

Scoring Rules and Survey Density Forecasts

Gianna Boero, Jeremy Smith and Kenneth F. Wallis*

Department of Economics
University of Warwick
Coventry CV4 7AL, UK

Revised, August 2009

Abstract This paper provides a practical evaluation of some leading density forecast scoring rules in the context of forecast surveys. We analyse the density forecasts of UK inflation obtained from the Bank of England's Survey of External Forecasters, considering both the survey average forecasts published in the Bank's quarterly *Inflation Report*, and the individual survey responses recently made available to researchers by the Bank. The density forecasts are collected in histogram format, and the ranked probability score (RPS) is shown to have clear advantages over other scoring rules. Missing observations are a feature of forecast surveys, and we introduce an adjustment to the RPS, based on the Yates decomposition, to improve its comparative measurement of forecaster performance in the face of differential non-response. The new measure, denoted RPS*, is recommended to analysts of forecast surveys.

Keywords: Density forecast evaluation; Brier (quadratic probability) score; Epstein (ranked probability) score; Logarithmic score; Bank of England Survey of External Forecasters; Missing data; Forecast comparison

Acknowledgments We are grateful to David Bessler and Gabriel Casillas-Olvera for supplying their code, Valentina Corradi and Norman Swanson for helpful suggestions, and Bank of England staff for assembling and helping to clean the survey dataset. Readers wishing to gain access to the data should write to the Publications Editor, Inflation Report and Bulletin Division, Bank of England, Threadneedle Street, London EC2R 8AH, UK.

*Corresponding author [K.F.Wallis@warwick.ac.uk]

1. Introduction

In many forecasting applications the focus of attention is the future value of a continuous random variable, and the presentation of a density forecast or predictive distribution – an estimate of the probability distribution of the possible future values of the variable – is becoming increasingly common. Tay and Wallis (2000) survey early applications in macroeconomics and finance, and more than half of the inflation targeting central banks, worldwide, now present density forecasts of inflation in the form of a fan chart. The best-known series of density forecasts in macroeconomics dates from 1968, when the American Statistical Association and the National Bureau of Economic Research jointly initiated a quarterly survey of macroeconomic forecasters in the United States, known as the ASA-NBER survey; Zarnowitz (1969) describes its original objectives. In 1990 the Federal Reserve Bank of Philadelphia assumed responsibility for the survey and changed its name to the Survey of Professional Forecasters (SPF). Survey respondents are asked not only to report their point forecasts of several variables but also to attach a probability to each of a number of pre-assigned intervals, or bins, into which output growth and inflation, this year and next year, might fall. In this way, respondents provide density forecasts of these two variables, in the form of histograms. The probabilities are then averaged over respondents to obtain survey average density forecasts, again in the form of histograms, which are published. More recently the Bank of England (since 1996) and the European Central Bank (since 1999) have conducted similar surveys with similar density forecast questions, and they also share the practice of the SPF in making the individual responses to the survey, suitably anonymised, available for research purposes. This paper considers methods for the comparative assessment of the quality of such forecasts, with the Bank of England Survey of External Forecasters (SEF) as the practical example. Other aspects of the SEF dataset are explored by Boero, Smith and Wallis (2008a,b,c).

A scoring rule measures the quality of a probability forecast by a numerical score based on the forecast and the eventual outcome, and can be used to rank competing forecasts. The earliest example of such a rule, introduced by Brier (1950) and subsequently bearing his name, concerns the situation in which an event can occur in only one of a small number of mutually exclusive and exhaustive categories, and a forecast consists of a set of probabilities, one for each category, that the event will occur in that category. The Brier score is then given as the sum of the squared differences between the forecast probabilities and an indicator

variable that takes the value 1 in the category in which the event occurred and 0 in all other categories. Much of the theoretical work underpinning probability forecast construction and evaluation originally appeared in the meteorological literature. The example in Brier's article concerned the verification of probability forecasts of rain or no-rain in given periods: this has only two categories and is sometimes called an event probability forecasting problem. The mathematical formulation adopted by Brier has also resulted in the use of the name "quadratic probability score" (QPS), which is used below, although it is potentially misleading, because a family of quadratic scoring rules exists, of which the Brier score is just one member (Stael von Holstein and Murphy, 1978).

