (2005), (2016) procedure. They find that most strategies studied in the extant literature are not significant after adjusting for multiple hypothesis testing. Our paper unifies these two literatures by providing the first evidence on the performance of time-series predictors formed from cross-sectional variables.

To further explore the robustness of our findings, in the Supplementary Material we also examine p -values calculated using the Benjamini and Yekutieli (2001) procedure, which controls the false discovery rate while allowing for arbitrary dependence. On average, false discovery rate methods have better power to reject false null hypotheses, but this comes at a cost: they are more likely to reject true null hypotheses. In other words, the Benjamini and Yekutieli (2001) procedure is less conservative than the Romano and Wolf (2016) procedure as it allows for more false positives.²¹ Nonetheless, in all analyses our main conclusions are unchanged when we use the Benjamini and Yekutieli (2001) procedure in place of the Romano and Wolf (2016) procedure.

²⁰See Romano and Wolf (2016) for a detailed discussion of the procedure and see Section C of the Supplementary Material for a detailed overview of our implementation of the procedure. We thank Allan Timmermann for helpful conversations about the White (2000) and Romano and Wolf (2005) procedures.

²¹Specifically, the Romano and Wolf (2016) procedure controls the probability of having *any* false positives while the Benjamini and Yekutieli (2001) procedure controls the *expected proportion* of false positives, so it allows for more false positives as you consider more predictors.

C. In-Sample Results Adjusted for Multiple Hypothesis Testing

Table 3 reports detailed estimates for individual predictors that are statistically significant before adjusting for multiple testing. For brevity, we present results for only our main specification, which examines 253 candidate predictors, however the results are similar for other specifications. The table presents coefficient estimates, R^2 values, and both raw and Romano and Wolf (2016) adjusted p -values. Here we reach the main conclusion of the paper: most of the cross-sectional variables that appear to be statistically significant when examined in isolation are no longer statistically significant when examined in the context of all cross-sectional predictors.

For example, at the 1-month and 3-month horizons, none of the predictors remain significant at the 10% level when computing Romano and Wolf (2016) adjusted p -values compared to 27 when computing individual p -values (Table 2). At the 6-month horizon, there are 2 predictors that are significant at the 10% level (compared to 36 in Table 2) and at the 12-month horizon there are 2 predictors significant at the 10% level (compared to 43 in Table 2). Moreover, the predictors with remarkable statistical significance when examined in isolation, no longer appear to be so stellar when examined among the set of 253 predictors.

In terms of economic significance, a number of predictors are noteworthy. Z_SCORE and $ASSET_TURNOVER$, with R^2 values of 4.6% and 4.7% shown in column 5 of Panel B in Table 2, are among the best predictors at the 3-month horizon. Moreover, several other predictors have R^2 values exceeding 3% including $\%OPERATING_ACCRUAL$ and $TARGET_PRICE$. Although these R^2 values might seem small in absolute magnitude, Zhou (2010) notes that R^2 values from predictive regressions are typically small in absolute magnitude as stock returns are difficult to forecast. Accordingly, Zhou (2010) builds on the insights of Ross (2005) to construct a mathematical bound on the maximum R^2 that can exist under no arbitrage conditions. Using consumption growth rates as a state variable, Zhou's bounds imply a maximum R^2 at the monthly horizon of between 0.079% and 0.177%²² and Huang and Zhou (2017) show that the quarterly R^2 is bound by at most 3.74%, depending on the specification, and in most cases it is less than 3%. In light of this, the return predictability of some Table 3 predictors is economically large.²³

In addition, we note that many of the best predictors at the 3-month horizon are also good predictors at the 6-month and 12-month horizons. At the 12-month horizon, a number of the predictors have impressive R^2 values of 10% or greater

²²The bounds developed in Zhou (2010) depend on the choice of a state variable. We do not take a stance on state variables in this article, we simply note that the R^2 values we document appear economically meaningful relative to the example bounds presented in Zhou (2010).

²³We also construct predictors based on sets of cross-sectional predictors using principal component analysis (PCA). PCA requires nonmissing observations for each predictor; as a result, we are not able to utilize the entire set of predictors until 1999. In untabulated results, the first principal component extracted from all equal-weighted predictors and value-weighted predictors exhibits strong return predictability over the period 1999–2017, before adjusting for multiple testing. However, if we estimate the PCA using the set of variables with data starting in 1980, the evidence becomes mixed; we find weak evidence of return predictability for value-weighted predictors, but no evidence of return predictability for equal-weighted predictors. We thank Bryan Kelly for this suggestion.

including Z_SCORE, ASSET_TURNOVER, and %OPERATING_ACCRUAL. For example, Z_SCORE exhibits a remarkable R^2 of 18.5% at the 12-month horizon. However, despite some impressive results for individual predictors, after applying the Romano and Wolf (2005), (2016) procedure only Z_SCORE and ASSET_TURNOVER are statistically significant.

Table 4 summarizes the results after applying the Romano and Wolf (2005), (2016) procedure. In Panel A, our main specification (RAW_PREDICTORS + FIRST_DIFF) exhibits 0, 0, 2, and 2 statistically significant predictors at the 1- month, 3- month, 6- month, and 12-month horizons. One possible critique is that these results include predictors that should not have been included, and this causes a failure to reject a false null hypothesis (a type II error). To address this concern, we also present results for a variety of different sets of predictors and subsamples. In our main specification we use only the raw or first-differenced version of each variable, yet some existing papers in the time-series return predictability literature examine stochastically detrended variables. Accordingly, in Panel A, we examine all possible predictors which add predictors that require stochastic detrending. Consistent with Rapach et al. (2016) we find that

