

## Accuracy and Properties of German Business Cycle Forecasts

By Steffen Osterloh\*

### Abstract

In this paper the accuracy of a wide range of German business cycle forecasters is assessed for the period from 1995 to 2005. For this purpose, a data set is used comprising forecasts published on a monthly basis by Consensus Economics. The application of several descriptive as well as statistical measures reveals that the accuracy of the 2-years forecasts is low relative to a simple naïve forecast. This observation can mainly be explained by a systematic overestimation of the growth rates by the forecasters. Moreover, the lack of accuracy can also be explained partly by insufficient information efficiency as well as imitation behaviour. Finally, it is shown that notwithstanding the common errors which affected the accuracy of all forecasters mainly because of their systematic overestimation, they differ significantly in their forecast accuracy.

*Keywords:* business cycle forecasting, forecast evaluation, Consensus forecasts

*JEL Classification:* C52, E32, E37

### 1. Introduction

Forecasting the business cycle is one of the activities of economists which are most critically observed by the public. This forecasting comprises variables like GDP growth and its components, prices, unemployment or interest rates and is conducted by a variety of institutions: by public authorities for budgeting, by central banks for monetary policy, by research institutes for policy consultancy and banks for planning investment strategies. The high number of institutions publishing business cycle forecasts has always aroused public as well as scientific interest regarding the evaluation of their forecast accuracy.<sup>1</sup>

---

\* Centre for European Economic Research (ZEW), P.O. Box 103443, D-68034 Mannheim, Germany; Phone: +49/621/1235-165, Fax: +49/621/1235-223; E-mail: osterloh@zew.de

I would like to thank the participants of the workshop “Makroökonomik und Konjunktur” in Dresden as well as Bernhard Boockmann, Friedrich Heinemann, Marcus Kappler, Waldemar Rotfuß and an anonymous referee for very helpful comments. All remaining errors are mine.

<sup>1</sup> A survey of the history of business cycle forecast evaluation in Germany can be found in Antholz (2006).

In this paper, the emphasis will be on short-term forecasts (with a forecast horizon up to 24 months) of German real GDP growth. This has also been the main object of investigation of most of German research regarding forecast evaluation. In most cases the focus was on forecasts of well-known German institutions like the six institutes and their joint forecast as well as the Council of Economic Advisors (Sachverständigenrat) which were analyzed regarding their accuracy and other desirable properties.<sup>2</sup> A recent and very detailed work in this field is the article of Döpke and Fritsche (2006), who analyse the forecasting performance of 14 German institutions over a long period of more than 30 years, comprising the public research institutes, the Council of Economic Advisors (Sachverständigenrat), international organizations (OECD, IMF, European Commission), some private institutes and the federal government. The authors do not find systematic differences between the forecasters regarding the relative accuracy, besides the date of the publication of the forecast. Thus, in working with intermittently published forecasts an important consideration arises: As they differ in their date of publication, differences in accuracy may not exist because of different economic models or abilities, but forecasters who publish later may be favoured because they possess more information.

This problem is addressed in this paper by using the data published in the Consensus Forecasts survey, which has the special feature of a standardized date of publication and a monthly update of forecasts. Moreover, it enlarges the group of participants by adding a high number of private sector forecasters (mainly from the research departments of banks) which were not involved in former research. Only few notable publications can be found applying comparable data to the evaluation of forecasts. Blix et al. (2001) compare the accuracy of GDP and inflation forecasts of the individual forecasters for 6 different countries. In a more detailed analysis, Harvey et al. (2001) analyze several properties of forecasts made by Consensus participants from the United Kingdom. More research can be found which analyzes pooled Consensus forecasts. Batchelor (2001) compares the accuracy of the Consensus forecasts for the G7 countries with those of OECD and IMF; Öller and Barot (2000) conduct a similar comparison with the forecasts made by OECD and national institutes in Western Europe. Loungani (2001) and Juhn and Loungani (2002) compare the accuracy as well as efficiency of Consensus forecasts with IMF's forecasts for a larger country sample. Most recently, Isiklar and Lahiri (2007) analyze the incorporation of news and Isiklar et al. (2006) investigate the efficiency of the Consensus forecasts for 18 countries.

In this article, first, several standard descriptive measures are applied and the individual forecasters are ranked according to their accuracy. The results for forecasts with a horizon of more than 12 months are surprising as a simple naïve forecast shows the by far highest accuracy when the total German forecasts of the peri-

---

<sup>2</sup> See, e.g., Döpke and Langfeldt (1995), Hagen and Kirchgässner (1997), Grömling (2002), Heilemann and Stekler (2003) and Heilemann (2004).

od from 1995 to 2005 are regarded. The modified Diebold Mariano-test confirms this result of the bad performance of the forecasters compared to the naïve forecast for the longer forecast horizons. However, the Diebold-Mariano-test shows for no forecaster a systematically different forecast error compared to the best participant for all horizons. It is shown that the relatively bad accuracy can mainly be explained by a large overestimation found for all forecasters at the longer forecast horizons which decreases slowly towards the end of the target year. This finding is also confirmed by an extended data set which in addition contains the West German GDP forecasts since 1990.

Moreover, the data allows one to apply a number of tests regarding efficiency and imitation behaviour which can not be conducted with the data commonly used in studies regarding German business cycle forecasts. Another source of inaccuracy may be due to weak information efficiency which is tested via the assumption of unpredictability of forecast revisions. The results suggest that negative news have been incorporated too slowly, which impaired the adjustment of the forecasts from earlier optimistic views to new information. A further approach following Batchelor and Dua (1992) and Gallo et al. (2002) analyzes behavioural biases of individual forecasters. The respective test shows that an imitation of the view of other forecasters can be confirmed empirically for the majority of the forecasters, though large differences in magnitude apply.

Finally, a nonparametric test is used to answer the question whether all forecasters were equal or if some forecasted better than the rest. This rank-sum test looks at the positions of the forecasters in a ranking which is calculated on the basis of the forecast errors for every forecasted year. It can be concluded for the period at hand that not all forecasters are equal. Some of them performed significantly better than a random distribution of the ranks would have suggested and showed a better forecast accuracy than other forecasters.

The paper is organised as follows. First, in Chapter 2 the data is introduced in detail. In Chapter 3, several standard descriptive measures are introduced and applied to the data. In Chapter 4, the differences in forecast accuracy are tested for empirical significance. In the subsequent Chapter 5, two conditions for good forecasts which help to explain the differences in forecast accuracy are tested empirically. In addition, the special features of the data allow us to conduct a specific empirical test which looks at the imitation behaviour of forecasters (Chapter 6). This is followed by a test which looks at systematic differences in the forecast accuracy of the individual forecasters in Chapter 7. The final Chapter 8 concludes.

## 2. The Data

### 2.1 Consensus Forecasts

The data used in this analysis is available from Consensus Forecasts, a monthly survey conducted by the London-based company Consensus Economics. This survey, which is conducted since 1989, publishes forecasts of a variety of forecasters for key macroeconomic variables (in addition to GDP these include its components, prices, industrial production, unemployment and interest rates) for currently 70 countries.

Every month, each forecaster submits two point-forecasts for these variables, one for the current and one for the following year. These questionnaires also contain a precise definition of the predicted variables to ensure comparability. Consensus Forecasts is published in the second week of each month, based on the survey of the panellist forecasts in the two weeks before, with a common deadline (Harvey et al. 2001). This standardization of the date of publication and the definitions for all forecasters guarantees a high degree of comparability of the forecasts.

In addition to the individual participants' forecasts, the arithmetic average of all forecasters is published for every predicted variable, which is widely known as the 'Consensus'. This pooled value is often used by forecasters in combination with their own forecast to communicate to the public their own view relative to the generally expected value.

The relatively high frequency of publications of Consensus Forecasts, which is much higher than the usual publication cycles of the well-known forecasters like IMF, OECD or the German institutes, who often only publish 2–4 forecasts a year, allows the application of a number of specialized empirical procedures. This special feature of the data set is called "fixed-event" forecast in the literature, as one fixed event (the GDP growth rate for the year T) is predicted at a high number of horizons. This contrasts with the usually analyzed "fixed-horizon" forecasts, where point estimates for many years are conducted at only one fixed horizon<sup>3</sup>.

In this case of a fixed-event forecast, the monthly forecasting cycle begins in January of the previous year and ends in December of the forecasted year. Forecasts are therefore provided monthly by each participant, moving from horizons of 24 months up to 1 month ahead, producing altogether 24 forecasts. In the following, forecasts for a given year which were made in the same year (at a horizon of 12 months or less) are considered as "current year" forecasts, forecasts produced in the year before (with a horizon of 13–24 months) are denoted "next year" forecasts.<sup>4</sup>

---

<sup>3</sup> An example for this is the Council of Economic Advisers (Sachverständigenrat), whose forecast for the next year is always published in November of the year before.

<sup>4</sup> For example, next year forecasts for the GDP growth in 2005 were produced in every month between January and December 2004, current year forecasts between January and December 2005.