When evaluating survey density forecasts, the distinct classes or categories of the Brier score's set-up are taken to be the set of histogram bins. However the ranking or ordering of the bins in terms of the values of the underlying continuous variable is neglected. For a four-bin histogram where the outcome falls in the bin that has been assigned probability 0.3 and the other bins have probability 0.5, 0.1 and 0.1, for example, the Brier score is indifferent to the location of these last three probabilities, but forecasts that had placed 0.5 in a bin adjacent to the bin in which the outcome fell would generally be regarded as better forecasts than those that had not. The Ranked Probability Score introduced by Epstein (1969), a second member of the class of quadratic scoring rules, takes account of the ordering of the categories. It appears not to have been previously used in the present area of research, although its extension to continuous distributions, the continuous ranked probability score (CRPS), has recently attracted attention in the meteorological literature (Gneiting and Raftery, 2007).

Gneiting and Raftery's (2007) review of scoring rules, their characterisations and properties, includes the leading alternative to the quadratic scores, namely the logarithmic score. Originally proposed by Good (1952), this is defined as

$$\log S(f, x_t) = \log f(x_t)$$

for density forecast f of the random variable X_t evaluated at the outcome x_t . The logarithmic score has many attractive features, and appears in the literature in many guises. To a Bayesian the logarithmic score is the predictive likelihood, and if two forecasts are being compared, the log Bayes factor is the difference in their logarithmic scores. The definition in terms of a continuous density is readily adapted to discrete distributions and

discretised continuous distributions, as in the present context, although there is then a potential difficulty. As seen below, it happens from time to time in the individual survey responses that the outcome falls in a histogram bin to which the respondent has assigned zero probability, whereupon the log score is undefined. To assign an arbitrary value to the score on such occasions is an unsatisfactory solution, since the ranking of competing forecasts is sensitive to the chosen value. On the other hand zero-probability forecast outcomes are readily accommodated by the quadratic scores.

In this paper we compare and contrast the Brier and Epstein rules, or QPS and RPS, and the logarithmic score, in applications to survey density forecasts of UK inflation. Section 2 contains the technical background to our study, comprising a formal presentation of the rules, consideration of the relevance of the various decompositions that have been proposed, and the statistical tests of predictive ability that we employ. The empirical analysis begins in Section 3 with a comparison of the published survey average density forecasts from the SEF and the density forecasts of the Bank of England’s Monetary Policy Committee (MPC). Section 4 turns to the individual SEF respondents and uses the two quadratic scoring rules to evaluate their forecast performance: it is seen that the RPS is preferred. Incomplete data are a feature of this survey, like all forecast surveys, and our adjusted score, RPS*, is found to provide more reliable rankings of forecasters in the face of missing observations caused by differential non-response. Section 5 concludes.

2. Scoring rules and their applications

2.1. *The Brier, Epstein and logarithmic rules*

We consider a categorical variable whose sample space consists of a finite number K of mutually exclusive events, and for which a probability forecast of the outcome at time t is a vector of probabilities (p_{1t}, \dots, p_{Kt}) . We have in mind applications in which the categories are the K bins of a histogram of a continuous random variable X , and we define indicator variables d_{kt} , $k = 1, \dots, K$, which take the value 1 if the outcome x_t falls in bin k , otherwise $d_{kt} = 0$. Also in mind are time series forecasting applications, in which each forecast of the outcome at times $t = 1, \dots, T$ is formed at some previous time. For a sample of forecasts and realisations of the categorical variable, the sample mean Brier score is given as

$$\text{QPS} = \frac{1}{T} \sum_{t=1}^T \sum_{k=1}^K (p_{kt} - d_{kt})^2. \quad (1)$$

It is said to have a negative orientation – smaller scores are better. The range is usually stated as $0 \leq \text{QPS} \leq 2$, although the extreme values are obtained in extreme circumstances in which, in every period, all the probability is assigned to a single bin and the outcome either does or does not fall into it. More generally, there is a non-zero lower bound that corresponds to a best fit. If the bin probabilities are constant over time, $p_{kt} = p_{ks} = p_k$, say, $t \neq s$, $k = 1, \dots, K$, this is obtained for any forecast sequence in which the relative bin frequencies

$$\bar{d}_k = \frac{1}{T} \sum_{t=1}^T d_{kt}, \quad k = 1, \dots, K$$

match the probabilities p_k , whereupon the score achieves its minimum value

$$\text{QPS}_{\min} = 1 - \sum_{k=1}^K p_k^2.$$

Note that this is indifferent to the ordering of the time series of observations.

