

Do GDP Forecasts Respond Efficiently to Changes in Interest Rates?

By DEAN CROUSHORE AND KATHERINE MARSTEN*

June 15, 2015

In this paper, we examine and extend the results of Ball and Croushore (2003) and Rudebusch and Williams (2009), who show that the output forecasts in the Survey of Professional Forecasters (SPF) are inefficient. Ball and Croushore show that the SPF output forecasts are inefficient with respect to changes in monetary policy, as measured by changes in real interest rates, while Rudebusch and Williams show that the forecasts are inefficient with respect to the yield spread. In this paper, we investigate the robustness of both claims of inefficiency, using real-time data and exploring the impact of alternative sample periods on the results. Keywords: real-time data, output forecasts, yield spread, monetary policy

Forecasters have incentives to create the best forecasts they can. Studies that examine the efficiency of forecasts from surveys (Survey of Professional Forecasters or Livingston Survey), such as Pesaran and Weale (2006) and Croushore (2010), suggest that such aggregate forecasts are reasonably accurate, but not completely efficient. Two recent papers, Ball and Croushore (2003) and Rudebusch and Williams (2009), show that the output forecasts in the Survey of Professional Forecasters (SPF) are inefficient. Ball and Croushore show that the SPF output forecasts are inefficient with respect to changes in monetary policy

* Croushore: University of Richmond, Robins School of Business, dcrousho@richmond.edu. Marsten: Federal Reserve Board of Governors, katherine.marsten@frb.gov. We thank the University of Richmond for summer research support.

(as measured by changes in real interest rates), while Rudebusch and Williams show that the forecasts are inefficient with respect to the yield spread. In this paper, we investigate the robustness of both claims of inefficiency, using real-time data and exploring the impact of alternative sample periods on the results.

I. Research on Forecast Inefficiency

Ball and Croushore show that changes in monetary policy lead professional forecasters to modify their output forecasts, but they do not do so efficiently. The forecasters change their output forecasts in response to a change in monetary policy, but not by a sufficient amount. (Ball and Croushore find that the forecasters do respond efficiently in their forecasts for inflation, however.) To illustrate this result, Ball and Croushore compare the SPF output growth forecasts both to actual output growth and to a forecast from a univariate forecasting model, which assumes that output growth follows an AR(1) process with a mean shift in 1973:Q2. They then examine the forecast errors for both the time-series model and the SPF, as well as the differences between the time-series model and the SPF forecasts.

To measure monetary policy, Ball and Croushore use the change over the previous year in the real federal funds rate ($FF1$), which is defined as the nominal federal funds rate minus the expected inflation rate over the next year. Because of lags in the effect of monetary policy, they also look at an additional one-year lag of that measure ($FF2$). The forecasters in the SPF would know the values of $FF1$ and $FF2$ at the time they make their forecasts of output growth for the coming year.

Ball and Croushore find that $FF1$ is correlated with the time-series forecast errors, which suggests that monetary policy has an impact on output because it moves output in a manner that differs from what the time-series model suggests. They also find that $FF1$ is correlated with the differences between the time-series forecast and the SPF, which suggests that the SPF participants use information

that is different from just a univariate time-series model. However, the SPF forecast errors turn out to be negatively correlated with $FF1$, suggesting that an increase in the real federal funds rate over the past year leads output growth to fall by more than the SPF participants believe it will. Thus, the $FF1$ measure could be used to improve upon the forecasts of the SPF. Ball and Croushore examine the robustness of their results to a more general lag structure, to potential regime shifts, to including output shocks in the regression, and to using changes in the nominal federal funds rate instead of the real federal funds rate as a measure of monetary policy. Their results are robust to all of these variations.

Rudebusch and Williams focus mainly on the use of the yield spread to predict recessions, comparing a yield spread probit model to the SPF probability forecasts of a recession, and their results are robust, as Croushore and Marsten (2014) show. In the last section of their paper, however, Rudebusch and Williams also examine the forecast errors of SPF real output forecasts and their correlation with the lagged yield spread. They find that the SPF real output growth forecast errors are negatively correlated with the yield spread at many horizons (all except current quarter forecasts), suggesting that the SPF forecasters do not efficiently use information from the yield curve. To some extent, this may be due to the early years of the SPF because their evidence is weaker for a sample that is restricted to data after 1987.

II. Data

The forecast data that are the center of both Ball and Croushore (2003) and Rudebusch and Williams (2009) come from the Survey of Professional Forecasters, which began in 1968 as a joint effort by the American Statistical Association and the National Bureau of Economic Research (Zarnowitz and Braun (1993)) and was called the Economic Outlook Survey. The Federal Reserve Bank of Philadelphia took over the survey in 1990 and renamed it the Survey of Professional Forecasters (Croushore (1993)). Participants must be professional forecasters ac-

tively engaged in the business who are capable of making quarterly forecasts of numerous macroeconomic variables.

Ball and Croushore examine the SPF forecast of the average growth rate of real output over the coming year. For example, in the 1991Q4 survey, the forecasters provided real output forecasts for each quarter from 1991Q4 to 1992Q4. Ball and Croushore examine the forecasted growth rate over that four-quarter period. More generally, they examine the one-year-ahead forecasts, y_t^e , defined as

$$(1) \quad y_t^e = \left(\frac{Y_{t+4}^e}{Y_t^e} - 1 \right) \times 100\%,$$

where Y_t^e is the level of the output forecast at date t . They compare the forecast with the actual growth rate over the same period, which is

$$(2) \quad y_t = \left(\frac{Y_{t+4}}{Y_t} - 1 \right) \times 100\%,$$

where Y_t is the level of actual output at date t , where the definition of “actual” is discussed below.

Figure 1 plots both y_t^e and y_t from 1968Q4 to 2013Q2. Note that the forecast is smoother than the actual, which is a property of optimal forecasts.

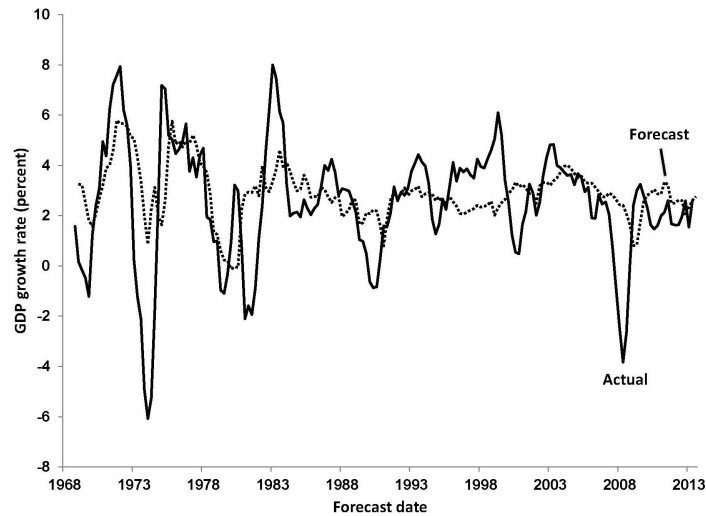
The forecast error for Ball and Croushore is

$$(3) \quad e_t = y_t - y_t^e.$$

Ball and Croushore run the regression given by this equation¹:

¹Note that the regressions do not contain constant terms. We also tested all of the regressions used in this paper to see if a constant term was ever statistically significant and it was not.