TABLE 4
Summary of In-Sample Performance Using Romano and Wolf p -Values

Table 4 displays a count of the number of predictive variables that are statistically significant at the 10% level or better, as a fraction of the total number of variables examined. We calculate statistical significance using Romano and Wolf adjusted p -values. For each anomaly, we estimate an in-sample predictive regression of the form:

$$r_{t:t+h} = \alpha + \beta x_t + \varepsilon_{t:t+h} \text{ for } t = 1, \dots, T - h,$$

where $r_{t:t+h} = (1/h)(r_{t+1} + \dots + r_{t+h})$, r_t is the continuously compounded S&P 500 return for month t from CRSP including dividends and excess of the monthly risk-free rate from Goyal and Welch (2008), h indicates the forecast horizon in months, and x_t is one of the 140 predictor variables. To construct time-series predictors out of cross-sectional predictors, we calculate the value-weighted and equal-weighted mean across all stocks on each date resulting in 280 possible predictors. In Panel A, we consider four different definitions: i) PREDICTORS_FROM_EXISTING_PAPERS uses only those variables that are used in the existing literature on time-series return predictability. ii) RAW_PREDICTORS examines every possible variable for which we reject the null that the raw variable is nonstationary. iii) RAW_PREDICTORS + FIRST_DIFF examines every possible variable however if a variable is not stationary in raw form, we then examine whether it is nonstationary in first-differenced form. If we fail to reject the null that the first differenced variable is nonstationary, we drop the variable. iv) ALL_POSSIBLE_PREDICTORS examines every possible variable. If we fail to reject the null that a variable is nonstationary, we calculate deviations from a linear trend model. If we fail to reject the null that the linearly detrended variable is nonstationary, we calculate the first-difference. If we fail to reject the null that the first differenced variable is nonstationary, we drop the variable. In Panel B, we examine subsamples of the variables in RAW_PREDICTORS + FIRST_DIFF formed on the 10 most statistically significant cross-sectional predictors (BEST_CROSS-SECTIONAL), and two different groupings based on the categories in McLean and Pontiff (2016): i) OPINION and ii) VALUATION. In Panel C, we examine EQUAL-WEIGHTED vs. VALUE-WEIGHTED_PREDICTORS for the variables in RAW_PREDICTORS + FIRST_DIFF

Predictive Variable	Return Horizon (h)			
	$h = 1$	$h = 3$	$h = 6$	$h = 12$
<i>Panel A. Candidate Predictors (Number Significant/Total Examined)</i>				
PREDICTORS_FROM_EXISTING_PAPERS	0/51	0/51	1/51	3/51
RAW_PREDICTORS	0/137	0/137	2/137	2/137
RAW_PREDICTORS + FIRST_DIFF	0/253	0/253	2/253	2/253
ALL_POSSIBLE_PREDICTORS	0/269	0/269	1/269	3/269
<i>Panel B. By Subcategory (Number Significant/Total Examined)</i>				
BEST_CROSS-SECTIONAL	0/20	0/20	1/20	1/20
OPINION	0/38	0/38	0/38	0/38
VALUATION	0/24	0/24	0/24	1/24
<i>Panel C. By Aggregation Method (Number Significant/Total Examined)</i>				
EQUAL-WEIGHTED_PREDICTORS	0/125	0/125	0/125	0/125
VALUE-WEIGHTED_PREDICTORS	0/128	0/128	3/128	2/128

detrended short interest is one of the best predictors at the 3-month, 6-month, and 12-month horizons.²⁴ However, after applying the Romano and Wolf procedure, we find that none of the 269 predictors are statistically significant at the 1-month and 3-month horizons and only 3 of the variables are statistically significant at the 12-month horizon.

Similarly, we examine two other sets of predictors (RAW_PREDICTORS and PREDICTORS_FROM_EXISTING_PAPERS) and find little evidence of return predictability. In sum, the three alternate sets of predictors in Panel A of Table 4 provide bounds on the possible set of variables to consider. At a minimum, we know the profession has examined the 51 variables in the PREDICTORS_FROM_EXISTING_PAPERS so these variables have to be included in a multiple testing framework. At a maximum, the 269 predictors in ALL_POSSIBLE_PREDICTORS represent all the cross-sectional variables that could possibly be examined. The fact that our results are unchanged across these two extremes shows that our conclusions are not sensitive to the number of variables considered. Moreover, when we examine the results by subcategory in Panel B, the results again show zero statistically significant predictors at the 1-month horizon and between 0 and 1 statistically significant predictors at the 12-month horizon.

Finally, Table IA.I of the Supplementary Material summarizes the results after applying the Benjamini and Yekutieli (2001) procedure. As expected, there are more statistically significant predictors in Table IA.I of the Supplementary Material than Table 4, but the overall conclusion is similar. Even though the Benjamini and Yekutieli (2001) procedure is more permissive than the Romano and Wolf (2016) procedure, there are relatively few predictors that remain significant. For example, when we examine predictors from existing papers, we find no evidence of predictability at the 1-month and 3-month horizons, and only 3 of the 60 predictors are significant at the 6-month horizon and 4 of the 60 predictors are significant at the 12-month horizon. Overall, across a wide variety of samples and methodologies, our conclusion remains unchanged: there is only weak in-sample evidence that cross-sectional predictors contain information about the systematic portion of returns.

D. Out-of-Sample Tests

A number of papers note that in-sample tests may overstate predictability due to the use of information that was not known *ex ante* (e.g., Cooper, Gutierrez, and Marcum (2005), Goyal and Welch (2008)). Accordingly, in this section, we revisit each of our tests using out-of-sample forecasting regressions. We again run predictive regressions of the form:

$$(7) \quad r_{t:t+h} = \alpha_t + \beta_t x_t + \varepsilon_{t:t+h} \quad \text{for } t = 1, \dots, T - h,$$

where $r_{t:t+h} = (1/h)(r_{t+1} + \dots + r_{t+h})$, r_t is the excess return on the S&P500, and x_t is the predictor variable. However, now we estimate the model separately for each

²⁴Detailed in-sample results for the set of all possible predictors are shown in Table IA.III of the Supplementary Material and out-of-sample results are shown in Table IA.IV of the Supplementary Material.