The Brier score is also indifferent to the fact that, in the histogram context, there is a natural ordering of the categories, or an implicit measure of the distance between them, which should be taken into account, as noted above. To do this, Epstein's (1969) proposal replaces the density functions implicit in the Brier score with their corresponding distribution functions (Murphy, 1971). Defining these as

$$P_{kt} = \sum_{j=1}^k p_{jt}, \quad D_{kt} = \sum_{j=1}^k d_{jt}, \quad k = 1, \dots, K,$$

with $P_{Kt} = D_{Kt} = 1$, the ranked probability score is

$$\text{RPS} = \frac{1}{T} \sum_{t=1}^T \sum_{k=1}^K (P_{kt} - D_{kt})^2. \quad (2)$$

The RPS penalises forecasts less severely when their probabilities are close to the actual outcome, and more severely when their probabilities are further from the actual outcome. Like the Brier score, its minimum value is 0, occurring in the same extreme circumstance of the outcomes falling in bins whose forecast probability is 1. Similarly, the maximum value of the RPS occurs when some $p_{kt} = 1$ and the outcome falls in a different bin, but the actual value depends on how far from the k th bin that is. In extremis, with the outcomes and the unit-probability bins located at opposite ends of the range, this value is $K - 1$.

Adapting the definition of the logarithmic score given above to the histogram context gives

$$\log S = \frac{1}{T} \sum_{t=1}^T \sum_{k=1}^K d_{kt} \log(p_{kt}).$$

This has a positive orientation – larger scores are better, and since $\log S$ typically takes negative values, scores with smaller absolute values are typically better.

2.2. *Decompositions of the quadratic scores*

Several decompositions or partitions of the Brier score and, by extension, the Epstein score have been proposed, with the aim of obtaining information about different aspects of forecast performance. Early contributions focused on the event probability forecasting problem and used a simplified version of the Brier score given in equation (1), which we denote QPSE, namely

$$\text{QPSE} = \frac{1}{T} \sum_{t=1}^T (p_t - d_t)^2. \quad (3)$$

Here p_t is the forecast probability, and $d_t = 1$ if the event occurs or zero if it does not. The QPSE score is equal to half of the value obtained from equation (1) with $K = 2$, since it neglects the complementary non-occurrence of the event, whose forecast probability is $1 - p_t$.

Sanders (1963) requires that all probabilities be expressed in tenths and partitions the T forecasts into eleven subsets of size T_j , say, in which the forecast probability is

$p_j = j/10, j = 0, \dots, 10$. To consider QPSE subset-by-subset we rearrange the summation in equation (3) as

$$\text{QPSE} = \frac{1}{T} \sum_{j=0}^{10} \sum_{t \in T_j} (p_j - d_{jt})^2.$$

Expanding the terms in the inner summation gives

$$\sum_{t \in T_j} (p_j - d_{jt})^2 = T_j (p_j - \bar{d}_j)^2 + \sum_{t \in T_j} (d_{jt} - \bar{d}_j)^2 = T_j \left[(p_j - \bar{d}_j)^2 + \bar{d}_j (1 - \bar{d}_j) \right],$$

where \bar{d}_j is the relative frequency of occurrence of the event over the T_j occasions on which the forecast probability is p_j . Summing the first term on the right-hand side over j and dividing by T gives the component of QPSE that measures what is variously called validity,

reliability or calibration. A plot of \bar{d}_j against p_j is called a reliability diagram or calibration curve: for a “well-calibrated” forecaster this is close to a diagonal line. The sum over j of the second term on the right-hand side, divided by T , involves only the outcome indicators but nevertheless reflects forecaster behaviour, because the indicators are sorted into classes according to the forecaster’s probabilities. Sanders (1963) refers to this term as a measure of the “sharpness” of the forecasts, using a term introduced by Bross (1953, Ch.3); “resolution” and “refinement” are also in use. Its maximum value is obtained when each \bar{d}_j is 0.5, that is, the forecaster’s probabilities have not succeeded in discriminating high-probability and low-probability occurrences of the event, and sharpness is lacking.