## 2.2 Participants

The focus of this paper is on a sample which is restricted to the GDP forecasts for the reunified Germany and comprises the time span of 1995 until 2005.<sup>5</sup> This means that the current year forecasts for 1995–2005 and the next year forecasts for 1996–2005 are available. During this period, on average, 30 forecasters participated in the survey each month. These various German participants can be grouped according to their background:

1. Public research institutes<sup>6</sup>: ifo, DIW, RWI, HWWA, IfW
2. Banks:
  - a. ‘Großbanken’ (Deutsche Bank, Commerzbank, Dresdner Bank, HVB)
  - b. ‘Landesbanken’ (e.g., Helaba, Bayerische LB, West LB)
  - c. Co-operative central banks (DZ Bank, WGZ Bank)
  - d. Private banks (e.g., Bank Julius Bär, Sal. Oppenheim)
  - e. Affiliates of foreign banks (e.g., HSBC Trinkaus & Burkardt, SEB Germany)
  - f. Foreign investment banks (e.g., JP Morgan, UBS Warburg)
3. Others:
  - a. Private institutes: FAZ Institut, IW, Economist Intelligence Unit
  - b. Industry: Hoechst AG

For some forecasters, due to missing data, the implementation of several empirical tests becomes impossible. That is why for many of the following empirical tests, different participants had to be dropped, in some cases because the number of forecasted years was too low, in other cases because they did not participate at a sufficient number of horizons<sup>7</sup>.

In addition, some analyses are complemented by the use of an extended data set which comprises in addition to the forecasts for total German GDP, which are only available as from 1995, the forecasts for West German GDP starting with the two-year forecasts made in 1990. However, this data does not allow the coverage of all statistical tests, because it is rather fragmentary due to the sporadic participation of several forecasters in the first years of the Consensus survey. Moreover, it is not perfectly comparable, as the forecasts cover a different geographic unit, which does not rule out structural breaks. However, for some tests it is feasible to apply this extended data set in order to assess the robustness of the results for a longer time

---

<sup>5</sup> Due to the unavailability of three issues of Consensus Forecasts, few values had to be interpolated without influencing the explanatory power of the results.

<sup>6</sup> Although these institutes usually publish their forecasts only every three months, they participate in the Consensus survey at a much higher frequency. Therefore their inclusion in the analyses is justified.

<sup>7</sup> Table 7 in the Appendix shows the problems due to data availability as well as the participation of the forecasters and mergers or acquisitions in detail.

span. In the empirical part it is referred to results from the application of this extended data set.

### 3. Forecast Accuracy

#### 3.1 Descriptive Measures

A first evaluation of the forecast accuracy of the individual forecasters and their pooled forecast, the Consensus, is made with some standard measures, which are briefly introduced (for an overview see Döpke and Fritsche 2006).

A central measure, the forecast error, is defined as  $e_t = F_t - R_t$ , with  $F_t$  representing the predicted value for the year  $t$ , and  $R_t$  the realisation of the value in  $t$ . One aspect discussed controversially in the literature considers the question which value should be used as the realisation. The first preliminary figures for German GDP growth are published in January of the following year, but these figures can also be regarded as estimates as they are usually revised in the course of the following months. It is therefore argued that later publications are more accurate than the first preliminary figures. But in the course of time GDP figures are also often revised due to changes in methodology, which makes them difficult to compare with the initially forecasted values. Therefore Batchelor (2001) proposes to use the actual values published in the middle of the following year as the values used for forecast error computation. In this paper, the values which have been published as realisation in the Consensus Forecasts issue in June of the following year are taken as actual values.<sup>8</sup>

For the descriptive analysis of the accuracy of the forecasters, the following standard measures are considered:

1. Mean absolute error:  $MAE = \frac{1}{T} \sum_{t=1}^T |e_t|$

The MAE averages the absolute errors over all periods, giving positive and negative deviations of the same size the same weight. Using this measure (as well as the following ones), one has to assume a symmetric loss function. This means that negative errors amount to the same loss as positive errors of the same magnitude, which is usually assumed in analysing business cycle forecasts.

2. Mean squared error:  $MSE = \frac{1}{T} \sum_{t=1}^T e_t^2$  and

3. Root mean squared error:  $RMSE = \sqrt{\frac{1}{T} \sum_{t=1}^T e_t^2}$

The idea behind squaring the forecast errors as in these two measures is that large errors should be weighted more than small errors. While an error of 2% is

---

<sup>8</sup> An exception was made in the year 1999 due to the introduction of the ESA95 in April which led to a revision of the 1998 value. There the March value was applied.

twice as severe as an error of 1 % using the MAE, it is four times as severe using the MSE. As it is a widely accepted fact that the main target of GDP forecasting should be to avoid large errors, the RMSE has become the main instrument for measuring forecast accuracy.

$$4. \text{ Theil's } U: U = \frac{RMSE(Model)}{RMSE(Alternative\_Model)}$$

Theil's U (inequality coefficient) has been introduced to provide comparability between forecasts of various variables which have different variances. This is usually the case if one of the variables is more difficult to predict than the other. For this purpose, the RMSE of the forecast of concern is usually divided by the RMSE of a 'naïve' forecast which is used as alternative model, e.g. a random walk model. A value of lower than 1 shows that the model performs better than the alternative model, whereas a value higher than 1 shows that the alternative model is better.

### 3.2 Application

The ranking of the forecasters according to their accuracy can take place by simply calculating the RMSE of every forecaster over all periods and all forecast horizons. This has been done by Blix et al. (2001) for several OECD countries. But this is problematic in the present case, as the data set contains many missing values as discussed in Chapter 2.2. This is critical because the difficulty to forecast differs both between horizons and target years. Figure 1 shows that the average RMSE (across all forecasters and target years) diminishes strongly as we get closer to the end of the predicted year. This is in line with the expected trend: Getting closer to the end of the predicted year, the forecasters have available more information, and the forecasting becomes easier and more accurate.

Figure 2 shows that the average RMSE (over all periods and forecasters) for some target years (1993, 2001–2003) have been more than five times as high as for other years. This confirms the assumption that some years are much more difficult to predict than others.

If some forecasters mainly published in the periods which were easy to predict, this would lead to a low RMSE, but would not say anything about the individual forecasting ability. To allow for this, we use a variant of Theil's U, in which as alternative model not a naïve forecast has been used, but the Consensus. Its RMSE is calculated only for the periods where the individual forecaster also participated. In other words, periods where values of an individual forecaster are missing are dropped out of the RMSE of the Consensus before calculating Theil's U. In this setup, Theil's U value has to be interpreted as the relative accuracy compared to the Consensus for all of the periods where the individual forecaster participated. A value of lower than 1 shows that the individual forecaster was more accurate than the Consensus, a value higher than 1 shows that the Consensus performed better.

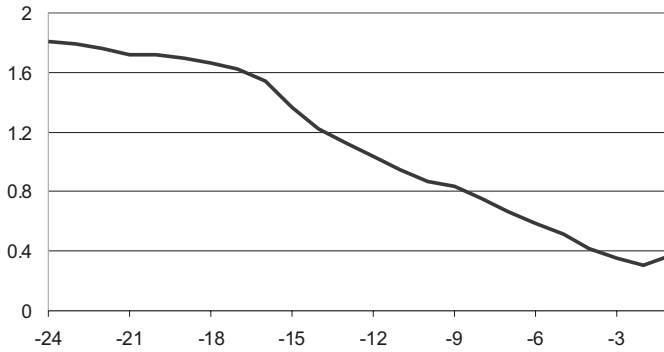


Figure 1: Average RMSE 1991 – 2005

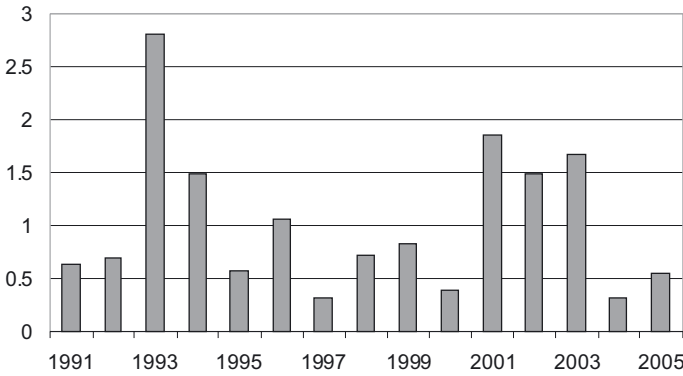


Figure 2: Average RMSE for Target Years

Tables 8 and 9 in the annex show the descriptive statistics for the individual forecasters for the time span 1995–2005, limited to those forecasters who participated in at least 6 years. In addition, the Consensus and a naïve forecast are published as benchmarks.

The choice of a naïve forecast is more or less arbitrary. One option is simply to use the last available actual value to account for short term trends (“no change” forecast). It is also imaginable to use a long-term growth average of sometimes 20 or more years to avoid a domination of short-time cyclical effects. Here, the naïve forecast has been calculated as the average of the published GDP growth rate of the respective past three years (rolling average). This is done to give the most recent growth rates a high weight, but avoid a domination of the cyclical component of the growth rate of a “no change” forecast.

The forecasters are ranked according to their Theil’s U compared to the Consensus as discussed above. In addition, the RMSE and MAE are provided, as well as



the “classical” Theil’s U coefficient which applies a “no change” forecast assuming the same growth rate as in the previous year. Moreover, a further measure shows the percentage of the periods where the forecaster’s value was closer to the realisation than the Consensus. It can easily be seen that in both rankings the different measures do not result in the exact same order, but still a ranking based on any other measure would be similar to the Theil’s U criteria.

The position of the Consensus is in both rankings above the average, which is consistent with other authors’ results. These results are also confirmed by theory, as McNees (1992) discusses. Extreme forecasts (which are very often wrong) cancel out, which lets McNees say that “many heads are better than one”.