FIGURE 1. DATA ON GDP FORECAST AND ACTUAL, BALL-CROUSHORE



Note: The figure shows the forecast for the average growth rate of GDP over the next four quarters with the forecast date shown on the horizontal axis, along with the actual value of GDP growth over those same four quarters.

$$(4) \quad e_t = \beta_1 FF1_t + \beta_2 FF2_t + \epsilon_t,$$

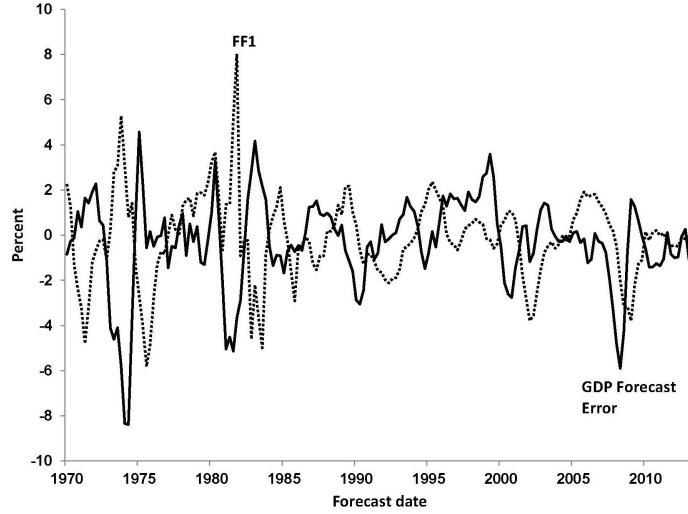
where each regression is run twice: once with β_2 set equal to zero and once with both coefficients allowed to be non-zero.

Figure 2 plots the forecast error and the FF1 variable used in this regression. A negative correlation is readily apparent in the data.

Rudebusch and Williams test the SPF forecast for output growth at various horizons. For a forecast made at date t , the forecasters have information on output at date $t - 1$ and make forecasts for horizons (h) of 0, 1, 2, 3, and 4 quarters, defined as

$$(5) \quad y_{t+h|t-1}^e = \left(\left(\frac{Y_{t+h|t-1}^e}{Y_{t+h-1|t-1}^e} \right)^4 - 1 \right) \times 100\%,$$

FIGURE 2. GDP FORECAST ERROR AND MEASURE OF MONETARY POLICY, BALL-CROUSHORE



Note: The figure shows the GDP forecast error over the next four quarters with the forecast date shown on the horizontal axis, along with the *FF1* measure of monetary policy known to the forecasters at that date.

where $h = 0, 1, 2, 3, 4$, and $Y_{t+h|t-1}^e$ is the level of the output forecast made at date t for date $t + h$, using data on output through date $t - 1$. Rudebusch and Williams test those forecasts against actual values, which are calculated as

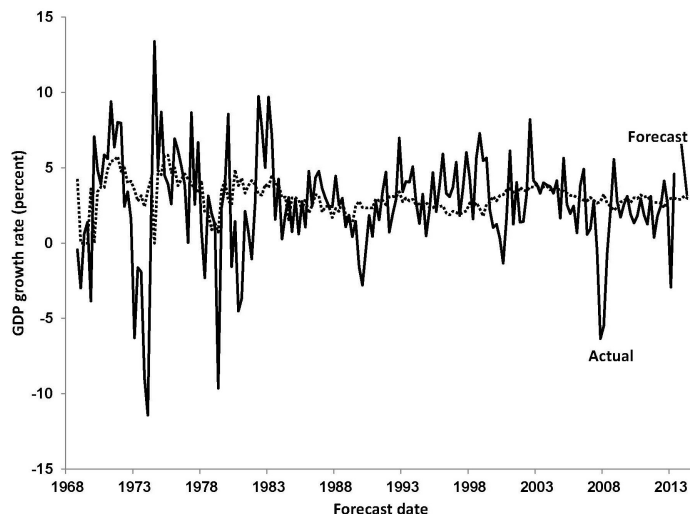
$$(6) \quad y_{t+h} = \left(\left(\frac{Y_{t+h}}{Y_{t+h-1}} \right)^4 - 1 \right) \times 100\%.$$

Figure 3 plots both $y_{t+4|t-1}^e$ and y_{t+4} from 1968Q4 to 2013Q2. Note that the four-quarter-ahead forecast is very smooth, especially since the start of the Great Moderation in the early 1980s.

The h -horizon forecast error is defined as

$$(7) \quad e_{t+h|t-1} = y_{t+h} - y_{t+h|t-1}^e.$$

FIGURE 3. DATA ON FOUR-QUARTER-AHEAD GDP FORECAST AND ACTUAL, RUDEBUSCH-WILLIAMS



Note: The figure shows the four-quarter-ahead forecast for the growth rate of GDP with the forecast date shown on the horizontal axis, along with the actual value of GDP growth in that quarter.

Rudebusch and Williams run the regression

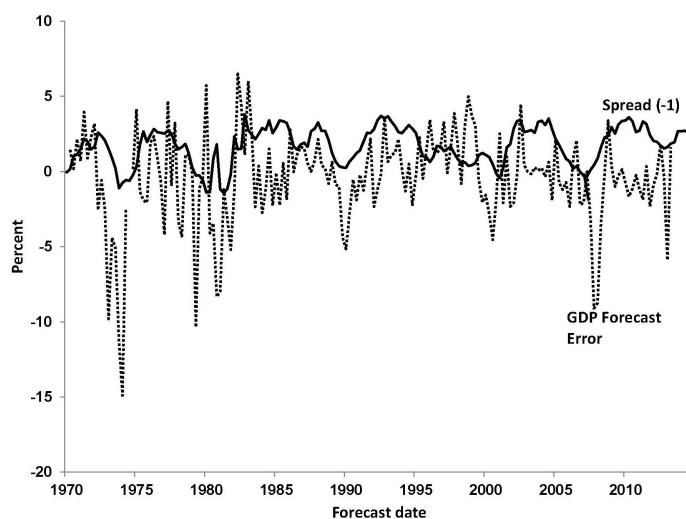
$$(8) \quad e_{t+h|t-1} = \alpha + \beta y_{t+h|t-1}^e + \gamma S_{t-1} + \epsilon_{t+h|t-1},$$

where S_{t-1} is the average yield spread in the quarter before the survey forecast made at time t . They run the regression three times for every time horizon: once with β and γ both equal to zero; once with γ set equal to zero; and once with all three coefficients allowed to be non-zero.

Figure 4 plots the four-quarter-ahead forecast error and the spread variable used in this regression. The relationship between the two measures is not clear because of the volatility of the forecast error.

To evaluate the forecasts, a researcher must know the data that were available to the SPF forecasters at the time they made their forecasts. For this purpose, we use the Real-Time Data Set for Macroeconomists (RTDSM), which was created by Croushore and Stark (2001) and made available on the web site of the Federal Reserve Bank of Philadelphia. The RTDSM provides information on real output

FIGURE 4. FOUR-QUARTER-AHEAD GDP FORECAST ERROR AND SPREAD, RUDEBUSCH-WILLIAMS



Note: The figure shows the four-quarter-ahead GDP forecast error with the forecast date shown on the horizontal axis, along with the spread variable known to the forecasters at that date.

(GNP before 1992, GDP since 1992) and other major macroeconomic variables, as someone standing at the middle of any month from November 1965 to today would have viewed the data. We call the information available to someone at a date in time a “vintage,” and the data are known as “real-time data.”

Ball and Croushore used a primitive type of real-time data to measure the “actual” value of output in their paper. The RTDSM did not yet exist, so they created a set of “first-final data,” which are the data on output released by the government at the end of the third month of the subsequent quarter. For example, for output in the first quarter of the year, the government produces the initial release of the data at the end of April, revises it at the end of May, and releases the first-final data at the end of June. So, all of Ball and Croushore’s results were based on using this first-final data release for each period. Rudebusch and Williams’s research also is based on the first-final data, but they also tested the robustness of their results to using the initial release of the data each quarter, as well as using the last vintage of data available to them in 2007.

We will use the real-time concepts that the two previous studies used, as well as some additional real-time concepts, in our analysis. As it turns out, much additional information about real output becomes available in the government’s annual revision of the data, which is usually released at the end of July each year. So, we will examine the robustness of the results to the use of the annual release of the data, which is actually a much more accurate measure of real output than the initial or first-final releases.