The second term in Sanders’ decomposition can be further partitioned as

$$\frac{1}{T} \sum_{j=0}^{10} T_j \bar{d}_j (1 - \bar{d}_j) = \bar{d} (1 - \bar{d}) - \frac{1}{T} \sum_{j=0}^{10} T_j (\bar{d}_j - \bar{d})^2,$$

where \bar{d} is the overall rate of occurrence of the event (Murphy, 1973). This separates out the variance or uncertainty of the indicator variable, $\bar{d} (1 - \bar{d})$, which depends only on nature’s determination of the occurrence or otherwise of the event. Murphy argues that the remainder can then more appropriately be called resolution, since it measures the degree to which the relative frequencies for the 11 subcollections of forecasts differ from the overall relative frequency of occurrence of the event: high resolution improves (lowers) the QPS.

This three-component decomposition is used in a study of the Bank of England Monetary Policy Committee’s density forecasts of inflation and growth by Galbraith and van Norden (2008). An event probability forecast is derived from a published density forecast by calculating the forecast probability that the variable in question exceeds a given threshold. The resulting probabilities take continuous values, rather than the discrete values assumed in the preceding derivations, and one could simply group the probabilities into bins. Instead, Galbraith and van Norden use a kernel estimator to obtain a smoothed calibration curve.

Calculating above-threshold and below-threshold probabilities from a density forecast in effect reduces the MPC’s density forecast, which has the two-piece normal functional form, to a two-bin histogram. The Bank’s forecast survey questionnaire most often specifies a six-bin histogram, and generalisations of these decompositions of the QPS for $K > 2$ are

available. However, they depend on similar discretisation and grouping of the forecasts to that used in the above derivations, although with six categories and probabilities stated in tenths (or similarly rounded) the number of possible forecasts is 3003, from Murphy's (1972) equation (1). Many of these possible configurations are of little practical relevance to the SEF individual dataset, where the forecast histograms are almost invariably unimodal, although the tail probabilities in the first and/or last open-ended bins are sometimes sufficiently large to give the impression of an additional local peak. Nevertheless the number of distinct configurations observed in the SEF histograms analysed in Section 4 is typically close to the time series sample size, and a decomposition of the individual scores into reasonable estimates of forecast reliability and resolution is not practicable.

A decomposition of the QPS which does not require such a grouping of forecasts into distinct subcollections is the covariance decomposition due to Yates (1982, 1988), obtained as follows:

$$\begin{aligned} \frac{1}{T} \sum_{t=1}^T \sum_{k=1}^K (p_{kt} - d_{kt})^2 &= \sum_{k=1}^K \frac{1}{T} \sum_{t=1}^T [(p_{kt} - \bar{p}_k) - (d_{kt} - \bar{d}_k) + (\bar{p}_k - \bar{d}_k)]^2 \\ &= \sum_{k=1}^K \left[\text{var}(p_k) + \text{var}(d_k) + (\bar{p}_k - \bar{d}_k)^2 - 2\text{cov}(p_k, d_k) \right]. \end{aligned} \quad (4)$$

Yates (1988) notes that the second term in this last expression, the sum of the outcome indicator sample variances $\text{var}(d_k) = \bar{d}_k(1 - \bar{d}_k)$, is outside the forecaster's influence, while the third term, the sum of squared "biases", indicates the miscalibration of the forecasts. He offers further algebraic rearrangement of the first and fourth terms, as in the initial event-probability derivation with $K = 2$ (Yates, 1982), although their interpretations do not readily generalise to the case $K > 2$.

The Yates decomposition is reported by Casillas-Olvera and Bessler (2006) in their comparative study of the MPC and SEF survey average density forecasts, which we extend in the next section. The contribution of the variance of d to the total QPS varies over subperiods, but is the same for the two forecasts under consideration, as indicated by the above derivation. When working with the forecasts supplied by individual respondents to the survey, however, we face the familiar problem of individual non-response, which differs across individuals, so that the data have the form of an unbalanced panel. Thus the individual scores are calculated over different subsamples of the maximum possible T observations, and

it is no longer the case that the contribution of the variance of d is the same for all individual forecasters. Since this term remains outside the forecasters' influence, to focus on their forecast performance in a comparable manner we standardise the contribution of the variance of d to their individual scores. In order that the score of an individual with no missing observations is unaltered, we do this by replacing the individual subsample outcome variance component of the QPS by the full-sample outcome variance.