Table 8 shows the results of the current year forecasts only (comprising the horizons of 12 until 1 months). Here, unsurprisingly, the naïve forecast reaches the last position. This is what one would expect, as this naïve forecast only contains the information of the past three years, but does not account at all for the development in the current year.

A very different picture emerges from the values of the next year forecasts (Table 9), which comprise the horizons of 24 until 13 months. Very surprisingly, the naïve forecast performs by far the best.<sup>9</sup> The weak performance of the forecasters is also reflected by the Theil’s U coefficients applying the “no change” forecast, which are in most cases higher than one. This striking result, which can not often be found in the forecast evaluation literature, has to be explained in the following chapters.<sup>10</sup>

Regarding the relative accuracy of the individual forecasters, it can be observed that in both rankings mainly less renowned forecasters rank at the first positions, confirming a similar result by Blix et al. (2001). The three forecasters who constantly show the highest accuracy are HSBC Trinkaus & Burkhardt, BfG Bank (later SEB Germany) and MM Warburg, while more prestigious forecasters like one research institute and one ‘Großbank’ can be found at the last ranks. But this result has to be regarded cautiously. As will be shown in the next chapters, this period seems to be quite special in its predictability. In addition, the results of the study conducted by Blix et al. for an earlier period show very different results. There, Trinkaus & Burkhardt was among the worst, while some underperformers of this study are among the best. Therefore, a generalization of the results should not be made.

---

<sup>9</sup> However, it has to be emphasized that the naïve forecast does not perform among the best when the sample extended by the West German forecasts since 1990 is considered. This is due a high over-estimation of the naïve forecast in the first years when the extraordinary growth rates because of the reunification boom are incorporated.

<sup>10</sup> Isiklar and Lahiri (2006) show similarly that for 11 out of 18 examined OECD countries, the 24-months-ahead forecast of the Consensus can not beat a naïve forecast. But the differences in accuracy in terms of Theil’s U are in most cases rather small.

#### 4. Test for Difference of Forecast Errors

In addition to the descriptive measures presented in the previous chapter, a statistical test introduced by Diebold and Mariano (1995) is used to test the differences of forecast errors for statistical significance. This test is an asymptotic  $t$ -test, generally testing for the null hypothesis that the difference of the mean squared errors of two forecast models A and B is zero for a given forecast horizon, i.e.  $\bar{d} = MSE_A - MSE_B = 0$ .<sup>11</sup> The test statistic of the Diebold-Mariano test is

$$DM = \frac{\bar{d}}{\sqrt{\hat{V}(\bar{d})}},$$

which follows a  $t$ -distribution with  $(N - 1)$  degrees of freedom,  $N =$  number of forecasts.

In this test, the variance  $\hat{V}$  is estimated robustly as in the work of Newey and West (1987). This allows to account for the autocorrelation problem due to overlapping periods. As for a two-step forecast an unpredicted event does not only affect the forecast error for the current year (which was predicted in the year before), but also for the next year (which has been predicted earlier in the year of the event), a positive autocorrelation between the two forecast errors can be expected.

In this paper, a modified version of the Diebold-Mariano test is used. This has been proposed for small samples by Harvey, Leybourne and Newbold (1997). The test statistic of this modified version is  $mDM = C \cdot DM$ , with the correction factor

$$C = \left( \frac{N + 1 - 2h + N^{-1}h(h - 1)}{N} \right)^{\frac{1}{2}},$$

with  $N =$  number of forecasts and  $h =$  steps of forecast. The calculated correction factors are for the 1995–2005 sample 0.849 for the next year forecasts and 0.953 for the current year forecasts (1990–2005 sample: 0.90 and 0.966, respectively).

Table 1 shows the results of the bilateral comparisons of the individual forecasters with the overall best forecaster who participated consistently (Affiliate 1) as the benchmark. The forecast horizons of 23, 18, 11 and 6 months have been used. Only forecasters have been considered who published forecasts in every year for the selected horizons<sup>12</sup>.

<sup>11</sup> Diebold and Mariano (1995) consider a general loss function. For the evaluation of business cycle forecasts, it is common to use a MSE loss differential.

<sup>12</sup> Some fields in Table 1 are empty due to missing data for the respective forecaster and horizon.

Table 1

**Diebold Mariano Test, Accuracy Compared to Best Forecaster**

	Total Germany 1995–2005				West Germany 1990–1994 + Total Germany 1995–2005			
	Horizon				Horizon			
	–23	–18	–11	–6	–23	–18	–11	–6
Institute 2		–0.03 (0.159)		–0.01 (0.022)				–0.12 (0.074)
Institute 4		0.06 (0.157)	0.39** (0.142)	0.33* (0.169)		–0.26 (0.335)	0.05 (0.174)	0.20 (0.149)
Institute 5			0.28* (0.124)					
Others 1	–0.16 (0.246)	0.30 (0.235)	0.16 (0.109)	0.14 (0.079)				
Großbank 1	0.29 (0.610)	0.62 (0.483)	–0.14 (0.152)	0.07 (0.087)	0.91* (0.445)	0.58 (0.371)	–0.16 (0.118)	–0.03 (0.111)
Großbank 2	0.32 (0.373)	0.40* (0.165)	0.04 (0.078)	0.07 (0.046)		0.45* (0.208)	0.19 (0.132)	0.05 (0.046)
Großbank 4	0.24 (0.213)	0.35 (0.240)	0.34 (0.189)	0.00 (0.091)		0.15 (0.285)	0.34 (0.191)	0.08 (0.125)
Co-operative 1	0.16 (0.193)	0.20 (0.225)	0.11* (0.058)	0.04 (0.040)	0.31 (0.276)	0.39 (0.268)	0.00 (0.104)	0.05 (0.044)
Co-operative 2	0.57 (0.417)	0.44 (0.294)	0.13 (0.170)	0.00 (0.081)	0.54 (0.320)	–0.02 (0.420)	0.17 (0.184)	–0.07 (0.105)
Landesbank 1	–0.35 (0.211)	–0.03 (0.187)	0.02 (0.057)	0.10 (0.062)		0.25 (0.280)	0.13 (0.085)	0.06 (0.054)
Landesbank 2	0.44 (0.514)	0.30 (0.411)	0.18 (0.096)	0.09** (0.033)	0.87 (0.605)	0.50 (0.343)	0.58 (0.385)	0.00 (0.082)
Landesbank 3	0.70 (0.669)	0.43 (0.357)	0.11 (0.087)	0.06 (0.056)	0.81 (0.437)	0.28 (0.291)	0.27* (0.147)	0.02 (0.070)
Landesbank 4	0.67 (0.402)	0.55 (0.318)	0.15 (0.194)	0.07 (0.059)	0.74** (0.274)	0.58* (0.266)	0.14 (0.143)	0.02 (0.066)
Landesbank 5	–0.12 (0.332)	–0.05 (0.225)	0.24 (0.172)	0.03 (0.060)	0.34 (0.549)	0.54 (0.617)	0.23 (0.152)	–0.07 (0.109)
Affiliate 1					benchmark			
Affiliate 2	–0.48 (0.383)	–0.28 (0.328)	0.03 (0.095)	0.01 (0.066)		0.03 (0.357)	0.12 (0.091)	–0.09 (0.105)
Affiliate 3	0.20 (0.286)	0.76* (0.350)	0.01 (0.167)	0.00 (0.030)		0.33 (0.397)	0.02 (0.137)	–0.03 (0.065)
Private 1	0.05 (0.266)	0.26 (0.299)	–0.21 (0.132)	–0.02 (0.056)				
Private 2	–0.13 (0.358)	0.08 (0.260)	–0.07 (0.047)	0.05 (0.046)				
Consensus	0.15 (0.205)	0.10 (0.176)	0.05 (0.060)	0.03 (0.046)	0.32 (0.262)	0.12 (0.154)	0.04 (0.046)	–0.06 (0.075)
naïve	–1.07** (0.439)	–0.65* (0.346)	0.58* (0.319)	0.92** (0.445)	0.66 (1.417)	0.89 (1.340)	2.37 (1.656)	2.80 (1.753)

The symbol \*\*\*, \*\*, and \* denotes rejection of the null hypothesis at the 1%, 5%, and 10% level respectively. Standard errors in parentheses.

It can be seen that mostly positive values (forecaster has higher MSE than the best forecaster) can be observed, but also negative values (forecaster has lower MSE) can be seen, and in most cases the null hypothesis that the mean squared forecast errors are equal can not be rejected. No forecaster has positive or negative results significantly different from zero for more than two horizons. The addition of the West German forecasts since 1990 yields overall similar results, but in some cases these differ highly from the results for the shorter sample. This reflects that due to its construction the DM-test reacts very sensitive to extreme forecast errors; in this case, the extraordinary high forecast errors for the year 1993 dominate the results. In that year the forecasters overestimated the growth rate 23 months ahead by more than 4%. All in all, it can be said that this test does not allow us to conclude that the overall best forecaster performed significantly better or worse than the others for all horizons.

Moreover, the modified DM-test was applied to a bilateral comparison with the naïve forecast introduced in Chapter 3.2. This test confirms the results of the descriptive statistics. For the two longer horizons (next year forecasts) of the 1995–2005 sample, the null hypothesis of smaller mean squared errors of the best forecaster compared to the naïve forecast can be rejected. The naïve forecast even performs significantly better than the best individual forecaster (negative values).