Rudebusch and Williams determine the yield spread as the difference between the interest rates on 10-year U.S. government Treasury bonds and the interest rate on three-month Treasury bills. We use the data on these two series from the FRED database maintained by the Federal Reserve Bank of St. Louis, using quarterly averages. Interest rates are not revised, so there is no need for the use of real-time data for the yield spread.

The yield spread can change for a variety of reasons, but one of the most important reasons is because of changes in monetary policy. A change in monetary policy causes current and expected future short-term interest rates to change, thus affecting the yield spread. The yield spread that Rudebusch and Williams use is strongly correlated with the $FF1$ variable used by Ball and Croushore. The simple contemporaneous correlation coefficient between the two measures from 1970 to 2014 is -0.54 . Thus, similar results should be obtained using the methods from either research method, though the results are likely to be quantitatively different because of the difference in interest-rate measures.

III. Replication Results

We begin our study by replicating the results of Ball and Croushore (2003). Their main result comes from a regression of the one-year-ahead forecast error $y_t - y_t^e$ on variable $FF1$ in one regression, and on variables $FF1$ and $FF2$ in a second regression. We use (for now) the first-final data for y_t , as in Ball and Croushore. We use the mean forecast across SPF participants for our measure

of y_t^e , which is slightly different than Ball and Croushore, who used the median forecast. Theory suggests that the mean forecast is the most appropriate concept to use for aggregating forecasts, as suggested by the literature on forecast encompassing and statistical theory; see Gneiting (2011). The results of the replication are shown in Table 1.

TABLE 1—BALL-CROUSHORE RESULTS AND REPLICATION

Regression: $e_t = \beta_1 FF1_t + \beta_2 FF2_t + \epsilon_t$

Original		
<i>FF1</i>	-0.464 (0.143)	-0.466 (0.155)
<i>FF2</i>		-0.138 (0.085)
χ^2 sig.	< 0.01	< 0.01
\bar{R}^2	0.20	0.21
Replication		
<i>FF1</i>	-0.489 (0.146)	-0.492 (0.158)
<i>FF2</i>		-0.137 (0.089)
χ^2 sig.	< 0.01	< 0.01
\bar{R}^2	0.21	0.22

Note: The table shows the original results reported in Ball and Croushore (2003) and our replication using mean SPF forecast data instead of median forecast data. Numbers in parentheses are HAC standard errors to adjust for overlapping observations. Bold numbers indicate coefficients that are statistically significantly different from zero at the 5 percent level. The sample is 1970Q1 to 1995Q2.

As Table 1 shows, we replicate the Ball and Croushore results almost exactly, which suggests that the differences between the mean and median SPF forecasts are very small indeed. This fact can also be seen by plotting the data, which we do not show here to conserve space. Our replication confirms the Ball and Croushore result that the real output growth forecast error, using first-final data as actuals, is negatively correlated with the lagged change in the real federal funds rate, *FF1*. So, the forecasters should have been able to make better output growth forecasts by using the data on monetary policy changes more effectively.

We perform the same type of replication exercise for the study by Rudebusch

and Williams. We regress the h -horizon forecast error for output growth (where $h = 0, 1, 2, 3,$ and 4) on a constant, the SPF forecast of output growth, and the yield spread from the quarter before the quarter in which the SPF forecast is made, so that it reflects the information available to the forecasters at the time they made their forecasts. In Table 2, we compare our results to those reported by Rudebusch and Williams. Bold numbers indicate coefficients that are significantly different from zero at the 5 percent level. HAC standard errors are used to account for overlapping observations, but are not shown to conserve space. The full sample begins in 1968Q4 and the post-1987 sample begins in 1988Q1; both samples end in 2007Q1.²

Though we are not able to replicate the Rudebusch and Williams results exactly, we broadly confirm their finding that there is significant in-sample fit for the yield spread in explaining SPF output forecast errors. However, the coefficient on the yield spread has the opposite sign of that found by Rudebusch and Williams. We find a positive coefficient, which means that a larger spread is correlated with a larger forecast error. So the forecasters did not increase their GDP forecasts enough as the spread increased (that is, as long-term interest rates rose relative to short-term interest rates). For every forecast horizon except for the current-quarter forecast, we find a significant coefficient for the yield spread. However, most of the F-tests for the overall regression have p-values that do not reject the null hypothesis that all the coefficients are zero. To confirm this finding, we also run the same type of exercise, leaving out the SPF forecast term, to see if the yield spread alone is significant. This is also important because we would expect the SPF forecast to be correlated with the yield spread, so these two terms on the right-hand side of the original regression equation are not likely to be independent. The results of dropping the SPF forecast from the regression are shown in Table 3.

²Because of missing observations, the full-sample four-quarter-ahead forecast sample begins in 1970Q2.

TABLE 2—RUDEBUSCH–WILLIAMS RESULTS AND REPLICATION

Regression: $\epsilon_{t+h t-1} = \alpha + \beta y_{t+h t-1}^e + \gamma S_{t-1} + \epsilon_{t+h t-1}$						
	Full sample	Post-1987 sample				
Original						
		Current-quarter forecast				
Constant	0.31	0.06	-0.04	0.47	0.67	0.43
SPF forecast		0.10	0.08		-0.08	-0.10
Yield spread			-0.10			0.16
F -test (p -value)	0.11	0.18	0.34	0.01	0.56	0.51
		One-quarter-ahead forecast				
Constant	-0.01	-0.16	-0.52	0.25	0.41	0.21
SPF forecast		0.05	-0.19		-0.06	-0.14
Yield spread			-0.65			-0.23
F -test (p -value)	0.96	0.66	0.01	0.32	0.80	0.48
		Two-quarter-ahead forecast				
Constant	-0.24	0.21	-0.14	0.19	0.82	0.84
SPF forecast		-0.15	-0.50		-0.24	-0.43
Yield spread			-0.88			-0.28
F -test (p -value)	0.47	0.40	0.00	0.53	0.43	0.34
		Three-quarter-ahead forecast				
Constant	-0.50	-0.13	-0.70	0.08	1.40	1.47
SPF forecast		-0.11	-0.31		-0.47	-0.71
Yield spread			-0.76			-0.32
F -test (p -value)	0.18	0.63	0.02	0.81	0.19	0.16
		Four-quarter-ahead forecast				
Constant	-0.41	0.54	-0.33	0.06	1.55	1.24
SPF forecast		-0.29	-0.37		-0.53	-0.69
Yield spread			-0.68			-0.43
F -test (p -value)	0.29	0.19	0.00	0.87	0.13	0.05
Replication						
		Current-quarter forecast				
Constant	0.36	0.03	-0.08	0.49	0.72	0.46
SPF forecast		0.14	0.11		-0.09	-0.11
Yield spread			0.11			0.17
F -test (p -value)	0.05	0.02	0.06	0.01	0.06	0.01
		One-quarter-ahead forecast				
Constant	0.08	-0.02	-0.32	0.29	0.50	0.37
SPF forecast		0.04	-0.21		-0.08	-0.15
Yield spread			0.63			0.18
F -test (p -value)	0.77	0.88	0.18	0.29	0.58	0.13
		Two-quarter-ahead forecast				
Constant	-0.13	-0.20	-0.38	0.22	0.29	0.32
SPF forecast		0.02	-0.27		-0.03	-0.15
Yield spread			0.66			0.17
F -test (p -value)	0.71	0.93	0.25	0.53	0.75	0.10
		Three-quarter-ahead forecast				
Constant	-0.34	0.51	0.24	0.10	1.73	1.89
SPF forecast		-0.27	-0.55		-0.58	-0.81
Yield spread			0.72			0.29
F -test (p -value)	0.39	0.55	0.06	0.80	0.56	0.13
		Four-quarter-ahead forecast				
Constant	-0.49	-0.16	-1.17	0.09	1.89	1.89
SPF forecast		-0.10	-0.20		-0.64	-0.89
Yield spread			0.82			0.40
F -test (p -value)	0.26	0.51	0.01	0.81	0.45	<0.01

Note: The table shows the original results reported in Rudebusch and Williams (2009) and our replication. Bold numbers indicate coefficients that are statistically significantly different from zero at the 5 percent level. HAC standard errors are used to account for overlapping observations, but are not shown to conserve space. The full sample begins in 1968Q4 (except for the four-quarter-ahead horizon case, when the sample begins in 1970Q2) and the post-1987 sample begins in 1988Q1; both samples end in 2007Q1.