We are not aware of a comparable covariance decomposition of the RPS, although the (mostly meteorological) literature contains much discussion of extensions of the previous reliability-resolution-uncertainty decomposition to the RPS and its continuous generalisation (Candille and Talagrand, 2005, for example). Nevertheless it is clear that the derivation in equation (4) applies equally well to the RPS given in equation (2), on replacing lower-case p and d by upper-case P and D . As a result, a similar variance of D term can be identified that is a function of the outcomes alone. For comparing forecast performance in the face of differential non-response we calculate an adjusted score, denoted RPS*, obtained by replacing the individual-specific measure of outcome variance in the RPS by its full-sample equivalent.

2.3. *Testing predictive ability*

To construct a formal test of equal predictive ability of two competing forecasts we follow Giacomini and White (2006). Their framework encompasses point, interval and density forecasts and a wide range of loss functions, and can be readily adapted to the present context, although it is an asymptotic test and our sample size is small. Adapting and simplifying the notation of their equations (4) and (6), the Wald-type test statistic is given as

$$W = T \left(T^{-1} \sum_{t=1}^T h_t \Delta L_t \right)' \Omega^{-1} \left(T^{-1} \sum_{t=1}^T h_t \Delta L_t \right)$$

where h_t is a $q \times 1$ vector of test functions, ΔL_t is the difference in the loss or score of the two forecasts in period t , and Ω is an appropriate heteroscedasticity and autocorrelation consistent (HAC) estimate of the asymptotic covariance matrix. The null hypothesis is

$$E(h_t \Delta L_t) = 0, \quad t = 1, \dots, T,$$

and this is rejected whenever W exceeds the relevant critical value of the χ_q^2 distribution.

The simplest test of equal forecast performance has $q = 1$ and $h_t = 1, t = 1, \dots, T$; in this case,

when using the logarithmic score, the test is equivalent to the (unweighted version of) the likelihood ratio test of Amisano and Giacomini (2007).

3. The SEF average and MPC density forecasts of inflation

Empirical study of the scoring rules begins with comparative evaluation of the average density forecasts of inflation, two years ahead, from the Survey of External Forecasters, and the Monetary Policy Committee's fan chart forecasts of inflation for the same horizon. Both forecasts are published in the Bank of England's quarterly *Inflation Report*, although to obtain numerical values of the parameters of the two-piece normal distribution on which the MPC's fan charts are based, it is necessary to consult the Bank's spreadsheets.

The Bank of England's quarterly Survey of External Forecasters began in 1996. The institutions covered in the survey include City firms, academic institutions and private consultancies, and are predominantly based in London. The sample changes from time to time as old respondents leave or new survey members are included, and not every institution responds every quarter, nor answers every question. Although there is no record of the response rate, the publication of summary results in the *Inflation Report* always includes the number of responses on which each reported statistic is based; typically this is in the low twenties.

Initially the SEF questionnaire asked for forecasts of inflation in the last quarter of the current and following years. Such questions eventually deliver sequences of fixed-event forecasts, analysed by Boero, Smith and Wallis (2008c), but not quarterly series of fixed-horizon forecasts, which are required for the present exercise. However, in 1998 a third question was added, asking for forecasts two years ahead, and this marks the start of the series analysed here. (At this time a second variable, GDP growth, was also added.) In May 2006 all three questions were switched to a fixed-horizon format, focusing on the corresponding quarter one, two and three years ahead, and so our series is continued via the second question of the new format. In the UK's inflation targeting policy regime, the Government chooses the targeted measure of inflation and its target value, and the SEF has sought forecasts of the same variable, namely the Retail Prices Index excluding mortgage interest payments (RPIX) until the end of 2003, then the Consumer Prices Index (CPI). Thus

forecasts collected in the eight quarters through 2002-3 have to be evaluated against outcomes in 2004-5 for the previous target variable, not the then-current target variable. At the time of writing, outcome data are available to the end of 2008, hence we use the surveys from 1998Q1 to 2006Q4, a total of 36. The histograms in the first five of these surveys have four bins (<1.5 , $1.5-2.5$, $2.5-3.5$, >3.5), then the two interior bins were further divided and from 1999Q2 there are six bins (<1.5 , $1.5-2$, $2-2.5$, $2.5-3$, $3-3.5$, >3.5); finally in 2004Q1 the whole grid was shifted downwards by 0.5, following the change in the target from 2.5% RPIX inflation to 2% CPI inflation. For comparative purposes we convert the MPC's fan chart forecasts at the two-year horizon to sets of probabilities for the same bins, using the MPC's parameterisation of the two-piece normal distribution (Wallis, 2004, Box A).