Analogous to the results of the descriptive statistics, this changes with the forecast horizons getting shorter as well as for the extended sample. The sign becomes positive for the two horizons of the current year forecasts (significantly smaller than zero at the 5% and 10% level, respectively), indicating that their errors are higher than the errors of the best forecaster. Considering the extended sample, the naïve forecast always performs worse than the best individual forecaster. This is due to highly overoptimistic naïve forecasts after the German reunification, which translated into large overestimations in the recession thereafter. This effect is also reflected in the large standard errors.

Summing up, the results of the modified DM-test confirm the finding of the descriptive statistics of a surprisingly bad performance of the individual forecasters as well as the Consensus compared to a simple naïve forecast for the next year forecasts of the total Germany forecasts.

## 5. Conditions for Good Forecasts

The following sections are intended to explain the findings of the surprisingly bad accuracy of the next year forecasts found in the previous chapters. For this purpose two main conditions for good forecasts are introduced and tested empirically, unbiasedness and information efficiency.

### 5.1 Unbiasedness

A first condition for good forecasts which is analyzed is unbiasedness. A forecast is considered to be biased if it is systematically too high or too low. If this is the case, a forecast is suboptimal because it can easily be improved on the basis of the bias known from past forecasts. In the case of an upward bias, this improvement could easily be made by subtracting the average overestimation from the forecast.

To verify this empirically, a simple  $t$ -test is used, regressing the forecast errors for a given forecast horizon on a constant:  $e_t = \alpha + u_t$  (Fildes and Stekler 2002). The null hypothesis, that the forecasts are unbiased, would hold if  $\alpha = 0$ . A negative value would show an underestimation, a positive value an overestimation. For this test a normal distribution of the forecast errors has to be assumed, which can not be rejected on the basis of the Jarque-Bera statistics (Bera and Jarque 1980).

In a first step, the participants were pooled according to the groups presented in section 2.2. Similarly to the Diebold-Mariano-test, robust standard errors (Period SUR) have to be used for the two-year forecasts because of possible autocorrelation (for details regarding the panel-corrected standard errors methodology, see Beck and Katz 1995). Table 2 shows the results for the four horizons (23, 18, 11 and 6 months) for the groups with several participants for each horizon. It can be seen that all groups of forecasters showed for all horizons positive biases which turn out to be significantly different from zero in all cases. The null hypothesis of unbiasedness can be rejected there for all forecast horizons. However, these biases do not differ considerably in their size between the groups.

Table 2

#### Pooled Test for Biasedness

	Großbank	Co-operative	Landesbank	Affiliate	Private
23 months	1.20*** (0.029)	1.19*** (0.029)	1.05*** (0.072)	0.98*** (0.089)	1.07*** (0.033)
18 months	1.11*** (0.025)	1.03*** (0.040)	0.92*** (0.049)	0.88*** (0.080)	0.95*** (0.036)
11 months	0.47*** (0.109)	0.52*** (0.133)	0.43*** (0.094)	0.38*** (0.110)	0.25* (0.123)
6 months	0.29*** (0.067)	0.29*** (0.076)	0.14** (0.066)	0.17** (0.069)	0.19** (0.087)

The symbol \*\*\*, \*\*, and \* denotes rejection of the null hypothesis at the 1%, 5%, and 10% level respectively. Standard errors in parentheses.

In Table 3 the results for all individual forecasters who participated in all of the years at the given horizons are depicted. In addition, the results of the data set which was expanded with the West German forecasts since 1990 are added, were

Table 3  
Test for Biasedness

	Total Germany 1995–2005				West Germany 1990–1994 + Total Germany 1995–2005			
	Horizon				Horizon			
	–23	–18	–11	–6	–23	–18	–11	–6
Institute 2		0.90** (0.328)		0.12 (0.100)				–0.02 (0.116)
Institute 4		0.97** (0.319)	0.49** (0.207)	0.27 (0.180)		1.07*** (0.298)	0.27 (0.211)	0.02 (0.214)
Institute 5			0.69*** (0.148)					
Others 1	1.08*** (0.296)	1.09*** (0.323)	0.51** (0.176)	0.21 (0.152)				
Großbank 1	1.15** (0.410)	1.08** (0.403)	0.32 (0.195)	0.26* (0.143)	1.27*** (0.310)	1.00*** (0.308)	0.27* (0.155)	0.06 (0.178)
Großbank 2	1.27*** (0.320)	1.17*** (0.316)	0.45** (0.174)	0.30** (0.100)		1.07*** (0.252)	0.26 (0.205)	0.08 (0.166)
Großbank 4	1.19*** (0.349)	1.08*** (0.295)	0.65*** (0.182)	0.30** (0.096)		0.89*** (0.205)	0.54*** (0.170)	0.22 (0.129)
Co-operative 1	1.15*** (0.337)	0.97** (0.339)	0.51** (0.171)	0.27** (0.109)	1.11*** (0.265)	0.91*** (0.290)	0.32* (0.179)	0.03 (0.190)
Co-operative 2	1.23** (0.399)	1.08** (0.359)	0.54** (0.173)	0.30** (0.115)	1.09*** (0.305)	0.97*** (0.255)	0.26 (0.192)	0.05 (0.192)
Landesbank 1	0.90** (0.352)	0.80* (0.373)	0.33 (0.190)	0.07 (0.190)		0.69** (0.302)	0.11 (0.203)	–0.11 (0.193)
Landesbank 2	1.08** (0.468)	1.00** (0.389)	0.47** (0.199)	0.03 (0.160)	1.14*** (0.356)	0.95*** (0.289)	0.37* (0.186)	–0.13 (0.158)
Landesbank 3	1.07** (0.510)	0.86* (0.424)	0.37 (0.218)	0.15 (0.160)	0.87* (0.405)	0.70** (0.322)	0.09 (0.248)	–0.03 (0.187)
Landesbank 4	1.32*** (0.386)	1.09** (0.359)	0.44** (0.190)	0.26** (0.104)	1.07*** (0.320)	0.95*** (0.288)	0.28 (0.181)	0.09 (0.138)
Landesbank 5	0.87* (0.429)	0.84** (0.350)	0.52** (0.217)	0.18 (0.123)	0.88** (0.329)	0.99*** (0.315)	0.45** (0.183)	0.09 (0.120)
Affiliate 1	1.00** (0.375)	0.91** (0.320)	0.29 (0.182)	0.15 (0.126)	0.85*** (0.280)	0.81*** (0.244)	0.17 (0.163)	–0.07 (0.160)
Affiliate 2	0.79** (0.308)	0.70** (0.295)	0.41* (0.191)	0.14 (0.115)		0.73** (0.257)	0.26 (0.235)	0.00 (0.128)
Affiliate 3	1.16*** (0.295)	1.03** (0.359)	0.45*** (0.136)	0.19*** (0.070)		0.92*** (0.264)	0.34*** (0.111)	0.10 (0.107)
Private 1	1.11*** (0.337)	1.00** (0.355)	0.27 (0.166)	0.18 (0.110)				
Private 2	1.02** (0.351)	0.90** (0.362)	0.24 (0.193)	0.19 (0.141)				
Consensus	1.11** (0.370)	0.95** (0.330)	0.45** (0.173)	0.20 (0.120)	1.01*** (0.286)	0.87*** (0.260)	0.28 (0.172)	0.01 (0.157)
naïve	0.21 (0.416)	0.21 (0.416)	0.05 (0.335)	0.05 (0.335)	0.72 (0.502)	0.72 (0.502)	0.52 (0.454)	0.52 (0.454)

The symbol \*\*\*, \*\*, and \* denotes rejection of the null hypothesis at the 1%, 5%, and 10% level respectively. Standard errors in parentheses.

available. Again, robust standard errors (according to Newey and West 1987) were applied. It can be observed that for both samples all forecasters showed for the two longest horizons highly positive biases which turn out to be highly significantly different from zero in all cases, so that the null hypothesis of unbiasedness can be rejected there. For the 23-months horizon, they show values of 0.79–1.32, which suggests that the GDP growth rate was overestimated on average by more than 1 percentage point by most forecasters, and by the Consensus forecast by 1.11 and 1.01 points, respectively. Getting closer to the end of the predicted year, this bias gets smaller in magnitude and loses significance in some cases. But even six months before the end of the target year, in many cases a significant positive bias can be observed.

This test has also been conducted for the naïve forecast introduced in Chapter 3.2. This naïve forecast shows for the period 1995 to 2005 only a very small positive bias, which is not significantly different from zero. The null hypothesis of unbiasedness can not be rejected here. Comparing this result with the results of the individual forecasters, a first explanation can be drawn for the bad accuracy of the forecasters for the next year forecasts in that period. It seems that the high positive bias in the early forecasts accounts for the much better accuracy of the naïve forecast for the two longer horizons as shown by the Diebold-Mariano test in Chapter 4.

## 5.2 Information Efficiency

Another condition of a good forecast is information efficiency. This means that an efficient forecast should efficiently incorporate a certain set of information which is known at the moment of the production of the forecast. If this is not the case, the forecast can easily be improved by incorporating the missing information.

### 5.2.1 Theory

Regarding the analysis of forecast efficiency, it has to be distinguished between its strong and its weak definition (see Kirchgässner and Müller 2006). The strong definition of information efficiency comprises all information which is available at the date of the production of a forecast. This means that a forecast is said to be strongly information efficient if it efficiently incorporates any available information. It is obviously very difficult to test for strong efficiency empirically as it is practically impossible to find all data available. For practical use, tests are usually restricted to some key variables like interest rates, oil prices or business surveys.