TABLE 5
Summary of Out-of-Sample Performance Using Unadjusted p -Values

Table 5 displays a count of the number of predictive variables that are statistically significant at the 10% level or better, as a fraction of the total number of variables examined. We calculate statistical significance using bootstrap p -values. For each anomaly, we estimate an out-of-sample predictive regression of the form:

$$r_{t:t+h} = \alpha + \beta x_t + \varepsilon_{t:t+h} \text{ for } t = 1, \dots, T - h,$$

where $r_{t:t+h} = (1/h)(r_{t+1} + \dots + r_{t+h})$, r_t is the continuously compounded S&P 500 return for month t from CRSP including dividends and excess of the monthly risk-free rate from Goyal and Welch (2008), h indicates the forecast horizon in months, and x_t is one of the 140 predictor variables. We estimate expanding rolling window regressions using only information available on each date. To construct time-series predictors out of cross-sectional predictors, we calculate the value-weighted and equal-weighted mean across all stocks on each date resulting in 280 possible predictors. In Panel A, we consider four different definitions: i) PREDICTORS_FROM_EXISTING_PAPERS uses only those variables that are used in the existing literature on time-series return predictability. ii) RAW_PREDICTORS examines every possible variable for which we reject the null that the raw variable is nonstationary. iii) RAW_PREDICTORS + FIRST_DIFF examines every possible variable however if a variable is not stationary in raw form, we then examine whether it is nonstationary in first-differenced form. If we fail to reject the null that the first differenced variable is nonstationary, we drop the variable. iv) ALL_POSSIBLE_PREDICTORS examines every possible variable. If we fail to reject the null that a variable is nonstationary, we calculate deviations from a linear trend model. If we fail to reject the null that the linearly detrended variable is nonstationary, we calculate the first-difference. If we fail to reject the null that the first differenced variable is nonstationary, we drop the variable. In Panel B, we examine subsamples of the variables in RAW_PREDICTORS + FIRST_DIFF formed on the 10 most statistically significant cross-sectional predictors (BEST_CROSS-SECTIONAL), and two different groupings based on the categories in McLean and Pontiff (2016): i) OPINION and ii) VALUATION. In Panel C, we examine EQUAL-WEIGHTED vs. VALUE-WEIGHTED_PREDICTORS for the variables in RAW_PREDICTORS + FIRST_DIFF

Predictive Variable	Return Horizon (h)			
	$h = 1$	$h = 3$	$h = 6$	$h = 12$
<i>Panel A. Candidate Predictors (Number Significant/Total Examined)</i>				
PREDICTORS_FROM_EXISTING_PAPERS	0/51	5/51	9/51	11/51
RAW_PREDICTORS	4/137	24/137	32/137	36/137
RAW_PREDICTORS + FIRST_DIFF	7/253	30/253	39/253	44/253
ALL_POSSIBLE_PREDICTORS	7/269	32/269	41/269	45/269
<i>Panel B. By Subcategory (Number Significant/Total Examined)</i>				
BEST_CROSS-SECTIONAL	0/20	2/20	3/20	3/20
OPINION	1/38	3/38	1/38	0/38
VALUATION	0/24	0/24	0/24	1/24
<i>Panel C. By Aggregation Method (Number Significant/Total Examined)</i>				
EQUAL-WEIGHTED_PREDICTORS	3/125	10/125	16/125	17/125
VALUE-WEIGHTED_PREDICTORS	4/128	20/128	23/128	27/128

time period, using only information that was available at each date. As such, we estimate new parameter estimates for α_t and β_t at each point in time. If the relation between the predictor variable and the equity risk premium is stable over time, then this out-of-sample approach should produce the same results as the in-sample analysis discussed in Section IV.A. If the relation between the predictor variable and the equity risk premium is not stable, the out-of-sample tests may lead to a different conclusion.

As previously discussed, the sample length for each predictive variable depends on data availability. For the out-of-sample tests, we use the first 10 years of data to train the model before we make our first forecast. As before, we start by summarizing the results across all specifications. Table 5 provides a summary of the performance of the candidate predictors. To make inferences, we calculate the out-of-sample R_{OS}^2 statistic defined as in Campbell and Thompson (2008).²⁵ To calculate the out-of-sample R_{OS}^2 , we use the prevailing mean equity risk

²⁵We use the unconstrained out-of-sample R^2 from Campbell and Thompson (2008) (i.e., we do not impose any sign restrictions).

premium, at each date, as our benchmark model and we use the Clark and West (2007) statistic to assess statistical significance.²⁶ Panel A summarizes statistical significance across the four different sets of predictors: i) PREDICTORS_FROM_EXISTING_PAPERS, ii) RAW_PREDICTORS, iii) RAW_PREDICTORS + FIRST_DIFF, and iv) ALL_POSSIBLE_PREDICTORS. Recall that in Table 2, the in-sample evidence was strongest at the annual horizons, however even at the 1-month horizon approximately 11% of the RAW_PREDICTORS + FIRST_DIFF predictors were statistically significant at the 10% level or better. In contrast, the out-of-sample evidence is weaker. At the 1-month horizon, less than 3% of the variables are statistically significant.

In Panel B of Table 5, we examine three subcategories of the 253 variables in our main specification (BEST_CROSS-SECTIONAL, OPINION, and VALUATION) and the results do not look much better. Again, the VALUATION and OPINION subcategories appear to be worse than the full sample of predictors, despite the economic motivation for VALUATION predictors and their popularity in the extant literature. Moreover, while the BEST_CROSS-SECTIONAL category continues to show some evidence of predictability at longer horizons, none of the predictors is significant at the 1-month horizon.

The remaining panels of Table 5 examine the out-of-sample results broken out by the different methodologies used to construct the aggregate predictor variables. Specifically, Panel C examines the results when we calculate the aggregate predictor using a value-weighted average or an equal-weighted average, respectively. Again, the results look similar regardless of the methodology used to construct the predictors. Only 4 of 128 value-weighted predictors are significant at the 1-month horizon and only 3 of 125 equal-weighted predictors are significant. Overall, the results show that cross-sectional variables contain some systematic information at longer horizons, but not at short horizons. Taken together, it is tempting to conclude that some cross-sectional predictors can be used to form aggregate predictors, suggesting they contain information about the systematic component of returns. However, our out-of-sample analyses consider more than 253 different specifications. In the next section, we revisit our results after accounting for the number of hypotheses tested.

E. Multiple Hypothesis Testing

As before, we ask whether our out-of-sample tests show evidence of return predictability after accounting for possible data snooping biases. To do this, we again use the Romano and Wolf (2016) procedure.²⁷ The procedure follows a similar process to the in-sample procedure, except we estimate rolling regressions and compare the

²⁶Formally, we test the null hypothesis that the mean square forecast error (MSFE) from the baseline model is less than or equal to the MSFE from the predictive model vs. the alternative hypothesis that the MSFE from the benchmark model is greater than the MSFE from the predictive model ($H_0: R_{OS}^2 \leq 0$ against $H_A: R_{OS}^2 > 0$).