Table 1. SEF average and MPC density forecasts of inflation: scores and test results

	SEF	MPC	p -value, $h_t = 1$	p -value, $h'_t = (1, \pi_{t-1})$
QPS	0.711	0.759	0.171	0.118
RPS	0.566	0.596	0.519	0.395
Log score	-1.465	-1.535	0.451	0.336

Note: $T = 36$ (two-year-ahead forecasts published 1998Q1-2006Q4)

The scores of the two forecasts are shown in Table 1. It is seen that the survey average forecast has a smaller QPS than the MPC forecast, which corresponds to Casillas-Olvera and Bessler's (2006) finding based on the first 14 of these 36 quarterly observations. The RPS and the log score give the same ranking of the two forecasts. The RPS values are smaller than the QPS values, because the forecast densities are unimodal and, most of the time, the outcomes fell towards the centre of these distributions: the positioning of relatively high probabilities close to the bin in which the outcome fell is acknowledged by the RPS, but not by the QPS.

The last two columns of Table 1 show the asymptotic p -values of the Giacomini-White test using two test functions, the first an intercept, the second also including the most recent inflation observation at the time the forecast was made, denoted π_{t-1} . We use a Newey-West estimate of Ω allowing for a moving average of order 8 in the forecast errors.

The results show that, although the three scores agree on the ranking of the two forecasts, in no case is the difference in the scores sufficiently great to reject the hypothesis of equal forecast performance. The p -value for each score is reduced when lagged inflation is added to the test function, providing weak evidence of differential use of this information by the two forecasts, but not sufficient to overturn the general conclusion.

To study the comparative behaviour of the quadratic scores in greater detail we turn to Figure 1, which illustrates observation-by-observation the components of the calculation of QPS and RPS, namely the histogram probabilities and the location of the inflation outcome, for the SEF average forecast and MPC forecast in the upper and lower panels respectively. The coloured segments of the vertical columns show, with reference to the left-hand scale, the allocation of forecast percentage probabilities to the histogram bins. For most of the period there are six bins, and the colours follow a rainbow array. The key to the figure records the RPIX inflation range for each bin; from 2004Q1 all these numbers should be reduced by 0.5, following the switch to CPI inflation. For the first five observations there are four bins, with the two interior bins combining, pairwise, the four interior bins of the six-bin grid, as described above: their colours are intermediate, in the same spectral sense, between the separate colours of their corresponding pairs. The large black dots show in which bin the inflation outcome, two years later, fell. There is no inflation scale in Figure 1, and the dots are simply placed in the centre of the probability range of the appropriate bin; this is the same bin for both forecasts, since we have calculated the MPC's probabilities as if the MPC was answering the SEF questionnaire, as noted above. Readers wishing to see a plot of actual inflation outcomes should consult Figure 2. The QPS and RPS for each observation are shown with reference to the right-hand scale; these points are joined by solid and dashed lines respectively, and their mean values over the 36 observations are reported in Table 1.

For most of the period, the inflation outcomes fell in one of the two central bins of the histograms, and the RPS is smaller than the QPS because it correctly acknowledges the appropriate unimodal shape of the densities, for both forecasts. The SEF scores are generally smaller than the MPC scores in these circumstances, because the SEF densities have smaller dispersion. However the last three forecasts provide an interesting contrast. The outcomes, with CPI inflation in excess of 3%, fell in the upper open-ended bin, and the MPC's greater tail probabilities result in its lower scores. The difference with the SEF is more marked in the case of the RPS, where the MPC correctly benefits from greater probabilities not only in the

upper bin, but also in the adjoining bin. However these three observations are not sufficient to offset the overall lower scores of the SEF average forecasts, as indicated by the sample means in Table 1. Nevertheless these different episodes illustrate the advantage of the RPS in better reflecting probability forecast performance in categorical problems which have a natural ordering, such as these density forecast histograms, and its continued use is recommended.