A more feasible definition is weak information efficiency. In this case the set of information is restricted to the past forecast errors. Weak information efficiency holds if a forecast efficiently incorporates all information about its past forecast errors. These errors therefore have to be unpredictable. For this purpose, the fol-

lowing regression is commonly used:  $e_t = \alpha + \beta e_{t-1} + u_t$ . The null hypothesis is that  $\beta = 0$ , which means that the current forecast error can not be explained by its past errors. Otherwise it would be possible to improve the forecast on the basis of the knowledge from past errors.

For the special case of a fixed-event forecast, Nordhaus (1987) proposes a specific test which does not aim at the unpredictability of the forecast errors, but at the unpredictability of forecast revisions. This test seems to be more appropriate for the data set used in this paper because of the low number of predicted years. In his approach, if weak efficiency holds, the forecast revision process should look like a random walk  ${}_jF_T = {}_kF_T + \sum_{s=k+1}^j \varepsilon_s, j > k, s = (k, \dots, j)$ . Each value predicted in later periods  ${}_jF_T$  should be a combination of the initial published value  ${}_kF_T$  and the sum of the revisions  $\varepsilon_s$  in the periods before. These revisions have an expected value of zero and should be identically independent distributed. Nordhaus (1987) gives the following intuition for this idea: “If I could look at your most recent forecasts and accurately say, ‘Your next forecast will be 2 percent lower than today’s,’ then you can surely improve your forecasts.”

The test for this random walk behaviour therefore looks at the correlation of forecast revisions, which are defined as  ${}_t v_T = {}_t F_T - {}_{t-1} F_T$ , i.e. the difference between the forecast in the current and the previous period for the same target year. Weak efficiency holds if these are unpredictable (white noise,  $\text{cov}({}_{t-1} v_T, {}_t v_T) = 0$ ). To test for this, the following model is estimated:  ${}_t v_T = \alpha({}_{t-1} v_T) + u_t$ , under the null hypothesis that the revisions are unpredictable:  $H_0 : \alpha = 0$ .

### 5.2.2 Application

Applying the original version of the Nordhaus test on the data at hand is problematic for this data set. The general assumption of an expected value of zero for the forecast revisions seems to be violated. For all the participants the number of negative revisions is much higher than the number of positive revisions. The Consensus shows since 1995 30 upward and 92 downward revisions. Its mean revision has a value of  $-0.06$  percentage points. This also holds for all of the individual forecasters, e.g. the forecasters with the most revisions, Landesbank 3 (34 upward, 75 downward) and Co-operative 2 (23 upward and 67 downward). The reason for this can be seen from Figure 3. This shows the development of the Consensus and the maximum and the minimum value of the individual forecasters in a year with a very high number of negative revisions. As it has been shown in Chapter 4.2, in an average year all forecasters started out with a large overestimation of the growth rate for the 24-months forecast. With decreasing horizons towards the end of the target year, more news became available contradicting their optimistic views. This caused a high necessity to revise the forecasts downwards, which can be seen in the downward trend of the Consensus forecast in the chart. This pattern can be observed for a high number of the years between 1995 and 2005.



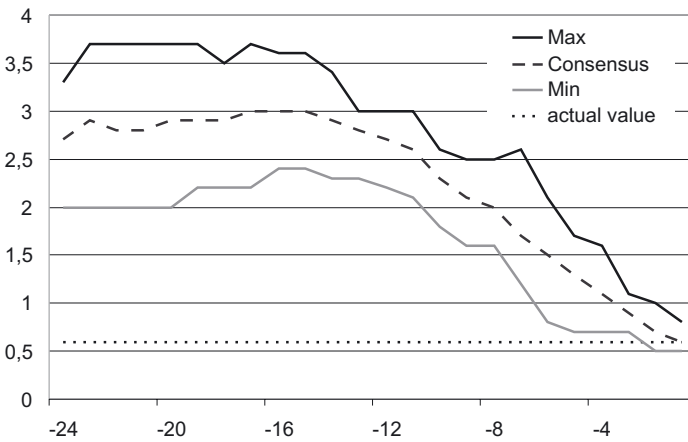


Figure 3: Forecast Revision Process (Target Year: 2001)

A measure introduced by Isiklar et al. (2006) based on the general principle of the Nordhaus test is used to account for the violation of the assumption of an expected value of zero for the forecast revisions. They compare the frequency distribution of the direction of the forecast revisions (upward or downward) with the expected distribution under the hypothesis of independence. Under the null hypothesis, the probability of the direction of a revision should be independent from the direction of the revision one month earlier and therefore be equal to its expected value.

This is applied to the monthly revisions of the Consensus and some selected individual forecasters (those who show the highest number of monthly revisions) in Table 4. Regarding positive revisions, for all forecasters a very low number of revisions in two subsequent months can be observed. This value is mostly similar to the expected value which is also very low because of the low number of positive revisions in general. But for the negative revisions, a much higher frequency of subsequent revisions can be seen for the Consensus than to be expected under the null hypothesis of independence. That is, given the hypothetical independent distribution, approximately 38 cases of two negative revisions in series would be expected, but actually 52 can be observed. This is also the case for the majority of the individual forecasters.

Without going further into an empirical analysis, these results point to a possible lack of information efficiency regarding the negative revisions. In most years, the forecasters started the forecasting cycle with an overestimation at the 24-months-ahead forecast. With the appearance of news, the forecasters revised their forecasts downwards. But the positive correlation and the higher than expected occurrence of subsequent periods with downward revisions show that these revisions were predictable. It seems that the news were not incorporated by conducting large negative revisions, but many small steps of subsequent negative revisions.

*Table 4*  
**Subsequent Forecast Revisions**

	Frequency/Expected frequency	
	${}_t v_T > 0$ and ${}_{t-1} v_T > 0$	${}_t v_T < 0$ and ${}_{t-1} v_T < 0$
Consensus	4 / 4.1	52 / 38.1
Others 1	1 / 1.2	28 / 18.9
Großbank 1	3 / 2.9	20 / 13.8
Landesbank 1	1 / 2.5	15 / 14.5
Landesbank 3	5 / 5.4	28 / 26.3
Landesbank 5	1 / 1.9	17 / 15.9
Co-operative 1	2 / 1.6	23 / 14.2
Co-operative 2	1 / 2.4	26 / 20.6
Affiliate 1	4 / 1.9	20 / 14.7
Affiliate 3	5 / 2.1	17 / 15.3
Foreign 5	2 / 2.0	19 / 20.9

### 5.2.3 Explanations

Nordhaus (1987) gives two possible explanations for this behaviour of the individual forecasters which he calls “forecast smoothing”. First, he argues that professional forecasters are fearful that “jumpy” forecasts will be treated as inconsistency by their bosses or customers. Therefore a high jump which is indicated by new data or events is distributed over several small jumps in subsequent months. They are especially reluctant to make a positive revision following a negative one or vice versa, even if the data points to this direction.

The other explanation is a psychological one, saying that forecasters tend to hold to their prior views too long, and therefore incorporate data opposing their views too slowly. Kirchgässner and Müller (2006) develop a specific loss function for forecasters based on these ideas of costly forecast revisions and find empirical support of unwillingness to revise for the forecasts of German institutes.

A possible explanation for low information efficiency of the Consensus can be found in a herding behaviour. If news is not incorporated by all forecasters at the same time, it is not reflected in a large revision of the Consensus, but spread via smaller revisions over several months. This will be discussed further in the next chapter.

## 6. Imitation Behaviour

One further possible source of the inaccuracy of forecasts is introduced by Gallo, Granger and Jeon (2002), who suspect an imitation behaviour of forecasters. This is the case when views expressed by other forecasters in the previous periods have an influence on an individual’s current forecast. They confirm this empiri-

cally by using the lagged mean value (the Consensus) as a proxy for the view of the other forecasters.

This imitation behaviour might reduce the forecast accuracy as it leads to a convergence to “a value which is not the “right” target”. Gallo et al. (2002) explain this behaviour with a possible aversion of the forecasters to produce extreme forecasts. If they see that their own forecast is too far away from the other forecasts (or the Consensus), they start wondering that they are possibly wrong in their view and revise towards the Consensus.

### 6.1 Methodology

Gallo et al. (2002) test the assumption of imitation behaviour empirically by running the following regression using OLS:  $F_{T,t}^i = \alpha + \beta F_{T,t-1}^i + \gamma Cons_{T,t-1} + \delta \sigma_{T,t-1} + u_{T,t}$ .

The current forecast of an individual forecaster is regressed on a constant, his own forecast published in the previous month ( $F_{T,t-1}^i$ ), the value of the Consensus which was published in the previous month, but which is only known in the current period ( $Cons_{T,t-1}$ ), and the standard deviation of all forecasts in the previous month. This is meant to capture the effect related to the forecasts moving closer together as the time-horizon decreases ( $\sigma_{T,t-1}$ ). A high  $\beta$  indicates a low likelihood that a forecaster changes his mind in subsequent periods. The sign of  $\gamma$  shows whether the movement of the individual is in agreement with the movement of the Consensus. A positive value shows that the forecaster tends to revise his forecasts in the same direction as the rest.

### 6.2 Results

In order to test for different behaviour of the forecaster groups as defined in Chapter 2.2, a panel estimation using pooled OLS was applied, in which a common gamma coefficient among the participants of each group was assumed. Similarly to the results of Gallo et al. (2002), a high explanatory capability can be observed ( $R^2$  ranging from 0.94 to 0.97). Also confirmed can be a high persistence effect, all coefficients are positive and highly significantly different from zero. The interesting factor regarding the hypothesis of imitation behaviour are the  $\gamma$  coefficients for which results are depicted in Table 5. It is apparent that for all groups of forecasters this factor is positive, but the coefficients differ a lot in their magnitude. The groups of the foreign banks and the largest German banks show the by far smallest coefficients, the groups of the “other” participants and the co-operative banks show the highest coefficients.