²⁷The Romano and Wolf procedure expands the White (2000) RCB to adjust p -values for each individual predictor. White (2000) shows this procedure is valid for both in-sample and out-of-sample tests and Sullivan et al. (1999) apply the White (2000) procedure to out-of-sample forecasting regressions.

prediction from these regressions to a rolling average market risk premium. Section B of the Supplementary Material provides a detailed overview of the procedure.

F. Out-of-Sample Results Adjusted for Multiple Hypothesis Testing

Table 6 displays detailed estimates for individual predictors that are statistically significant before adjusting for multiple testing. The table presents the time-series mean of each coefficient estimate, the Campbell and Thompson (2008) out-of-sample R^2 , and both raw and Romano and Wolf (2016) adjusted p -values. Interestingly, the out-of-sample R^2 values in Table 6 highlight many of the same predictors that performed well in Table 3 in the in-sample analyses. Analogous to Table 3, we report the time-series mean of the betas from these regressions in addition to the out-of-sample R^2_{OS} statistic. Indeed, Z_SCORE and ASSET_TURNOVER have out-of-sample R^2 values that are positive and

TABLE 6
Best Out-of-Sample Predictive Regression Results Using Romano and Wolf p -Values

Table 6 reports the mean of the ordinary least squares estimate of β , p -values, and the Campbell and Thompson (2008) R^2_{OS} statistic from out-of-sample predictive regression models of the form:

$$r_{t+h} = \alpha + \beta x_t + \epsilon_{t,t+h} \text{ for } t = 1, \dots, T-h,$$

where $r_{t+h} = (1/h)(r_{t+1} + \dots + r_{t+h})$, r_t is the continuously compounded S&P 500 return for month t from CRSP including dividends and excess of the monthly risk-free rate from Goyal and Welch (2008), h indicates the forecast horizon in months, and x_t is the predictor variable in the first column. $\hat{\beta}$ (column 4) is the time-series mean of the coefficient estimates for each predictor. The Campbell and Thompson R^2_{OS} statistic (columns 5 and 12) is calculated as 1 minus the proportional reduction in mean squared forecast error (MSFE) at the h -month horizon for a predictive regression forecast of the S&P 500 log excess return based on the predictor variable in the first column vis-a-vis the prevailing mean benchmark forecast. For each horizon, we run 253 out-of-sample regressions using the variables in the RAW_PREDICTORS + FIRST_DIFF set of predictors. Panel A displays results for the 1-month horizon, Panel B shows the 3-month horizon, Panel C shows the 6-month horizon, and Panel D shows the 12-month horizon. Within each panel, predictors are sorted by their Romano and Wolf p -value and then their unadjusted p -value. We report all predictors that have unadjusted p -values less than 10% for a given horizon. Unadjusted p -values are shown in columns 6 and 13 and Romano and Wolf (2016) adjusted p -values are shown in columns 7 and 14.

Rank	Predictor	EW or VW	Raw RW				Rank	Predictor	EW or VW	Raw RW			
			$\hat{\beta}$	(%) R^2	p -Value	p -Value				$\hat{\beta}$	(%) R^2	p -Value	p -Value
1	2	3	4	5	6	7	8	9	10	11	12	13	14
<i>Panel A. 1 Month Horizon</i>													
1	ASSET_TURNOVER	VW	0.0070	0.3	0.04	1.00							
2	Δ REC + ACCRUAL	VW	-0.0029	1.1	0.04	1.00							
3	Z_SCORE	EW	-0.0043	0.5	0.04	1.00							
4	Z_SCORE	VW	-0.0049	0.7	0.05	1.00							
5	Δ TAX_TO_ASSETS	EW	0.0065	-0.9	0.06	1.00							
6	PENSION_FUNDING	VW	-0.0043	-0.4	0.08	1.00							
7	COSKEWNESS	EW	0.0045	0.0	0.09	1.00							
<i>Panel B. 3 Month Horizon</i>													
1	IDIO_RISK	VW	-0.0015	1.11	0.00	0.99	16	Z_SCORE	EW	-0.0039	1.43	0.04	1.00
2	VOLUME	EW	-0.0005	0.57	0.00	1.00	17	EXCHANGE_SWITCH	VW	-0.0034	-0.91	0.04	1.00
3	ASSET_TURNOVER	VW	0.0073	3.83	0.00	1.00	18	Δ #INSTITUT_OWNERS	VW	0.0036	-13.43	0.04	1.00
4	ST_REVERSAL	EW	-0.0002	0.39	0.00	1.00	19	CASH_FLOW_VARIANCE	VW	0.0042	-0.03	0.05	1.00
5	ST_REVERSAL	VW	-0.0007	0.30	0.00	1.00	20	LT_REVERSAL	VW	-0.0066	-0.61	0.05	1.00
6	Z_SCORE	VW	-0.0048	3.25	0.01	1.00	21	CASH_FLOW_VARIANCE	EW	0.0020	1.09	0.05	1.00
7	LAG_MOMENT	VW	-0.0013	0.79	0.01	1.00	22	AGE	EW	0.0163	-106.82	0.07	1.00
8	MAX	VW	-0.0014	1.15	0.01	1.00	23	MOMENT_LT_REVERSAL	VW	0.0021	0.34	0.08	1.00
9	ABN_ANALYST_INTENSE	EW	-0.0073	-0.11	0.02	1.00	24	AGE	VW	0.0019	-0.07	0.08	1.00
10	LAG_MOMENT	EW	-0.0010	0.61	0.02	1.00	25	REVERSE_SPLIT	VW	0.0048	-0.10	0.08	1.00
11	Δ NC_OP_ASSETS	VW	-0.0013	1.63	0.02	1.00	26	Δ TAX_TO_ASSETS	EW	0.0051	2.18	0.08	1.00
12	PRICE	VW	-0.0040	0.97	0.02	1.00	27	Δ REC + ACCRUAL	VW	-0.0030	1.45	0.09	1.00
13	VOLUME	VW	-0.0018	0.69	0.03	1.00	28	CAPEX_GROWTH	VW	-0.0024	1.29	0.09	1.00
14	SIZE	EW	0.0000	0.30	0.03	1.00	29	MAX	EW	-0.0001	0.36	0.09	1.00
15	COSKEWNESS	VW	0.0055	1.43	0.03	1.00	30	MOMENT-REVERSE	VW	-0.0009	0.50	0.10	1.00