TABLE 3—RUDEBUSCH–WILLIAMS RESULTS WITH YIELD SPREAD ALONE

Regression: $e_{t+h|t-1} = \alpha + \gamma S_{t-1} + \epsilon_{t+h|t-1}$

	Full sample	Post-1987 sample
	Current-quarter forecast	
Constant	0.05	0.22
Yield spread	0.20	0.15
F -test (p -value)	0.04	<0.01
	One-quarter-ahead forecast	
Constant	-0.63	0.03
Yield spread	0.46	0.15
F -test (p -value)	0.16	0.07
	Two-quarter-ahead forecast	
Constant	-0.93	-0.02
Yield spread	0.51	0.14
F -test (p -value)	0.19	0.23
	Three-quarter-ahead forecast	
Constant	-1.24	-0.11
Yield spread	0.58	0.12
F -test (p -value)	0.05	0.56
	Four-quarter-ahead forecast	
Constant	-1.76	-0.36
Yield spread	0.81	0.26
F -test (p -value)	<0.01	0.15

Note: The table shows the rationality test used by Rudebusch and Williams (2009) but with only the yield spread used in the regression. Bold numbers indicate coefficients that are statistically significantly different from zero at the 5 percent level. HAC standard errors are used to account for overlapping observations, but are not shown to conserve space. The full sample begins in 1968Q4 and the post-1987 sample begins in 1988Q1; both samples end in 2007Q1.

The results show that there is some evidence that the yield spread itself is significant in explaining real GDP forecast errors, at least in the full sample. In Table 3, we see that the coefficient on the yield spread is significantly different from zero for the three- and four-quarter horizons for the full sample period. However, the evidence suggests that the relationship has changed and that since 1987, the forecasters no longer respond inefficiently to changes in the yield spread.

IV. Robustness Exercises

A. Extending the Sample

We now have 18 years of additional data since Ball and Croushore ran their research and eight years of additional data since Rudebusch and Williams ran theirs. In addition, the past seven years include some very difficult times for forecasters, with a deep recession followed by a very weak recovery. So, we update the data to see if the same results hold true for both studies.

Table 4 shows the results of extending the Ball and Croushore results from Table 1 to 2013Q2. In addition, the table shows an additional permutation, with the results of using the nominal federal funds rate instead of the real federal funds rate in measuring monetary policy. This provides an additional test of the robustness of the results.

TABLE 4—BALL-CROUSHORE RESULTS AND REPLICATION, EXTENDED SAMPLE

Regression: $e_t = \beta_1 FF1_t + \beta_2 FF2_t + \epsilon_t$				
Interest rate	Real		Nominal	
$FF1$	-0.372 (0.142)	-0.367 (0.151)	-0.333 (0.131)	-0.324 (0.138)
$FF2$		-0.110 (0.081)		-0.096 (0.066)
χ^2 sig.	< 0.01	0.02	0.01	0.01
\bar{R}^2	0.10	0.11	0.12	0.12

Note: The table shows the extension to 2013Q2 of the results reported in Table 1, using both the real federal funds rate and the nominal federal funds rate to measure changes in monetary policy. Bold numbers indicate coefficients that are statistically significantly different from zero at the 5 percent level.

The results show that the original Ball and Croushore results hold up reasonably well to changing the ending date of the sample from 1995Q2 to 2013Q2, which is an additional 18 years of data. The coefficients on monetary policy become slightly smaller, but the $FF1$ term remains significant in all regressions. Using the nominal interest rate instead of the real interest rate in measuring changes in monetary policy does not matter very much, suggesting that the results are

robust to the choice of proxy variable for measuring monetary policy.

Table 5 shows the results of extending the Rudebusch and Williams method from Table 3 to 2013Q2. The results confirm those for the shorter sample: There is evidence that the SPF forecasts are not efficient in the full sample that begins in 1968, but not for the modern sample that begins in 1987 and ends in 2013. The results suggest that in the early years of the sample, the SPF forecasters were not efficient in using the information about the yield spread in forecasting output growth, but they became more efficient in doing so, beginning in the late 1980s.

TABLE 5—RUDEBUSCH–WILLIAMS RESULTS WITH YIELD SPREAD ALONE, EXTENDED SAMPLE

Regression: $e_{t+h t-1} = \alpha + \gamma S_{t-1} + \epsilon_{t+h t-1}$		
	Full sample	Post-1987 sample
Current-quarter forecast		
Constant	0.10	0.34
Yield spread	0.12	0.00
F -test (p -value)	0.11	0.13
One-quarter-ahead forecast		
Constant	-0.61	0.00
Yield spread	0.36	0.05
F -test (p -value)	0.30	0.75
Two-quarter-ahead forecast		
Constant	-0.99	-0.30
Yield spread	0.45	0.13
F -test (p -value)	0.27	0.86
Three-quarter-ahead forecast		
Constant	- 1.37	-0.57
Yield spread	0.54	0.17
F -test (p -value)	0.08	0.78
Four-quarter-ahead forecast		
Constant	- 1.90	-0.89
Yield spread	0.77	0.34
F -test (p -value)	0.01	0.62

Note: The table shows the rationality test used by Rudebusch and Williams (2009) but with only the yield spread used in the regression. Bold numbers indicate coefficients that are statistically significantly different from zero at the 5 percent level. The sample ends in 2013Q2.

B. *Alternative Starting Dates*

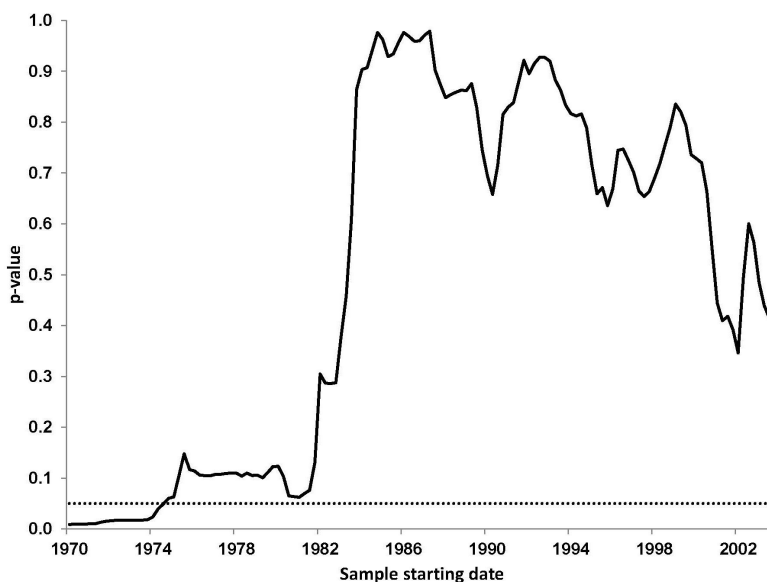
Because of the evidence in the preceding section that suggests a change in the efficiency of the SPF output forecasts beginning in the early 1980s, we look for additional evidence in support of that view by changing the starting dates of the samples. The idea is not to make a formal statistical test (because we will run many different tests that are not independent) but rather to look for evidence that something fundamental has changed in either the way the forecasters use data on monetary policy or the yield spread, or that the predictive power of the variables has changed over time. This method was used by Croushore (2012) to explain why researchers find differing outcomes of bias tests, depending on their choice of which survey of forecasts to examine because each different survey began at a different date. We will summarize the results by examining how the key p -values reported in Tables 4 and 5 change over time as we change the starting dates of the sample.