This result does not seem to be a coincidence, as it may be expected that the foreign investment banks in their forecasts put more emphasis on global factors and less on the views of other German forecasters. Moreover, the researchers of the four largest German banks (Deutsche Bank, Dresdner Bank, Commerzbank, HVB) probably have the highest reputation among the private sector forecasters.

Table 5

**Pooled Test for Imitation Behaviour**

	Institute	Groß- bank	Co- operative	Landes- bank	Affiliate	Private	Foreign	Other
Gamma	0.26*** (0.028)	0.16*** (0.026)	0.39*** (0.041)	0.31*** (0.022)	0.22*** (0.028)	0.25*** (0.023)	0.13*** (0.025)	0.45*** (0.034)

The symbol \*\*\*, \*\*, and \* denotes rejection of the null hypothesis at the 1%, 5%, and 10% level respectively. Standard errors in parentheses.

In addition in Table 10 in the Appendix the results for the individual forecasters are depicted. Again, the three lowest values come from three U.S. investment banks, and two of the values are even not significantly different from zero. Meanwhile some other forecasters, often from smaller research departments, reach very high values. Some of them are even of the same magnitude as the values of their own lagged forecasts.

Summing up, the results indicate that for almost all of the forecasters the hypothesis of imitation behaviour can not be rejected empirically, but this effect differs markedly in its magnitude between the groups of forecasters.

## 7. Test for Differences in Forecast Accuracy

Chapter 4 indicated that the Diebold-Mariano test was not able to show that the best forecaster was significantly better than the other participants for all horizons in terms of mean squared errors. However, this test does not answer the question whether all forecasters were jointly equal their in accuracy because of the following two reasons: First, the DM test only allows for bilateral comparisons for a fixed horizon, and does not consider all participants of the survey and all horizons. Second, the literature as discussed in Stekler (1987) shows that tests based on root mean squared errors (such as the DM-test) are not an appropriate tool to answer the question at hand, because such tests overweight target years with higher forecasting uncertainty.

Therefore, in this chapter, an alternative nonparametric test proposed by Stekler (1987) is presented to answer empirically the question whether all forecasters were equal or if some forecasted better than the rest.

### 7.1 Methodology

A first idea to test for individual differences in accuracy would be to conduct an *F*-test of the forecast errors across the forecasters, but Batchelor (1990) argues that this would not be legitimate. As was shown in Figure 2, some years were more

difficult to predict than others. If the average forecast errors were used for the analysis, years which were hard to forecast (with a higher variance of the errors) would dominate the results. Therefore, Stekler proposes a rank-sum-test, which constitutes a nonparametric test. This test only looks at the positions of the forecasters in a ranking which is calculated on the basis of the forecast errors for every forecasted year.

In a first step, all forecasters  $(1, \dots, n)$  are ranked for every target year according to their forecast error, with the value of 1 given to the best and the last rank  $(n)$  to the worst. Stekler (1987) uses the RMSE over all forecast horizons to rank the forecasters for every target year. But this approach does not seem to be appropriate for this data set given its many missing values as discussed in Chapter 2. The usage of the RMSE would favour the forecasters who mainly participated at the short horizons. Therefore, in this paper, the forecasters have been ranked according to their Theil's U compared to the Consensus as introduced in Chapter 3.2. Then, for every forecaster  $i$ , the rank sum  $r_i = \sum_{t=1}^T r_{it}$  over all predicted years  $(1, \dots, T)$  was calculated.

For this procedure the sample has been restricted to the 22 forecasters who took part in the survey in all years between 1995 and 2005. Table 11 in the annex shows their ranks for Theil's U calculated for both the next and the current year forecasts (comprising all 24 horizons). It appears that there are big differences in forecast accuracy, with some forecasters consistently performing better than the average and others consistently ranking lowest.

The assumption that forecasters are not equal in their forecast accuracy – i.e. that their positions are not random – was tested empirically. The null hypothesis claims that all forecasters are equal. This means that looking at different periods, there should be no systematic differences in the rank distribution. Thus, for every year the rank of an individual forecaster has to be identical to the average rank, i.e.  $(n + 1)/2$ . For 22 forecasters this average rank is 11.5. Therefore it is claimed that the null hypothesis is  $H_0 : r_i = T(n + 1)/2$ , such that for every forecaster the rank sum over all target years  $T$  should be identical to its expected value.

To test this empirically, Batchelor (1990) developed the following test statistic, which follows a  $\chi_{n-1}^2$  distribution under the null hypothesis:

$$f = \sum_{i=1}^n \frac{\{r_i - T(n + 1)/2\}^2}{Tn(n + 1)/12}.$$

The nominator shows the squared deviation of each forecaster's rank sum from its expected value. In the denominator, the variance of an individual rank statistic  $[n(n + 1)/12]$  for the sum of  $T$  individual ranks is used.

## 7.2 Results

Three test statistics have been calculated: one comprising all 24 forecast horizons, one only for the next year forecasts and one for the current year forecasts. Table 6 presents the results of the test statistics for the period 1995–2005. The result of the test statistic for all 24 forecast horizons shows the highest significance: The null hypothesis that the forecast accuracy of all participants is equal is rejected at the 1% level. The values of the current year forecasts are significant at the 5% level, the significance of the next year forecasts misses slightly the 10% level. This is contrary to the results of Batchelor (1990), who did not find any significant differences between the accuracy of the forecasters in the United States<sup>13</sup>.

Table 6

### Results of the Rank Sum Test

Whole period	Current year	Next year
40.14***	37.12**	29.17

Critical values for 21 degrees of freedom: 1%: 38.98; 5%: 32.67; 10%: 29.62.

The symbol \*\*\*, \*\*, and \* denotes rejection of the null hypothesis at the 1%, 5%, and 10% level respectively.

Therefore, it can be concluded that the different positions found in the rankings are not random. Not all forecasters are equal. Some of them performed significantly better than a random distribution of the ranks would have suggested and showed a better forecast accuracy than other forecasters.

## 8. Conclusions

The empirical analysis of the data from the Consensus Forecasts survey allows us to draw a number of interesting conclusions regarding the performance of German business cycle forecasters in the period from 1995 to 2005.

The most striking result is the weak forecast accuracy of all forecasters for the next year forecasts, especially relative to a simple naïve forecast. As it has been shown, this can mainly be explained by the large positive bias which all forecasters show for the longer forecast horizons in the analyzed time period. This bias could be confirmed empirically for all forecasters, showing that they started off with a systematic overestimation at the longest forecast horizon which decreased slowly

<sup>13</sup> It has to be noted that the application of this test to the extended panel does not lead to a rejection of the null hypothesis. However, this is due to the unavoidable omission of 5 of the 22 forecasters who did not participate since the beginning of the survey. A closer inspection showed that these 5 forecasters had a major impact on the results, as they constituted two of the worst and one of the best participants.

while getting closer to the end of the target year. This finding is also confirmed by the longer panel which additionally covered the West German forecasts from 1990 to 1994

Although the comprehensive discussion of the reasons for this bias is beyond the scope of this paper, two aspects seem to be important. The most obvious explanation for the overestimation appears to be the assumption that the decline of the German trend growth rate since the mid-90s was not expected and therefore not incorporated into the forecasts. But this explanation may be dissatisfying because the reasons underlying the decline of potential growth were known by the forecasters for some time before the decline happened. Moreover it should not be neglected that these 10 years were dominated by a number of unpredictable negative macroeconomic shocks, like the bursting of the new economy bubble or the terrorist attacks of 9/11. A more detailed analysis of these explanatory factors should be undertaken in future research.

Another source of inaccuracy may be associated with weak information efficiency. Although this is much more difficult to confirm empirically than bias, it can be assumed from the analyses that negative news have been incorporated too slowly. This impaired the adjustment of the forecasts from earlier optimistic views to new information. As a final source of inaccuracy an imitation effect could be identified. An imitation of the view of other forecasters can be confirmed empirically for the majority of the forecasters and differs highly in its magnitude.

Notwithstanding these common errors which affected the forecast accuracy of all forecasters, the results of a rank-sum test indicate that they differ significantly in their forecast accuracy. This makes the rankings based on different descriptive measures interesting. Still, it remains an open question whether the best performers, who are mainly less renowned forecasters, reach their relatively better results because of their generally better models and abilities, or if they simply cope better with the challenges of this specific period which was apparently very difficult to predict.

## References

- Antholz, B.* (2006): "Geschichte der quantitativen Konjunkturprognose-Evaluation in Deutschland," *Vierteljahrshefte zur Wirtschaftsforschung* 75, 12–33.
- Batchelor, R.A.* (1990): "All forecasters are equal," *Journal of Business and Economic Statistics* 8, 143–144.
- (2001): "How useful are the forecasts of intergovernmental agencies? The IMF and OECD versus the consensus," *Applied Economics* 33, 225–235.
- Batchelor, R. A./Dua, P.* (1992): "Conservatism and Consensus-seeking among Economic Forecasters," *Journal of Forecasting*, 11, 169–181.
- Beck, N./Katz, J. N.* (1995): "What to do (and not to do) with Time-Series Cross-Section Data," *The American Political Science Review* 89, 634–647.