The inclusion of Figure 2 for the benefit of readers who are unfamiliar with the UK's inflationary experience over this period also allows us to relate a further comparison between the SEF average forecasts and the MPC's forecasts. Figure 2 shows the inflation outcomes, 2000Q1-2008Q4, together with point forecasts made two years earlier, namely the MPC density forecast means as published on the Bank's spreadsheets and the corresponding means calculated from the SEF average histograms. (We apply the standard formula assuming that the open-ended bins have twice the width of the interior bins; it makes no difference whether the probabilities are assumed to be concentrated at the mid-points of the respective bins, or spread uniformly across each bin.) The general tendency of the external forecasts to stay close to the inflation target irrespective of the inflation experience at the time the forecasts were made is often taken to be an indication of the credibility of the MPC and the inflation targeting policy regime. Viewed simply as forecasts, however, as in the analysis of the MPC's forecasts by Groen, Kapetanios and Price (2009), we find that their respective forecast RMSEs are 0.65 (MPC) and 0.61 (SEF), which matches the ranking of these forecasts given in Table 1 by the scoring rules.

4. Scoring the individual SEF respondents

4.1. QPS and RPS for regular respondents

The dataset of individual SEF responses made available by the Bank of England gives each respondent an identification number, so that their individual responses, including non-response, can be tracked over time, and their answers to different questions can be matched. The total number of respondents appearing in the dataset is 48, but there has been frequent entry and exit, as in other forecast surveys, and no-one has answered every question since the beginning. To avoid complications caused by long gaps in the data, and to maintain degrees of freedom at a reasonable level, we follow the practice of US SPF researchers and conduct

our analyses of individual forecasters on a subsample of regular respondents. For the present purpose we define “regular” as “more than two-thirds of the time”, which gives us a subsample of 16 respondents, who each provided between 25 (two respondents) and 36 (one respondent) of the 36 possible two-year-ahead density forecasts of inflation over the 1998Q1-2006Q4 surveys.

Although the survey average forecasts always have non-zero probabilities in every bin, as seen in Figure 1, many individual forecasts use fewer of the available bins. Moreover, it happens in 12 individual forecasts, made by four respondents, that inflation falls in a bin which has forecast probability of zero, hence for these four respondents a logarithmic score cannot be calculated. Rather than exclude these four respondents from further consideration, we prefer to set the logarithmic score aside, and in this section consider only the quadratic scores.

We first extend the QPS-RPS comparison of the previous section to the individual forecasters. For each regular respondent both scores are calculated from their available forecasts and outcomes, thus T in equations (1) and (2) takes values between 25 and 36 for the different respondents. A scatter diagram of the results is presented in Figure 3, which also includes the SEF average density forecast as a point of reference, plotted at the values given in Table 1. Bearing in mind the difference in scales, it is seen that all 16 points lie below the “45°” line, thus Section 3’s finding for the SEF average forecast that the RPS is less than the QPS extends to these individual forecasters, for the same general reasons discussed above. The scatter of points is positively sloped, although there is less than perfect agreement between the rankings: the rank correlation coefficient between the QPS and RPS of the regular respondents is 0.76. There are several ambiguous pairwise comparisons: whenever the line joining two points has a negative slope, the QPS and RPS disagree about the relative performance of the corresponding forecasters.

For detailed individual scrutiny we first pick out individual 26, who is the only ever-present regular respondent, and is highly ranked (3rd) on both scores, and is an outlier in one further respect. Whereas almost three-quarters of all the individual forecasts in the sample (357 out of 485) utilise all available histogram bins, there are 21 forecasts which have non-zero entries in only two bins, and 17 of these are individual 26’s forecasts. The upper panel of Figure 4 shows the observation-by-observation components of the score calculations for

individual 26 as in Figure 1; on the five occasions when inflation fell in outer bins with zero forecast probabilities, the large black dots are placed on the boundary of the grids. These include two quarters with inflation below 2% (the 2000Q2,Q3 forecasts) and two with inflation above 3% (the 2006Q2,Q4 forecasts). These are of especial interest, because for each of these four observations the QPS takes approximately the same value, in the range 1.50-1.58, suggesting that the four forecasts are of approximately equal quality. On the other hand the RPS gives a well-separated ranking of these forecasts: 2006Q2 is clearly worst, followed by 2006Q4, whereas 2000Q2,Q3 are rather better. Study of the location of the various probabilities forming the histograms shows that this alternative view of the comparative quality of these forecasts is correct, and the QPS's indifference to this question again emphasises its inadequacy as an indicator of the quality of these density forecasts. In the following section we in turn set the QPS aside.