(continued on next page)

TABLE 6 (continued)
 Best Out-of-Sample Predictive Regression Results Using Romano and Wolf p -Values

Rank	Predictor	EW or VW	Raw RW					Rank	Predictor	EW or VW	Raw RW				
			$\hat{\beta}$	(%) R^2	p -Value		$\hat{\beta}$				(%) R^2	p -Value			
1	2	3	4	5	6	7	8	9	10	11	12	13	14		
<i>Panel C. 6 Month Horizon</i>															
1	ST_REVERSAL	EW	0.0000	0.28	0.00	0.92	21	CASH_FLOW_VARIANCE	EW	0.0018	2.39	0.03	0.86		
2	VOLUME	EW	-0.0003	0.90	0.00	1.00	22	LT_REVERSAL	EW	-0.0028	1.45	0.03	1.00		
3	IDIO_RISK	VW	-0.0020	1.56	0.00	0.98	23	EXCHANGE_SWITCH	VW	-0.0041	-1.94	0.04	1.00		
4	SIZE	VW	0.0004	0.39	0.00	1.00	24	REVERSE_SPLIT	VW	0.0043	-1.32	0.04	1.00		
5	ST_REVERSAL	VW	0.0003	0.34	0.00	1.00	25	SPREADS	VW	0.0080	-2.32	0.04	1.00		
6	ASSET_TURNOVER	VW	0.0075	9.63	0.00	1.00	26	IDIO_RISK	EW	-0.0001	1.31	0.04	1.00		
7	MOMENT_LT_REVERSAL	VW	0.0008	1.18	0.00	1.00	27	MOMENT-REVERSE	EW	-0.0011	1.19	0.04	1.00		
8	VOLUME	VW	-0.0017	0.69	0.00	1.00	28	AGE	VW	0.0019	-0.10	0.05	1.00		
9	MOMENT-REVERSE	VW	-0.0024	2.19	0.00	1.00	29	COSKEWNESS	VW	0.0046	1.42	0.05	1.00		
10	PRICE	VW	-0.0036	1.88	0.01	1.00	30	DIVIDENDS	EW	0.0000	1.01	0.06	1.00		
11	MAX	EW	-0.0014	0.80	0.01	1.00	31	Z_SCORE	EW	-0.0021	1.91	0.06	1.00		
12	MAX	VW	-0.0024	1.46	0.01	1.00	32	ΔCAPEX-ΔIND_CAPEX	VW	-0.0037	0.02	0.06	1.00		
13	SIZE	EW	0.0005	0.23	0.01	1.00	33	PRICE	EW	-0.0023	0.48	0.06	1.00		
14	Z_SCORE	VW	-0.0050	1.37	0.01	0.95	34	DIVIDENDS	VW	0.0000	0.97	0.06	1.00		
15	LAG_MOMENT	VW	-0.0015	0.94	0.01	1.00	35	ΔCAPEX-ΔIND_CAPEX	EW	-0.0038	-0.07	0.06	1.00		
16	LT_REVERSAL	VW	-0.0064	1.15	0.01	1.00	36	CAPEX_GROWTH	VW	-0.0026	3.07	0.07	1.00		
17	SPREADS	EW	0.0120	-1.61	0.02	1.00	37	ΔSALES-ΔINVENTORY	EW	-0.0042	-4.29	0.08	1.00		
18	ΔNC_OP_ASSETS	VW	-0.0009	3.32	0.02	1.00	38	ASSET_GROWTH	VW	0.0003	0.97	0.09	1.00		
19	LAG_MOMENT	EW	-0.0012	1.57	0.02	1.00	39	ABN_ANALYST_INTENSE	EW	-0.0043	-0.66	0.09	1.00		
20	CASH_FLOW_VARIANCE	VW	0.0057	0.26	0.03	1.00									
<i>Panel D. 12 Month Horizon</i>															
1	VOLUME	EW	0.0000	0.25	0.00	0.95	23	MAX	EW	-0.0003	-0.31	0.03	1.00		
2	VOLUME	VW	-0.0002	0.15	0.00	0.99	24	DIVIDENDS	EW	0.0000	1.69	0.03	1.00		
3	LAG_MOMENT	VW	-0.0020	4.36	0.00	1.00	25	DIVIDENDS	VW	0.0000	1.68	0.03	1.00		
4	SIZE	VW	0.0009	-0.19	0.00	1.00	26	Z_SCORE	VW	-0.0053	15.01	0.03	0.91		
5	MOMENT_LT_REVERSAL	EW	0.0009	3.85	0.00	0.98	27	ΔNC_OP_ASSETS	VW	-0.0008	7.35	0.03	1.00		
6	MOMENT_LT_REVERSAL	VW	0.0013	4.00	0.00	0.94	28	AGE	VW	0.0021	-1.17	0.03	1.00		
7	REVERSE_SPLIT	VW	0.0028	-2.50	0.00	1.00	29	EXCHANGE_SWITCH	VW	-0.0040	-2.26	0.04	1.00		
8	LAG_MOMENT	EW	-0.0011	2.90	0.00	1.00	30	MOMENT-REVERSE	EW	-0.0006	1.06	0.05	1.00		
9	PRICE	VW	-0.0047	2.99	0.00	1.00	31	SUSTAINABLE_GROWTH	VW	-0.0021	2.31	0.05	1.00		
10	MOMENT-REVERSE	VW	-0.0020	2.70	0.00	1.00	32	ΔCAPEX-ΔIND_CAPEX	EW	-0.0025	2.46	0.05	1.00		
11	SPREADS	EW	0.0096	1.51	0.00	1.00	33	ΔCAPEX-ΔIND_CAPEX	VW	-0.0030	0.41	0.07	1.00		
12	MAX	VW	-0.0013	1.08	0.00	1.00	34	CASH_FLOW_VARIANCE	VW	0.0062	-0.76	0.07	1.00		
13	ASSET_TURNOVER	VW	0.0079	17.46	0.00	1.00	35	ORG_CAPITAL	EW	0.0027	4.62	0.07	1.00		
14	ST_REVERSAL	VW	0.0011	-0.95	0.01	1.00	36	VOLUME/MV	VW	-0.0012	0.02	0.09	1.00		
15	ST_REVERSAL	EW	0.0010	-1.01	0.01	1.00	37	PRICE	EW	-0.0030	-0.83	0.09	1.00		
16	SPREADS	VW	0.0067	1.58	0.01	1.00	38	VOLUME_TREND	EW	0.0031	2.61	0.09	1.00		
17	SIZE	EW	0.0010	-1.30	0.01	1.00	39	MOMENTUM	EW	0.0007	-1.55	0.09	1.00		
18	LT_REVERSAL	VW	-0.0052	6.09	0.01	1.00	40	ASSET_GROWTH	VW	0.0001	3.07	0.09	1.00		
19	LT_REVERSAL	EW	-0.0025	3.61	0.01	1.00	41	ΔSALES-ΔINVENTORY	EW	-0.0024	-4.57	0.09	1.00		
20	IDIO_RISK	VW	0.0005	-0.84	0.01	1.00	42	COSKEWNESS	VW	0.0039	0.18	0.09	1.00		
21	CAPEX_GROWTH	VW	-0.0034	7.55	0.02	1.00	43	%OPERATING_ACCRUAL	VW	-0.0046	9.75	0.10	1.00		
22	CASH_FLOW_VARIANCE	EW	0.0013	4.44	0.02	0.75	44	IPOS	VW	-0.0018	-1.26	0.10	1.00		