Ball and Croushore began their sample using the survey forecast in 1968Q4. Using the real federal funds rate as a measure of monetary policy, with just the $FF1$ term, which is the first column in Table 4, and changing the sample starting date, we obtain a set of p -values from 1968Q4 to 2004Q4. The results are shown in Figure 5.

The results suggest that only samples beginning in the late 1960s feature inefficiency. Had the SPF begun later, we would find no evidence of inefficiency. In particular, in the early 1980s, the p -values jump up dramatically, suggesting a change in the behavior of monetary policy, its impact on output, or the SPF at that time.

Doing a similar exercise for the Rudebusch and Williams experiment, we find mixed results, as shown in Figure 6. For starting dates early in the sample, the yield spread is significant and leads to the rejection of the test for forecast efficiency. As we start the sample in the 2000s, however, it is a significant constant term, rather than the yield spread, that leads to the rejection of efficiency. The

FIGURE 5. ALTERNATIVE SAMPLE STARTING DATES, BALL-CROUSHORE



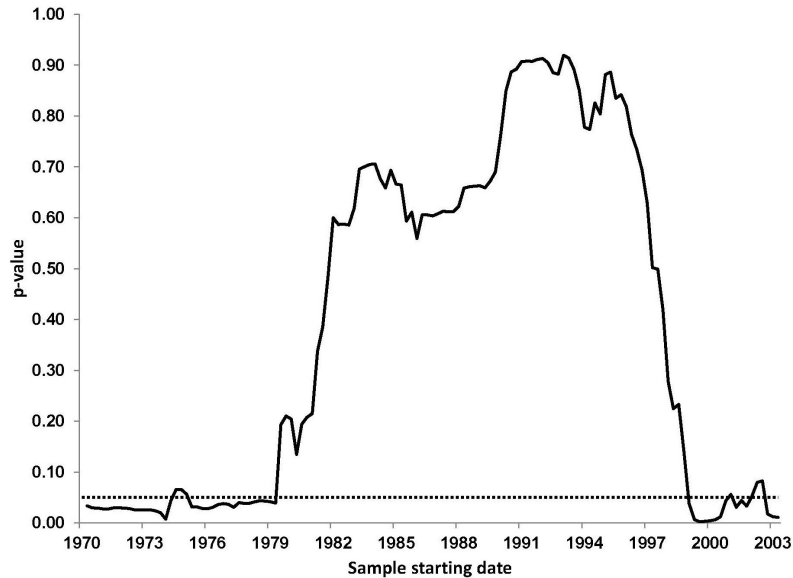
Note: The figure shows the p -value of the test for forecast efficiency, as in the first column of Table 4, but with different sample starting dates. The p -values are not overall tests because of multiple testing, but show what a researcher who started the sample at the date shown on the horizontal axis would have found for a sample that ends in 2013Q2.

significant constant term suggests that the forecasts were biased in that period, perhaps due to the deep recession associated with the financial crisis that began in 2008.

C. Alternative Ending Dates

Given that the starting date for the evaluation of forecast efficiency matters, what if we change the end date of the sample, instead? The idea of this experiment is to imagine researchers at different points in time using the SPF to evaluate forecast efficiency. If the results are the same, regardless of the ending date of the sample, we could argue that the finding of inefficiency is robust. But if the results of the test change over time, then perhaps the outcome is not robust but is special to a particular sample period, which would make it difficult to generalize from the results.

FIGURE 6. ALTERNATIVE SAMPLE STARTING DATES, RUDEBUSCH-WILLIAMS



Note: The figure shows the p -value of the test for forecast efficiency, as in the first column of Table 5 for the four-quarter horizon, but with different sample starting dates. The p -values are not overall tests because of multiple testing, but show what a researcher who started the sample at the date shown on the horizontal axis would have found for a sample that ends in 2013Q2.

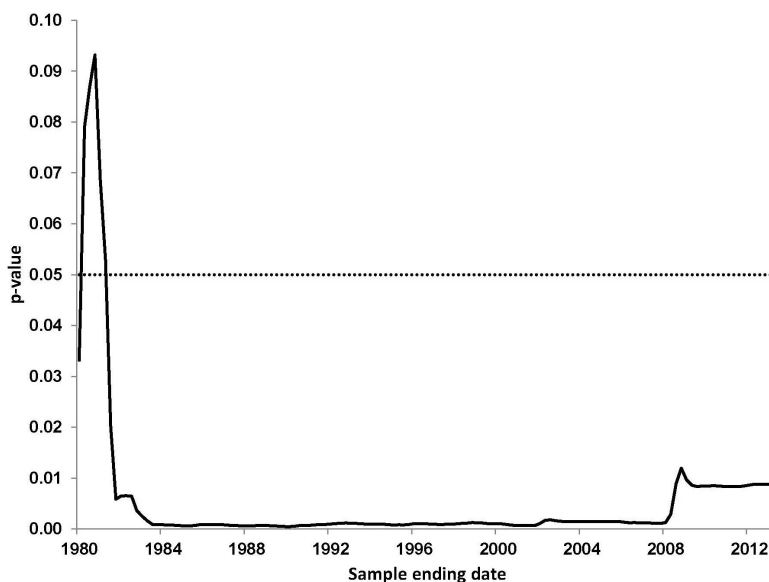
We begin in Figure 7 by rerunning the Ball and Croushore empirical work for a sample that begins in 1970Q1 and ends at various dates from 1980Q1 to 2012Q3. The results show that for almost every sample ending date, we reject the null hypothesis of forecast efficiency, so the result is quite general.

Engaging in the same exercise for the Rudebusch and Williams paper, we find again that the results are consistent with fairly robust conclusions against the efficiency of the SPF forecasts, because for almost every sample ending date, there is evidence of inefficiency (Figure 8).

D. Alternative Actuals

Does the choice of variable being used as “actual” matter? In everything we have done so far, we have used the so-called first-final data as the actual value that is used in assessing forecast accuracy. But are the results sensitive to that

FIGURE 7. ALTERNATIVE SAMPLE ENDING DATES, BALL-CROUSHORE



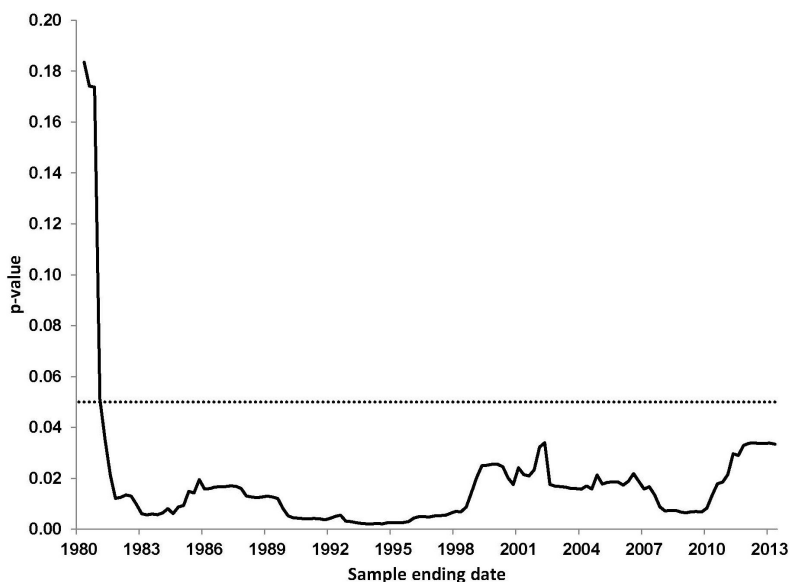
Note: The figure shows the p -value of the test for forecast efficiency, as in the first column of Table 4, but with different sample ending dates. The p -values are not overall tests because of multiple testing, but show what a researcher who started the sample in 1970Q1 and ending at the date shown on the horizontal axis would have found.

choice? To investigate, we run the same analysis that we showed in Table 4 and Table 5, comparing the results using first-final data (under the heading First), annual data, and pre-benchmark data (Pre-bench). The annual data are those coming from the first annual revision, which is usually released at the end of July of the following year, while the pre-benchmark data come from the last release of the data prior to a benchmark revision in which major methodological changes are made.