- Bera, A. K. / Jarque, C. M. (1980): "Efficient tests for normality, homoscedasticity and serial independence of regression residuals," *Economics Letters* 6, 255–259.
- Blix, M. / Wadefjord, J. / Wienecke, U. / Ådahl, Martin (2001): "How Good Is the Forecasting Performance of Major Institutions?," *Sveriges Riksbank Economic Review* 2001, 38–68.
- Diebold, F. X. / Mariano, R. S (1995): "Comparing Predictive Accuracy," *Journal of Business and Economic Statistics* 13, 253–264.
- Döpke, J. / Fritsche, U. (2006): "Growth and Inflation Forecasts for Germany. A Panel-based Assessment of Accuracy and Efficiency," *Empirical Economics* 31, 777–798.
- Döpke, J. / Langfeldt, E. (1995): "Zur Qualität von Konjunkturprognosen für Westdeutschland 1976–1994," *Kieler Diskussionsbeiträge* 247, IfW Kiel.
- Fildes, R. / Stekler, H. (2002): "The State of Macroeconomic Forecasting," *Journal of Macroeconomics* 24, 435–468.
- Gallo, G. M. / Granger, C. W. J. / Jeon, Y. (2002): "Copycats and Common Swings: The impact of the use of forecasts in information sets," *IMF Staff Papers* 49, 4–21.
- Grömling, M. (2002): "Evaluation and Accuracy of Economic Forecasts," *Historical Social Research* 27, 242–255.
- Hagen, H. M. / Kirchgässner, G. (1997): "Interest Rate Based Forecast vs. Professional Forecasts of German Economic Growth: A Comparison," *Universität St. Gallen Discussion Paper no. 9714*.
- Harvey, D. I. / Leybourne, S. J. / Newbold, P. (2001): "Analysis of a panel of UK macroeconomic forecasts," *Econometrics Journal* 4, 37–55.
- (1997): "Testing the Equality of Prediction Mean Squared Errors," *International Journal of Forecasting* 13, 281–291.
- Heilemann, U. (2004): "Besser geht's nicht – Genauigkeitsgrenzen von Konjunkturprognosen," *Jahrbücher für Nationalökonomie und Statistik* 224, 51–64.
- Heilemann, U. / Stekler, H. O. (2003): "Has the accuracy of German macroeconomic forecasts improved?," *Discussion Paper of the German Research Council's Research Centre*, 475, 31/03.
- Isiklar, G. / Lahiri, K. (2007): "How Far Ahead Can We Forecast? Evidence from Cross-country Surveys," *International Journal of Forecasting* 23, 167–187.
- Isiklar, G. / Lahiri, K. / Loungani, P. (2006): "How quickly do forecasters incorporate news? Evidence from cross-country surveys," *Journal of Applied Econometrics* 21, 703–725.
- Juhn, G. / Loungani, P. (2002): "Further Cross-Country Evidence on the Accuracy of the Private Sector's Output Forecasts," *IMF Staff Papers* 49, 49–64.
- Kirchgässner, G. / Müller, U. K. (2006): "Are Forecasters Reluctant to Revise their Predictions? Some German Evidence," *Journal of Forecasting* 25, 401–413.
- Loungani, P. (2001): "How Accurate Are Private Sector Forecasts? Cross-Country Evidence from Consensus Forecasts of Output Growth," *International Journal of Forecasting* 17, 419–432.
- McNees, S. K. (1992): "The Uses and Abuses of 'Consensus' Forecasts," *Journal of Forecasting* 11, 703–710.



- Newey, W. K. / West, K. D. (1987): "A simple positive semi-definite heteroskedasticity and autocorrelation-consistent covariance matrix,"* *Econometrica* 55, 703–708.
- Nordhaus, W. D. (1987): "Forecasting Efficiency: Concepts and Applications,"* *The Review of Economics and Statistics* 69, 667–674.
- Öller, L.-E. / Barot, B. (2000): "The accuracy of European growth and inflation forecasts,"* *International Journal of Forecasting* 16, 293–315.
- Stekler, H. O. (1987): "Who forecasts better?,"* *Journal of Business and Economic Statistics* 5, 155–158.

## Appendix

### Data Availability

A problem which appeared when working with the data set at hand concerns mergers and acquisitions making the composition of the Consensus Economics panel change several times. In order to arrive at a continuous time series for the banks whose research departments of banks changed their owner and their name, the consistency of the forecasts and staffs were used as criteria. Applying this, in the cases of the mergers of DG Bank and GZ Bank to DZ Bank or the acquisition of BfG Bank by SEB, continuous time series could be generated. However, in some cases the continuation of the time series was not possible, as in the case of the merger of HYPO Bank and Bayerische Vereinsbank to HypoVereinsbank.

Another problem was caused by missing values because of the discontinuous participation of the respective forecasters in the survey. This was mostly the case when a forecaster simply forgot to report his forecasts to Consensus Economics. If this happens, the previous forecast is not inserted by Consensus Economics, but no value is published. Yet, for some forecasters missing values appeared more generally. This mainly concerned the German research institutes, whose forecasts are available for almost all of the predicted years, but in most cases they show many missing values at single forecast horizons. This reflects the fact that these institutes usually do not revise their forecasts as often as banks do.

Another group of forecasters published regularly for all horizons, but either stopped participating after some years (very often due to an acquisition) or did not start in 1995 but later. In both cases the implementation of several empirical tests for the respective forecasters becomes impossible.

Table 7

## Availability of Forecasters

	Group	Availability	Mergers / Acquisitions
DIW	Institute	07 / 1990 – 12 / 2005	
HWWA	Institute	02 / 1996 – 12 / 2005	
Ifo	Institute	02 / 1990 – 12 / 2005	
RWI	Institute	05 / 1994 – 12 / 2005	
IfW	Institute	02 / 1990 – 12 / 2005	
FAZ Institut	Others	01 / 1992 – 11 / 2005	
Deutsche Bank	Großbank	02 / 1990 – 12 / 2005	
Commerzbank	Großbank	03 / 1990 – 12 / 2005	
Dresdner Bank	Großbank	02 / 1990 – 12 / 2005	
DZ Bank	Co-operative	02 / 1990 – 12 / 2005	Until 12 / 2001: DG Bank
WGZ Bank	Co-operative	02 / 1990 – 12 / 2005	
Bayerische Vereinsbank	–	02 / 1990 – 08 / 1998	
HYPO Bank	–	02 / 1990 – 07 / 1998	
Bankgesellschaft Berlin	Landesbank	02 / 1990 – 12 / 2005	Until 02 / 1994: Berliner Bank
Bayerische Landesbank	Landesbank	02 / 1990 – 12 / 2005	
Westdeutsche Landesbank	Landesbank	02 / 1990 – 12 / 2005	
DekaBank	Landesbank	02 / 1990 – 12 / 2005	Until 12 / 1998: Deutsche Girozentrale
BfG Bank	Affiliate	02 / 1990 – 12 / 2005	Since 04 / 2001: SEB
BHF-Bank	Affiliate	01 / 1995 – 12 / 2005	
Bank Julius Bär	Private	04 / 1994 – 12 / 2005	
Delbruck & Co	Private	02 / 1990 – 04 / 2003	
Sal Oppenheim	Private	02 / 1990 – 12 / 2005	
MM Warburg	Private	04 / 1993 – 12 / 2005	
HSBC Trinkaus & Burkhardt	Affiliate	02 / 1990 – 12 / 2005	Until 11 / 1998: Trinkaus & Burkhardt
JP Morgan	Foreign	04 / 1994 – 12 / 2005	
Morgan Stanley	Foreign	06 / 1996 – 12 / 2005	
Invesco Bank	Foreign	08 / 1998 – 10 / 2004	
Merrill Lynch	Foreign	07 / 1998 – 11 / 2002	
UBS Warburg	Foreign	05 / 1998 – 12 / 2005	Until 04 / 2000: Warburg Dillon Reed; since 07 / 2003: UBS
HypoVereinsbank	Großbank	09 / 1998 – 12 / 2005	
IW Köln	Others	12 / 1999 – 12 / 2005	
Lehman Brothers	Foreign	03 / 2002 – 12 / 2005	
Industriekreditbank	–	02 / 1990 – 07 / 1992	
UBS Frankfurt	–	05 / 1994 – 03 / 1998	
SMH Bank	–	02 / 1990 – 05 / 1998	
Bank in Liechtenstein	Foreign	02 / 1990 – 08 / 1998	

	Group	Availability	Mergers / Acquisitions
Hoechst AG	Others	02/1990–04/1999	
Economist Intelligence Unit	–	10/2003–12/2005	
Goldman Sachs	–	10/2003–12/2005	
Bank of America	–	10/2003–12/2005	
Global Insight	–	10/2004–12/2005	
Citigroup	–	05/2004–12/2005	