economically meaningful at the 3-month, 6-month, and 12-month horizons. Once again, ASSET_TURNOVER exhibits impressive results with an out-of-sample R^2 value of 17.5% at the 12-month horizon. Because the out-of-sample R^2_{OS} statistic measures the proportional reduction in MSFE that results from using the predictor (relative to the benchmark model), its magnitude is also useful for interpreting the economic significance of these findings. For the variables listed above, the out-of-sample R^2_{OS} statistic suggests economically large return predictability.

TABLE 7
Summary of Out-of-Sample Performance Using Romano and Wolf p -Values

Table 7 displays a count of the number of predictive variables that are statistically significant at the 10% level or better, as a fraction of the total number of variables examined. We calculate statistical significance using Romano and Wolf adjusted p -values. For each anomaly, we estimate an out-of-sample predictive regression of the form:

$$r_{t:t+h} = \alpha + \beta x_t + \varepsilon_{t:t+h} \text{ for } t = 1, \dots, T - h,$$

where $r_{t:t+h} = (1/h)(r_{t+1} + \dots + r_{t+h})$, r_t is the continuously compounded S&P 500 return for month t from CRSP including dividends and excess of the monthly risk-free rate from Goyal and Welch (2008), h indicates the forecast horizon in months, and x_t is one of the 140 predictor variables. We estimate expanding rolling window regressions using only information available on each date. To construct time-series predictors out of cross-sectional predictors, we calculate the value-weighted and equal-weighted mean across all stocks on each date resulting in 280 possible predictors. In Panel A, we consider four different definitions: i) PREDICTORS_FROM_EXISTING_PAPERS uses only those variables that are used in the existing literature on time-series return predictability. ii) RAW_PREDICTORS examines every possible variable for which we reject the null that the raw variable is nonstationary. iii) RAW_PREDICTORS + FIRST_DIFF examines every possible variable however if a variable is not stationary in raw form, we then examine whether it is nonstationary in first-differenced form. If we fail to reject the null that the first differenced variable is nonstationary, we drop the variable. iv) ALL_POSSIBLE_PREDICTORS examines every possible variable. If we fail to reject the null that a variable is nonstationary, we calculate deviations from a linear trend model. If we fail to reject the null that the linearly detrended variable is nonstationary, we calculate the first-difference. If we fail to reject the null that the first differenced variable is nonstationary, we drop the variable. In Panel B, we examine subsamples of the variables in RAW_PREDICTORS + FIRST_DIFF formed on the 10 most statistically significant cross-sectional predictors (BEST_CROSS-SECTIONAL), and two different groupings based on the categories in McLean and Pontiff (2016): i) OPINION and ii) VALUATION. In Panel C, we examine EQUAL-WEIGHTED vs. VALUE-WEIGHTED_PREDICTORS for the variables in RAW_PREDICTORS + FIRST_DIFF

Predictive Variable	Return Horizon (h)			
	$h = 1$	$h = 3$	$h = 6$	$h = 12$
<i>Panel A. Candidate Predictors (Number Significant/Total Examined)</i>				
PREDICTORS_FROM_EXISTING_PAPERS	0/51	0/51	0/51	0/51
RAW_PREDICTORS	0/137	0/137	0/137	0/137
RAW_PREDICTORS + FIRST_DIFF	0/253	0/253	0/253	0/253
ALL_POSSIBLE_PREDICTORS	0/269	0/269	0/269	0/269
<i>Panel B. By Subcategory (Number Significant/Total Examined)</i>				
BEST_CROSS-SECTIONAL	0/20	0/20	0/20	0/20
OPINION	0/38	0/38	0/38	0/38
VALUATION	0/24	0/24	0/24	0/24
<i>Panel C. By Aggregation Method (Number Significant/Total Examined)</i>				
EQUAL-WEIGHTED_PREDICTORS	0/125	0/125	0/125	0/125
VALUE-WEIGHTED_PREDICTORS	0/128	0/128	0/128	0/128

However, many of these conclusions change after adjusting for multiple testing. Table 7 summarizes the results. We fail to reject the null of no predictability at all forecasting horizons for all predictors using the Romano and Wolf (2016) stepdown procedure. In Panel A, the results are consistent across all of the alternate samples (PREDICTORS_FROM_EXISTING_PAPERS, RAW_PREDICTORS, RAW_PREDICTORS + FIRST_DIFF, and ALL_POSSIBLE_PREDICTORS). As a result, our conclusions are not sensitive to the number of variables considered in the Romano and Wolf calculation. In Panels B and C, we again find no predictability. The conclusions are not significantly different when we examine false discovery rates, shown in Table IA.II of the Supplementary Material. At the 1-month and 3-month horizons, none of the predictors are significant in any of the samples. At longer horizons, we find limited evidence of predictability. For example, at the 12-month horizon, we find that 2 out of the predictors from the existing literature remain significant. Overall, the results in this section echo the conclusions of Section IV.B. Many predictors exhibit strong out-of-sample t -statistics. However, once multiple hypothesis testing is considered, the results

are weaker and we find little evidence that cross-sectional predictors contain systematic information.