For the Ball and Croushore results, shown in Table 6, we see that the results are quite robust to the choice of actual variable. The coefficient on the $FF1$ term hardly changes at all, as is also true of the significance level and \overline{R}^2 statistic. So, the results of Ball and Croushore hold up very well to alternative choices of measuring the actual value of GDP growth.

The Rudebusch and Williams results are slightly more sensitive than the Ball

FIGURE 8. ALTERNATIVE SAMPLE ENDING DATES, RUDEBUSCH-WILLIAMS



Note: The figure shows the p -value of the test for forecast efficiency, as in the first column of Table 5 for the four-quarter horizon, but with different sample ending dates. The p -values are not overall tests because of multiple testing, but show what a researcher who started the sample in 1970Q2 and ending at the date shown on the horizontal axis would have found.

and Croushore results to the choice of actual measure of GDP, but not dramatically so. For either the annual concept of actual or the prebenchmark concept, the results are nearly identical to the results using the first-final concept, as Table 7 shows. The only change in the significance of any result is that with the annual concept of actual, the yield spread is significantly different from zero for the two-quarter-ahead forecast. The overall significance of the regression does change slightly for some horizons, especially for the three-quarter-ahead forecast. But, in general, the conclusions do not change, in that the only significant results occur for the full sample, and no variables are ever significant in the post-1987 sample.

E. Rolling Windows

Given that the results appear to depend on the exact sample period chosen, especially for the Rudebusch and Williams results, a useful technique to explore

TABLE 6—BALL-CROUSHORE RESULTS, ALTERNATIVE ACTUALS

Regression: $e_t = \beta FF1_t + \epsilon_t$			
Interest rate	Real		
	First	Annual	Pre-bench
$FF1$	-0.372	-0.367	-0.357
	(0.142)	(0.141)	(0.155)
χ^2 sig.	< 0.01	0.01	0.02
\bar{R}^2	0.10	0.09	0.08
Interest rate	Nominal		
	First	Annual	Pre-bench
$FF1$	-0.333	-0.323	-0.326
	(0.131)	(0.125)	(0.138)
χ^2 sig.	0.01	0.01	0.02
\bar{R}^2	0.12	0.10	0.10

Note: The table shows the results of using three alternative measures of actual GDP growth: First, which is the first final (used earlier), Annual, the annual revision, and Pre-bench, the prebenchmark revision. As in Table 1, we compare the results using both the real federal funds rate and the nominal federal funds rate to measure changes in monetary policy. Bold numbers indicate coefficients that are statistically significantly different from zero at the 5 percent level.

the sensitivity of the results is to examine rolling windows of forecasts to see if there are periods when the forecasts were inefficient and other periods when the forecasts were efficient. We chose to examine ten-year rolling windows of forecasts and examine the p -value for forecast efficiency in the regression equation. For the Ball and Croushore case, we are examining the first column of Table 4, but with rolling 10-year samples. Figure 9 shows the results.

The results of the rolling-window exercise are interesting and somewhat surprising. For ten-year samples that begin before about 1982, there is evidence of inefficiency, but for samples that begin after that, there is no evidence of inefficiency for a single ten-year window. The results suggest that something has changed in the nature of the forecasts in the early 1980s.

Performing the same type of analysis with the Rudebusch and Williams study, shown in Figure 10, we find similar results. After samples that begin in the early 1980s, there are only a few ten-year windows for which the p -value falls below 0.05. This again suggests that something about the forecasts changed in the early

TABLE 7—RUDEBUSCH–WILLIAMS RESULTS WITH YIELD SPREAD ALONE, ALTERNATIVE ACTUALS

Regression: $e_{t+h|t-1} = \alpha + \gamma S_{t-1} + \epsilon_{t+h|t-1}$

	Full sample			Post-1987 sample		
	First	Annual	Pre-bench	First	Annual	Pre-bench
	Current-quarter forecast					
Constant	0.10	-0.03	-0.07	0.34	-0.10	0.00
Yield spread	0.12	0.19	0.18	0.00	0.19	0.10
<i>F</i> -test (<i>p</i> -value)	0.11	0.08	0.13	0.16	0.08	0.47
	One-quarter-ahead forecast					
Constant	-0.61	-0.70	-0.76	0.00	-0.41	-0.28
Yield spread	0.36	0.41	0.42	0.05	0.21	0.10
<i>F</i> -test (<i>p</i> -value)	0.30	0.12	0.21	0.75	0.31	0.89
	Two-quarter-ahead forecast					
Constant	-0.99	-1.15	-1.19	-0.30	-0.66	-0.55
Yield spread	0.45	0.54	0.53	0.13	0.26	0.16
<i>F</i> -test (<i>p</i> -value)	0.27	0.08	0.17	0.86	0.38	0.73
	Three-quarter-ahead forecast					
Constant	-1.37	-1.55	-1.54	-0.57	-0.88	-0.75
Yield spread	0.54	0.66	0.63	0.17	0.30	0.20
<i>F</i> -test (<i>p</i> -value)	0.08	0.01	0.05	0.78	0.64	0.55
	Four-quarter-ahead forecast					
Constant	-1.90	-1.76	-1.72	-0.89	-1.12	-1.07
Yield spread	0.77	0.75	0.69	0.34	0.41	0.35
<i>F</i> -test (<i>p</i> -value)	0.01	0.00	0.00	0.62	0.52	0.48

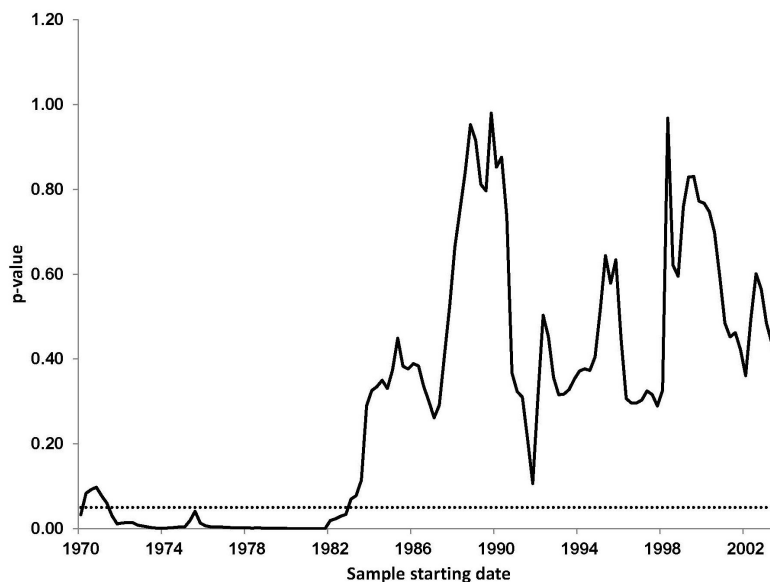
Note: The table shows the rationality test used by Rudebusch and Williams (2009) but with only the yield spread used in the regression, using three alternative measures of actual GDP growth: First, which is the first final (used earlier), Annual, the annual revision, and Pre-bench, the prebenchmark revision. Bold numbers indicate coefficients that are statistically significantly different from zero at the 5 percent level. The sample ends in 2013Q2.

1980s.