Table 8

**Descriptive Statistics: Current Year Forecasts, 1995–2005**

Forecaster	Theil's U (Consensus)	Theil's U ("no change")	RMSE	MAE	Percentage better than Consensus
1 Foreign 1	0.86	0.36	0.42	0.29	0.55
2 Private 1	0.87	0.40	0.47	0.31	0.55
3 Foreign 2	0.89	0.40	0.53	0.34	0.63
4 Affiliate 1	0.93	0.40	0.45	0.31	0.52
5 Affiliate 2	0.95	0.41	0.48	0.32	0.54
6 Affiliate 3	0.95	0.44	0.50	0.35	0.46
7 Institute 1	1.00	0.48	0.55	0.38	0.45
8 Institute 2	1.00	0.53	0.57	0.38	0.47
9 Großbank 1	1.00	0.46	0.53	0.35	0.46
10 Consensus	1.00	0.45	0.52	0.34	–
11 Großbank 2	1.01	0.46	0.53	0.38	0.39
12 Institute 3	1.04	0.45	0.56	0.35	0.44
13 Landesbank 1	1.04	0.47	0.54	0.40	0.34
14 Foreign 3	1.05	0.46	0.57	0.39	0.34
15 Co-operative 1	1.05	0.50	0.57	0.38	0.38
16 Co-operative 2	1.07	0.48	0.56	0.36	0.50
17 Landesbank 2	1.07	0.49	0.56	0.41	0.32
18 Private 2	1.08	0.48	0.53	0.33	0.56
19 Private 3	1.08	0.51	0.58	0.38	0.35
20 Landesbank 3	1.10	0.50	0.58	0.39	0.45
21 Landesbank 4	1.11	0.51	0.58	0.40	0.40
22 Großbank 3	1.11	0.47	0.58	0.34	0.51
23 Landesbank 5	1.15	0.53	0.61	0.40	0.28
24 Others 1	1.16	0.53	0.61	0.41	0.31
25 Großbank 4	1.18	0.52	0.61	0.42	0.33
26 Private 4	1.20	0.68	0.74	0.51	0.11
27 Others 2	1.21	0.52	0.75	0.53	0.18
28 Foreign 5	1.21	0.49	0.56	0.39	0.38
29 Institute 5	1.23	0.60	0.64	0.46	0.30
30 Institute 4	1.25	0.58	0.68	0.59	0.20
31 naïve	2.05	0.93	1.07	0.94	0.13

Table 9

**Descriptive Statistics: Next Year Forecasts, 1995–2005**

Forecaster	Theil's U (Consensus)	Theil's U ("no change")	RMSE	MAE	Percentage better than Consensus
1 naïve	0.80	0.85	1.11	0.90	0.54
2 Foreign 2	0.84	0.88	1.30	0.97	0.82
3 Foreign 3	0.88	0.96	1.36	1.11	0.68
4 Affiliate 2	0.90	0.94	1.25	1.01	0.57
5 Großbank 3	0.93	1.02	1.41	1.12	0.53
6 Landesbank 1	0.95	0.97	1.30	1.05	0.52
7 Institute 1	0.96	0.96	1.37	1.13	0.51
8 Affiliate 1	0.97	1.02	1.31	1.03	0.52
9 Private 4	0.97	1.13	1.45	1.14	0.52
10 Institute 3	0.98	1.02	1.48	1.17	0.59
11 Private 1	0.99	1.02	1.39	1.09	0.55
12 Institute 4	0.99	0.99	1.37	1.11	0.36
13 Others 1	1.00	1.04	1.39	1.10	0.37
14 Consensus	1.00	1.06	1.38	1.07	–
15 Landesbank 5	1.00	1.05	1.41	1.11	0.44
16 Private 2	1.00	1.04	1.35	1.04	0.42
17 Institute 2	1.00	1.01	1.42	1.15	0.46
18 Foreign 1	1.02	1.06	1.40	1.02	0.51
19 Co-operative 1	1.03	1.05	1.41	1.13	0.29
20 Landesbank 2	1.04	1.09	1.45	1.12	0.40
21 Co-operative 2	1.06	1.09	1.45	1.09	0.47
22 Großbank 2	1.06	1.11	1.48	1.20	0.29
23 Landesbank 3	1.07	1.11	1.49	1.14	0.41
24 Others 2	1.07	1.10	1.80	1.55	0.20
25 Großbank 4	1.07	1.10	1.47	1.22	0.21
26 Großbank 1	1.08	1.10	1.48	1.07	0.50
27 Affiliate 3	1.09	1.13	1.52	1.23	0.22
28 Landesbank 4	1.09	1.13	1.51	1.20	0.26
29 Private 3	1.11	1.13	1.51	1.20	0.23
30 Institute 5	1.12	1.22	1.57	1.29	0.15
31 Foreign 5	1.16	1.07	1.41	1.10	0.30

Table 10

**Test for Imitation Behaviour**

	Beta	Gamma	Delta
Institute 1	0.56*** (0.072)	0.42*** (0.071)	-1.47*** (0.475)
Institute 2	0.75*** (0.078)	0.25*** (0.076)	0.49 (0.367)
Institute 3	0.47*** (0.107)	0.58*** (0.109)	-0.52 (0.460)
Institute 4	0.78*** (0.052)	0.23*** (0.052)	-0.25 (0.337)
Institute 5	0.74*** (0.066)	0.30*** (0.068)	-0.09 (0.375)
Großbank 1	0.90*** (0.052)	0.11* (0.060)	-0.36 (0.318)
Großbank 2	0.80*** (0.058)	0.22*** (0.060)	0.02 (0.263)
Großbank 3	0.78*** (0.095)	0.23** (0.105)	0.06 (0.707)
Großbank 4	0.81*** (0.056)	0.20*** (0.055)	-0.06 (0.282)
Co-operative 1	0.60*** (0.062)	0.41*** (0.061)	-0.08 (0.216)
Co-operative 2	0.62*** (0.058)	0.40*** (0.060)	-0.22 (0.301)
Landesbank 1	0.59*** (0.057)	0.39*** (0.051)	-0.47* (0.244)
Landesbank 2	0.63*** (0.052)	0.40*** (0.052)	-0.83*** (0.267)
Landesbank 3	0.78*** (0.056)	0.24*** (0.062)	-0.23 (0.301)
Landesbank 4	0.77*** (0.054)	0.24*** (0.053)	-0.08 (0.280)
Landesbank 5	0.57*** (0.053)	0.44*** (0.053)	-0.61** (0.243)
Affiliate 1	0.68*** (0.060)	0.35*** (0.061)	-0.49** (0.206)
Affiliate 2	0.77*** (0.037)	0.26*** (0.039)	-0.76*** (0.256)
Affiliate 3	0.87*** (0.053)	0.15** (0.067)	-0.03 (0.380)
Private 1	0.84*** (0.062)	0.19*** (0.065)	-0.19 (0.262)
Private 2	0.81*** (0.052)	0.20*** (0.054)	-0.79*** (0.291)
Private 3	0.63*** (0.058)	0.41*** (0.062)	-0.19 (0.272)
Private 4	0.71*** (0.046)	0.30*** (0.043)	-0.34 (0.302)

*Continued Table 10*

	Beta	Gamma	Delta
Foreign 1	0.92*** (0.052)	0.08 (0.056)	-0.44 (0.349)
Foreign 2	0.82*** (0.063)	0.20*** (0.068)	-1.14** (0.441)
Foreign 3	0.66*** (0.090)	0.32*** (0.086)	-1.35** (0.565)
Foreign 4	0.90*** (0.061)	0.12* (0.067)	-0.41 (0.506)
Foreign 5	0.96*** (0.057)	0.03 (0.072)	-0.01 (0.439)
Foreign 6	0.70*** (0.098)	0.20** (0.088)	0.63 (0.620)
Foreign 7	0.60*** (0.122)	0.42*** (0.143)	-0.67 (0.859)
Others 1	0.64*** (0.048)	0.40*** (0.048)	-0.44** (0.202)
Others 2	0.51*** (0.065)	0.51*** (0.061)	-0.10 (0.320)
Others 3	0.52*** (0.105)	0.53*** (0.123)	-0.43 (0.511)

The symbol \*\*\*, \*\*, and \* denotes rejection of the null hypothesis at the 1%, 5%, and 10% level respectively. Standard errors in parentheses.

Table 11: Rank Sum Test

	Institute 1	Institute 2	Institute 4	Institute 5	Groß-bank 1	Groß-bank 2	Groß-bank 4	Co-operative 1	Co-operative 2	Landes-bank 1	Landes-bank 2	Landes-bank 3	Landes-bank 4	Landes-bank 5	Private 1	Private 2	Private 3	Affiliate 1	Affiliate 2	Affiliate 3	Foreign 5	Others 1
1995	1	7	21	17	11	16	5	10	18	2	12	6	13	8	3	15	14	9	4	19	22	20
1996	8	20	14	21	6	17	10	9	16	3	7	5	18	1	2	4	15	12	11	22	19	13
1997	22	4	15	6	2	7	19	11	1	21	8	20	14	16	10	3	17	12	13	9	18	5
1998	8	6	17	19	13	15	21	12	1	11	4	10	18	16	7	14	5	3	2	20	22	9
1999	6	10	9	22	12	17	20	18	16	7	2	5	4	14	11	1	13	3	21	19	8	15
2000	14	10	20	6	2	4	13	19	7	21	18	17	15	16	5	12	1	11	9	22	8	3
2001	2	4	3	15	20	6	19	14	10	5	7	21	11	13	9	16	17	8	1	18	22	12
2002	9	2	16	14	11	5	13	7	22	8	18	20	17	21	12	15	19	4	1	6	3	10
2003	4	13	18	15	17	20	10	14	11	5	21	12	22	9	8	7	16	3	2	1	19	6
2004	14	3	13	15	22	21	18	2	12	1	10	9	16	4	6	11	7	8	20	17	19	5
2005	14	11	5	12	3	18	7	9	4	17	15	2	16	6	10	20	22	8	1	19	21	13
Sum	102	90	151	162	119	146	155	125	118	101	122	127	164	124	83	118	146	81	85	172	181	111

Expected Sum = 126.5.

Reproduced with permission of the copyright owner. Further reproduction prohibited without permission.