V. Conclusion

There is a large literature examining the cross-sectional determinants of stock returns. Similarly, many time-series variables have been proposed as predictors of the equity risk premium. While these literatures have evolved largely independently, at least 26 papers have proposed certain cross-sectional variables as candidates for time-series predictability. Using various samples of cross-sectional predictors and accounting for the number of predictors and their interdependence, we examine the link between cross-sectional and time-series predictability. Our analyses provide new information on the nature of return predictability.

We find plenty of evidence that, in isolation, certain cross-sectional variables make great time-series predictors. Some of these variables, like `Z_SCORE` and `ASSET_TURNOVER`, have never been proposed as time-series variables. However, these results largely disappear once we account for the data snooping bias arising from the plethora of predictive variables considered. Moreover, when we examine out-of-sample forecasting regressions, we continue to find little evidence of return predictability.

If each of our 140 cross-sectional predictors were examined by different econometricians, it is likely that several articles would be written discussing the powerful in-sample time-series information in cross-sectional variables. Claims of predictability in these articles would likely be bolstered by out-of-sample Goyal and Welch (2008) tests. Indeed, several such articles exist. In this article, we take a different approach. Once we consider the set of all existing cross-sectional variables documented in the extant literature, the difference in conclusions is stark. The evidence no longer suggests that cross-sectional variables contain information about the systematic component of returns.

Supplementary Material

To view supplementary material for this article, please visit <http://doi.org/10.1017/S0022109022000266>.

References

- Baker, M., and J. Wurgler. "The Equity Share in New Issues and Aggregate Stock Returns." *Journal of Finance*, 55 (2000), 2219–2257.
- Baker, M., and J. Wurgler. "Investor Sentiment and the Cross-Section of Stock Returns." *Journal of Finance*, 61 (2006), 1645–1680.
- Baker, M., and J. Wurgler. "Investor Sentiment in the Stock Market." *Journal of Economic Perspectives*, 21 (2007), 129–152.
- Barberis, N., and R. Thaler. "A Survey of Behavioral Finance." *Handbook of the Economics of Finance*, 1 (2003), 1053–1128.
- Bartsch, V.-S.; H. Dichtl; W. Drobetz; and A. Neuhierl. "Data Snooping in Equity Premium Prediction." *Journal of International Forecasting*, 37 (2021), 72–94.
- Benjamini, Y., and D. Yekutieli. "The Control of the False Discovery Rate in Multiple Testing Under Dependency." *Annals of Statistics*, 29 (2001), 1165–1188.

- Bossaerts, P., and P. Hillion. "Implementing Statistical Criteria to Select Return Forecasting Models: What Do We Learn?" *Review of Financial Studies*, 12 (1999), 405–428.
- Campbell, J. Y. "A Variance Decomposition for Stock Returns." *Economic Journal*, 101 (1991), 157–179.
- Campbell, J. Y., and R. J. Shiller. "The Dividend-Price Ratio and Expectations of Future Dividends and Discount Factors." *Review of Financial Studies*, 1 (1988), 195–228.
- Campbell, J. Y., and S. Thompson. "Predicting Excess Stock Returns Out of Sample: Can Anything Beat the Historical Average?" *Review of Financial Studies*, 21 (2008), 1509–1531.
- Chen, A. Y. "The Limits of p-Hacking: Some Thought Experiments." *Journal of Finance*, 76 (2021), 2447–2480.
- Chen, A. Y., and T. Zimmerman. "Publication Bias and the Cross-Section of Stock Returns." *Review of Asset Pricing Studies*, 10 (2020), 249–289.
- Chordia, T.; A. Goyal; and A. Saretto. "Anomalies and False Rejections." *Review of Financial Studies*, 33 (2020), 2134–2179.
- Chordia, T.; R. Roll; and A. Subrahmanyam. "Order Imbalance, Liquidity and Market Returns." *Journal of Financial Economics*, 65 (2002), 111–130.
- Clark, T. E., and K. D. West. "Approximately Normal Tests for Equal Predictive Accuracy in Nested Models." *Journal of Econometrics*, 138 (2007), 291–311.
- Cochrane, J. H. "Volatility Tests and Efficient Markets: A Review Essay." *Journal of Monetary Economics*, 27 (1991), 463–485.
- Cooper, M., and H. Gulen. "Is Time-Series-Based Predictability Evident in Real Time?" *Journal of Business*, 79 (2006), 1263–1292.
- Cooper, M.; R. Gutierrez; and B. Marcum. "On the Predictability of Stock Returns in Real Time." *Journal of Business*, 78 (2005), 469–500.
- Dickey, D. A., and W. Fuller. "Distribution of the Estimators for Autoregressive Time Series with a Unit Root." *Journal of the American Statistical Association*, 74 (1979), 427–431.
- Dow, C. "Scientific Stock Speculation." *The Magazine of Wall Street* (1920).
- Dunn, O. J. "Multiple Comparisons Among Means." *Journal of the American Statistical Association*, 56 (1961), 52–64.
- Engelberg, J.; R. D. McLean; and J. Pontiff. "Anomalies and News." *Journal of Finance*, 73 (2018), 1971–2001.
- Fama, E. F., and K. R. French. "The Cross-Section of Expected Stock Returns." *Journal of Finance*, 47 (1992), 427–465.
- Fama, E. F., and K. R. French. "A Five-Factor Asset Pricing Model." *Journal of Financial Economics*, 116 (2015), 1–22.
- Foster, F. D.; T. Smith; and R. E. Whaley. "Assessing Goodness-of-Fit of Asset Pricing Models: The Distribution of the Maximal R^2 ." *Journal of Finance*, 52 (1997), 591–607.
- Gibson, T. *The Pitfalls of Speculation*. New York: The Moody Corporation (1906).
- Goetzmann, W., and P. Jorion. "Testing the Predictive Power of Dividend Yields." *Journal of Finance*, 48 (1993), 663–679.
- Goyal, A., and P. Santa-Clara. "Idiosyncratic Risk Matters!" *Journal of Finance*, 58 (2003), 975–1007.
- Goyal, A., and I. Welch. "A Comprehensive Look at the Empirical Performance of Equity Premium Prediction." *Review of Financial Studies*, 21 (2008), 1455–1508.
- Green, J.; J. R. Hand; and X. F. Zhang. "The Characteristics that Provide Independent Information About Average U.S. Monthly Stock Returns." *Review of Financial Studies*, 30 (2017), 4389–4436.
- Harvey, C. R.; Y. Liu; and A. Saretto. "An Evaluation of Alternative Multiple Testing Methods for Finance Applications." *Review of Asset Pricing Studies*, 10 (2020), 199–248.
- Harvey, C. R.; Y. Liu; and H. Zhu. "...and the Cross-Section of Expected Returns." *Review of Financial Studies*, 29 (2016), 5–68.
- Hirshleifer, D.; K. Hou; and S. H. Teoh. "Accruals, Cash Flows, and Aggregate Stock Returns." *Journal of Financial Economics*, 91 (2009), 389–406.
- Hodrick, R. "Dividend Yields and Expected Stock Returns: Alternative Procedures for Inference and Measurement." *Review of Financial Studies*, 5 (1992), 357–386.
- Holm, S. "A Simple Sequentially Rejective Multiple Test Procedure." *Scandinavian Journal of Statistics*, 6 (1979), 65–70.
- Hou, K.; C. Xue; and L. Zhang. "Digesting Anomalies: An Investment Approach." *Review of Financial Studies*, 28 (2015), 650–705.
- Huang, D., and G. Zhou. "Upper Bounds on Return Predictability." *Journal of Financial and Quantitative Analysis*, 52 (2017), 401–425.
- Jacobs, H., and S. Müller. "Anomalies Across the Globe: Once Public, No Longer Existing?" *Journal of Financial Economics*, 135 (2020), 213–230.