V. Did the nature of the SPF change in 1981Q3?

The results of the starting-date and rolling-window analysis suggest that the SPF forecast errors before the early 1980s are the key source of the inefficiency results for both studies. The *p*-values jump up sharply for sample periods that begin after 1981. As it happens, the SPF survey was revamped substantially in the third quarter of 1981, as Zarnowitz and Braun (1993) discuss. Zarnowitz and Braun show that after this change, the number of occasional forecasters dropped and the proportion of truly “professional” forecasters increased. They

FIGURE 9. ROLLING WINDOWS, BALL-CROUSHORE

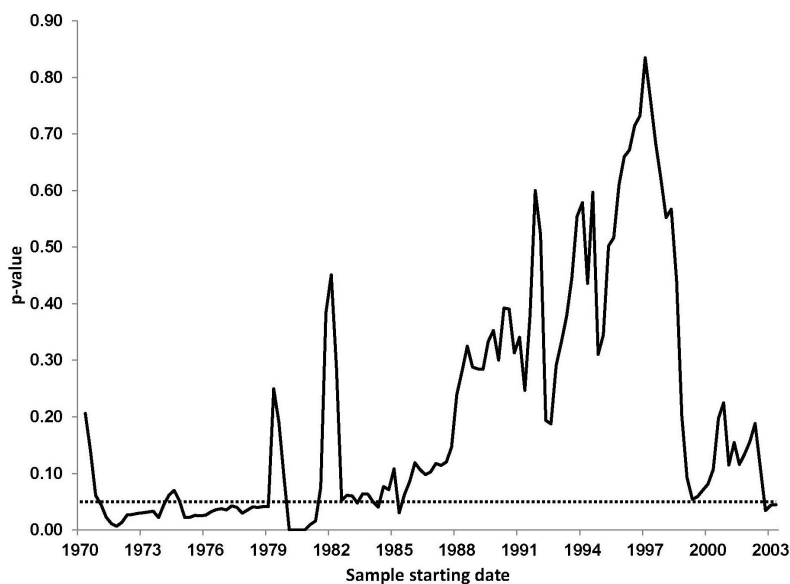


Note: The figure shows the p -value of the test for forecast efficiency, as in the first column of Table 4, but with rolling 10-year samples, beginning with the date shown on the horizontal axis. The p -values are not overall tests because of multiple testing, but show what a researcher who started the sample at the date shown on the horizontal axis would have found for a sample that ends ten years later.

note that “A large number of individuals participated in the earliest surveys but many were not prepared to fill out a detailed questionnaire each quarter and soon dropped out.” (p. 13) This result may be because the early surveys were sent to the entire Business and Economic Statistics section of the American Statistical Association, rather than just to professional forecasters. Certainly, since the Philadelphia Fed took over the survey in 1991, a participant must undergo a detailed screen to ensure that the participant is indeed a professional forecaster capable of providing detailed quarterly forecasts of many major macroeconomic variables. And a forecaster can be removed from the survey if he or she proves to be not up to the task of rigorous and detailed forecasting.

Indeed, if we rerun the empirical results for Table 4 using only the SPF sample beginning in 1981Q3, we find that the inefficiency found over the full sample is somewhat lower, as Table 8 shows. Comparing the results in Table 4 to those in

FIGURE 10. ROLLING WINDOWS, RUDEBUSCH-WILLIAMS



Note: The figure shows the p -value of the test for forecast efficiency, as in the first column of Table 5, but with rolling 10-year samples, beginning with the date shown on the horizontal axis. The p -values are not overall tests because of multiple testing, but show what a researcher who started the sample at the date shown on the horizontal axis would have found for a sample that ends ten years later.

Table 8, we see lower coefficient estimates, higher p -values, and lower \bar{R}^2 statistics. However, there is still evidence that the forecasters are not using information on monetary policy efficiently, especially for using changes in the nominal federal funds rate as a measure of monetary policy.

TABLE 8—BALL-CROUSHORE RESULTS AND REPLICATION, SAMPLE BEGINNING 1981Q3

Regression: $e_t = \beta_1 FF1_t + \beta_2 FF2_t + \epsilon_t$

Interest rate	Real		Nominal	
$FF1$	-0.275 (0.155)	-0.275 (0.162)	-0.250 (0.123)	-0.244 (0.134)
$FF2$		-0.086 (0.108)		-0.117 (0.081)
χ^2 sig.	0.08	0.03	0.04	<0.01
\bar{R}^2	0.07	0.07	0.09	0.10

Note: The table shows the results reported in Table 4, but starting the sample in 1981Q3. Bold numbers indicate coefficients that are statistically significantly different from zero at the 5 percent level.

Performing the same exercise for the Rudebusch and Williams study, we find in Table 9 that the SPF sample beginning in 1981Q3 shows no evidence of the inefficiency that we found for the sample beginning in 1968Q4. We had found the same lack of inefficiency in the forecasts for the sample that began in 1988; these results suggest that the change in the SPF in 1981Q3 may have had a lot to do with the decline in inefficiency found in our extension of the Rudebusch and Williams results.

TABLE 9—RUDEBUSCH–WILLIAMS RESULTS WITH YIELD SPREAD ALONE, SAMPLE BEGINNING 1981Q3

Regression: $e_{t+h t-1} = \alpha + \gamma S_{t-1} + \epsilon_{t+h t-1}$		
	1968Q4 to 2013Q2	1981Q3 to 2013Q2
	Current-quarter forecast	
Constant	0.10	0.45
Yield spread	0.12	-0.07
F -test (p -value)	0.11	0.15
	One-quarter-ahead forecast	
Constant	-0.61	-0.41
Yield spread	0.36	0.23
F -test (p -value)	0.30	0.57
	Two-quarter-ahead forecast	
Constant	-0.99	-0.61
Yield spread	0.45	0.27
F -test (p -value)	0.27	0.60
	Three-quarter-ahead forecast	
Constant	-1.37	-0.58
Yield spread	0.54	0.20
F -test (p -value)	0.08	0.76
	Four-quarter-ahead forecast	
Constant	-1.90	-0.89
Yield spread	0.77	0.39
F -test (p -value)	0.01	0.38

Note: The table shows the rationality test from Table 5, comparing the results from starting the sample in 1968Q4 to the results from starting the sample in 1981Q3. Bold numbers indicate coefficients that are statistically significantly different from zero at the 5 percent level. The sample ends in 2013Q2.

VI. Forecast-Improvement Exercises

The true test of inefficiency in forecasting is to see if inefficiencies in a forecast can be improved in real time. Such forecast-improvement exercises are rare in the literature but very convincing, as in Faust, Rogers and Wright (2005). The idea is to simulate a real-time forecast, using only real-time data and without peeking at future data, to show how one could use regressions that show inefficiency period-by-period, modifying the SPF forecast to make a better forecast. It is an out-of-sample exercise that suggests how one could exploit the forecast inefficiency. Because the Ball and Croushore results showed signs of inefficiency even after the change in the SPF in 1981Q3, we might be able to exploit such inefficiency in real time. But given that the yield spread in the Rudebusch and Williams exercise is not significant after 1981Q3, it seems unlikely to be exploitable in real time, but we will nonetheless run the forecast-improvement exercise to see what might be possible.

To run a forecast-improvement exercise, we take the regression results from the tests for inefficiency and apply them in real-time to modify the SPF survey forecast. For example, in the Ball and Croushore study, the main regression with significant coefficients was:

$$(9) \quad e_t = \alpha FF1_t + \epsilon_t.$$

Taking the estimated $\hat{\alpha}$ and using it in the regression, and recalling that $e_t = y_t - y_t^e$, we create, at each date t , an improved forecast y_t^f , where

$$(10) \quad y_t^f = y_t^e + \hat{\alpha} FF1_t.$$

Thus we are using the regression to modify the original SPF forecast to try to

make a new forecast that is closer to the realized value, based on the past relationship between forecast errors and the most recently observed change in monetary policy. We need a number of years of SPF forecasts before we can reasonably engage in this exercise, so that the sampling error is not too large. So, we will start the forecast-improvement exercise in 1980Q1, using the forecast errors (based on first-final actuals) from 1970Q1 to 1978Q4 to estimate Equation 9, then use the estimated coefficient to improve on the four-quarter-ahead SPF forecast made in 1980Q1. Then, step forward one quarter, use the forecast-error data from surveys from 1970Q1 to 1979Q1, re-estimate Equation 9, then use the estimated coefficient to improve on the four-quarter-ahead SPF forecast made in 1980Q2. Keep repeating this process through to 2013Q2. At the end of the process, compare the root-mean-squared forecast error (RMSFE) for these improved forecasts to the RMSFE for the original SPF forecasts to see if indeed the forecasts have been improved upon or not. We can do the same type of exercise using five-year and ten-year rolling windows to see if that works better.

The results of the forecast-improvement exercises for the Ball and Croushore analysis are shown in Table 10. The results show that with rolling windows, there is some ability to improve on the SPF survey results, but it is never very large (at most a reduction in the RMSE of 7.6 percent) and is never statistically significant.

We do the same exercise for the Rudebusch and Williams study. In this case, we base our analysis on the equation

$$(11) \quad e_{t+h|t-1} = \alpha + \gamma S_{t-1} + \epsilon_{t+h|t-1}.$$

We try to improve the SPF forecast in an analogous fashion to that in the Ball and Croushore study:

$$(12) \quad y_{t+h|t-1}^f = y_{t+h|t-1}^e + \hat{\alpha} + \hat{\gamma} S_{t-1}.$$

TABLE 10—BALL–CROUSHORE FORECAST IMPROVEMENT EXERCISE

Interest rate	Real	Nominal
Expanding Window		
RMSE of Survey	1.762	1.762
RMSE of Improvement to Survey	1.769	1.775
RMSE Ratio	1.004	1.008
<i>p</i> -value	0.96	0.93
Ten-Year Rolling Window		
RMSE of Survey	1.762	1.762
RMSE of Improvement to Survey	1.656	1.686
RMSE Ratio	0.940	0.957
<i>p</i> -value	0.38	0.52
Five-Year Rolling Window		
RMSE of Survey	1.762	1.762
RMSE of Improvement to Survey	1.627	1.640
RMSE Ratio	0.924	0.931
<i>p</i> -value	0.12	0.20

Note: The table shows results of the forecast-improvement exercise. The root-mean-squared error (RMSE) of both the survey and the improvement to the survey are reported in the first two rows of results. The RMSE ratio is the RMSE of the improvement to the survey divided by the RMSE of the survey, so a number less than one indicates that the attempt to improve the survey was successful, while a number greater than one indicates that the attempt to improve upon the survey failed. The *p*-value is the result of the Harvey, Leybourne and Newbold (1997) variation of the Diebold and Mariano (1995) test for significant forecast differences.

The results of the forecast-improvement exercise are shown in Table 11. We find *RMSE* ratios always greater than one, which means the attempt to improve on the survey makes the forecasts worse. The results are statistically significantly worse at the current-quarter horizon.

VII. Summary and Conclusions

There is some evidence in both the Ball and Croushore, as well as the Rudebusch and Williams, studies that the SPF forecasts of GDP growth are not efficient with respect to the yield spread and to measures of monetary policy. But the evidence suggests that the main inefficiencies arose in the early years of the survey and not since the survey redesign in 1981. In addition, the forecast-improvement exercises suggest that the inefficiencies are not large enough to be exploitable in real time.

TABLE 11—RUDEBUSCH–WILLIAMS FORECAST IMPROVEMENT EXERCISE

Horizon	0Q	1Q	2Q	3Q	4Q
RMSE of Survey	2.023	2.427	2.502	2.670	2.688
RMSE of Improvement to Survey	2.538	2.767	2.795	2.781	2.872
RMSE Ratio	1.254	1.140	1.117	1.041	1.068
<i>p</i> -value	0.02	0.05	0.10	0.06	0.17
Ten-Year Rolling Window					
RMSE of Survey	2.023	2.427	2.502	2.670	2.688
RMSE of Improvement to Survey	2.538	2.767	2.795	2.781	2.872
RMSE Ratio	1.254	1.140	1.117	1.041	1.068
<i>p</i> -value	0.02	0.05	0.10	0.06	0.17
Five-Year Rolling Window					
RMSE of Survey	2.023	2.427	2.502	2.670	2.688
RMSE of Improvement to Survey	2.538	2.767	2.795	2.781	2.872
RMSE Ratio	1.254	1.140	1.117	1.041	1.068
<i>p</i> -value	0.02	0.05	0.10	0.06	0.17

Note: The table shows results of the forecast-improvement exercise. The root-mean-squared error (RMSE) of both the survey and the improvement to the survey are reported in the first two rows of results. The RMSE ratio is the RMSE of the improvement to the survey divided by the RMSE of the survey, so a number less than one indicates that the attempt to improve the survey was successful, while a number greater than one indicates that the attempt to improve upon the survey failed. The *p*-value is the result of the Harvey, Leybourne and Newbold (1997) variation of the Diebold and Mariano (1995) test for significant forecast differences.

More analysis of the change in the survey design is warranted to verify whether it was the change in the survey design, or some other cause, that generated the inefficiencies that we observe.

REFERENCES

- Ball, Laurence, and Dean Croushore.** 2003. "Expectations and the Effects of Monetary Policy." *Journal of Money, Credit, and Banking*, 35: 473–484.
- Croushore, Dean.** 1993. "Introducing: The Survey of Professional Forecasters." *Federal Reserve Bank of Philadelphia Business Review*, 3–13.
- Croushore, Dean.** 2010. "Philadelphia Fed Forecasting Surveys: Their Value for Research." *Federal Reserve Bank of Philadelphia Business Review*, 1–11.
- Croushore, Dean.** 2012. "Forecast Bias in Two Dimensions." *Federal Reserve Bank of Philadelphia Working Paper*, 12-9.
- Croushore, Dean, and Katherine Marsten.** 2014. "The Continuing Power of the Yield Spread in Forecasting Recessions." *Federal Reserve Bank of Philadelphia Working Paper*, 14-5.
- Croushore, Dean, and Tom Stark.** 2001. "A Real-Time Data Set for Macroeconomists." *Journal of Econometrics*, 105: 111–130.
- Diebold, Francis X., and Roberto S. Mariano.** 1995. "Comparing Predictive Accuracy." *Journal of Business and Economic Statistics*, 13: 253–263.
- Faust, Jon, John H. Rogers, and Jonathan H. Wright.** 2005. "News and Noise in G-7 GDP Announcements." *Journal of Money, Credit, and Banking*, 37: 403–419.
- Gneiting, Tilmann.** 2011. "Making and Evaluating Point Forecasts." *Journal of the American Statistical Association*, 106(494): 746–762.
- Harvey, David, Stephen Leybourne, and Paul Newbold.** 1997. "Testing the Equality of Prediction Mean Squared Errors." *International Journal of Forecasting*, 13: 281–291.

Pesaran, M. Hashem, and Martin Weale. 2006. "Survey Expectations." In *Handbook of Economic Forecasting*. Vol. 1, , ed. Graham Elliott, Clive W.J. Granger and Allan Timmermann, 715–776. Elsevier.

Rudebusch, Glenn D., and John C. Williams. 2009. "Forecasting Recessions: The Puzzle of the Enduring Power of the Yield Curve." *Journal of Business and Economic Statistics*, 27(4): 492–503.

Zarnowitz, Victor, and Phillip Braun. 1993. "Twenty-two Years of the NBER-ASA Quarterly Economic Outlook Surveys: Aspects and Comparisons of Forecasting Performance." In *Business Cycles, Indicators and Forecasting*. , ed. James H. Stock and Mark W. Watson, 11–94. University of Chicago Press.