- Jegadeesh, N. "Seasonality in Stock Price Mean Reversion: Evidence from the U.S. and the U.K." *Journal of Finance*, 46 (1991), 1427–1444.
- Johnson, T. "A Fresh Look at Return Predictability Using a More Efficient Estimator." *Review of Asset Pricing Studies*, 9 (2019), 1–46.
- Kacperczyk, M.; S. Van Nieuwerburgh; and L. Veldkamp, "A Rational Theory of Mutual Funds' Attention Allocation." *Econometrica*, 84 (2016), 571–626.
- Kelly, B., and S. Pruitt. "Market Expectations in the Cross-Section of Present Values." *Journal of Finance*, 68 (2013), 1721–1756.
- Kothari, S. P.; J. Lewellen; and J. Warner. "Stock Returns, Aggregate Earnings Surprises, and Behavioral Finance." *Journal of Financial Economics*, 79 (2006), 537–568.
- Lewellen, J. "Predicting Returns with Financial Ratios." *Journal of Financial Economics*, 74 (2004), 209–235.
- Linnainmaa, J., and M. Roberts. "The History of the Cross-Section of Stock Returns." *Review of Financial Studies*, 31 (2018), 2606–2649.
- Lintner, J. "The Valuation of Risky Assets and the Selection of Risky Investments in Stock Portfolios and Capital Budgets." *Review of Economics and Statistics*, 47 (1965), 13–37.
- McLean, R. D., and J. Pontiff. "Does Academic Publication Destroy Stock Return Predictability?" *Journal of Finance*, 71 (2016), 5–32.
- Neely, C.; D. Rapach; J. Tu; and G. Zhou. "Forecasting the Equity Risk Premium: The Role of Technical Indicators." *Management Science*, 60 (2014), 1772–1791.
- Nelson, C., and M. Kim. "Predictable Stock Returns: The Role of Small Sample Bias." *Journal of Finance* 48 (1993), 641–661.
- Politis, D., and J. Romano. "The Stationary Bootstrap." *Journal of the American Statistical Association*, 89 (1994), 1303–1313.
- Pontiff, J. "Costly Arbitrage: Evidence from Closed-End Funds." *Quarterly Journal of Economics*, 111 (1996), 1135–1151.
- Pontiff, J. "Costly Arbitrage and the Myth of Idiosyncratic Risk." *Journal of Accounting and Economics*, 42 (2006), 35–52.
- Pontiff, J., and L. Schall. "Book-to-Market as a Predictor of Market Returns." *Journal of Financial Economics*, 49 (1998), 141–160.
- Rapach, D. E.; M. C. Ringgenberg; and G. Zhou. "Aggregate Short Interest and Return Predictability." *Journal of Financial Economics*, 121 (2016), 46–65.
- Romano, J.; A. Shaikh; and M. Wolf. "Formalized Data Snooping Based on Generalized Error Rates." *Econometric Theory*, 24 (2008), 404–447.
- Romano, J. P., and M. Wolf. "Stepwise Multiple Testing as Formalized Data Snooping." *Econometrica*, 73 (2005), 1237–1282.
- Romano, J., and M. Wolf. "Efficient Computation of Adjusted p -Values for Resampling-Based Step-down Multiple Testing." *Statistics and Probability Letters*, 113 (2016), 38–40.
- Ross, S. *Neoclassical Finance*. Princeton: Princeton University Press (2005).
- Seyhun, H. N. "The Information Content of Aggregate Insider Trading." *Journal of Business*, 61 (1988), 1–24.
- Sharpe, W. "Capital Asset Prices: A Theory of Market Equilibrium Under Conditions of Risk." *Journal of Finance*, 19 (1964), 425–442.
- Sloan, R. "Do Stock Prices Fully Reflect Information in Accruals and Cash Flows About Future Earnings?" *Accounting Review*, 71 (1996), 289–315.
- Stambaugh, R. "Predictive Regressions." *Journal of Financial Economics*, 54 (1999), 375–421.
- Sullivan, R.; A. Timmermann; and H. White. "Data-Snooping, Technical Trade Rule Performance, and the Bootstrap." *Journal of Finance*, 54 (1999), 1647–1691.
- Welch, I. "Reproducing, Extending, Updating, Replicating, Reexamining, and Reconciling." *Critical Finance Review*, 8 (2019), 301–304.
- Wen, Q. "Asset Growth and Stock Market Returns: A Time-Series Analysis" *Review of Finance*, 23 (2019), 599–628.
- White, H. "A Reality Check for Data Snooping." *Econometrica*, 68 (2000), 1097–1126.
- Zhou, G. "How Much Stock Return Predictability Can We Expect from an Asset Pricing Model?" *Economic Letters*, 108 (2010), 184–186.