4.2. *Missing data*

For comparison we include in the lower panel of Figure 4 the corresponding data for individual 25, who has the best RPS result, as shown in Figure 3. Although the first seven forecasts do not score as well as those of individual 26, the local peaks in the latter's RPS at the zero-probability outcomes have much diminished counterparts in individual 25's scores. Also very noticeable, however, is that individual 25's last two forecasts are missing, whereas these observations make relatively high contributions to individual 26's overall RPS.

To place such comparisons on an equal basis, one might consider calculating the scores over the subsample of observations common to both forecasters, thus in the above case simply using the first 34 datapoints for both respondents. However this neglects available information on the performance of the forecaster who has responded more often. Moreover to make multiple comparisons among our 16 regular respondents this is not a practical solution. Although none of these respondents is missing more than 11 of the 36 possible forecasts, the incidence of the missing forecasts shown in Figure 5 is such that there are only three occasions when all 16 individual forecasts are available. Overall, 91 of the possible $16 \times 36 = 576$ forecasts are missing, comprising 77 cases of complete non-response to the questionnaire, and 14 cases of an incomplete questionnaire being returned, known as *item non-response* to survey practitioners. Also shown in Figure 5 is the latest inflation data available at the time the forecasts were prepared (π_{t-1} of Section 3), and there is no evidence

of a relation between the process leading to missing forecasts and this variable (nor any other variable we have considered). Accordingly we treat the missing data as *missing at random* and the observed data as *observed at random*, using terms introduced by Rubin (see Little and Rubin, 2002). Neither imputation-based methods nor model-based methods for handling incomplete data, as discussed by Little and Rubin, appear relevant to the present context of forecast comparison, although we note an interesting application to the construction of a combined point forecast in the face of missing data in the US SPF by Capistran and Timmermann (2007).

Instead, as discussed at the end of Section 2, we focus on the components of the score that reflect forecaster performance, by correcting the score for variation in the outcome variance term identified in the Yates decomposition (equation (4), generalised to the RPS). To retain comparability with the uncorrected score, we replace the outcome variance calculated over an individual's subsample by the full-sample outcome variance. Thus the score for individual 26, who has no missing observations, does not change. (To calculate the variance of d or D over the full sample we assume six histogram bins throughout and assign the first five inflation outcomes accordingly, even though those forecasts had only four bins.)

The results are shown in Figure 6, as a scatter diagram of RPS and adjusted RPS (denoted RPS*) values. As in Figure 3, the two scores give different rankings of forecasters, with a rank correlation coefficient of 0.72. Points lying above the 45° line represent individuals whose score has increased as a result of the adjustment, and their previous lower score might be considered to be the result of having missed some hard-to-forecast occasions. This adjective certainly applies to the last three inflation outcomes in our sample, and individuals 2, 25 and 27 did not respond on two of these occasions, while individual 8 missed all three. The adjustment corrects for the smaller outcome variance in their respective subsamples and increases their scores, resulting in a more accurate picture of their relative forecast performance. In particular, the adjustment moves individual 25 from 1st to 4th position in the ranking, and individual 8 from 8th to 14th.

For a final illustration at the individual level we present the data for the two respondents whose scores are decreased most as a result of the adjustment in Figure 7. Individual 9, in the upper panel, has the same number of missing observations – ten – as individual 8, but these correspond to outcomes that fell in the central bins of the histograms.

Thus the subsample outcome variance is greater than the full-sample variance and the adjustment reduces the score. Nevertheless individual 9 remains ranked in last place, as a result of the excessive dispersion of the forecast histograms, in particular the high probabilities attached to forecast outcomes in the lowest, open-ended bin, which did not materialise. On the other hand for individual 31, in the lower panel of Figure 6, who has eleven missing observations similarly distributed, the adjustment changes the ranking, from 6th on RPS to the top ranked position on RPS*. The scores for the four forecasts made between 2002Q3 and 2003Q2 are unusually small, as a result of placing rather high probabilities in the bins into which inflation duly fell, and zeroes in the outer bins. Throughout, unlike individual 9, individual 31 placed small, or zero, probabilities in the lower open-ended bin, and the latter's relative scores benefited from this choice, except in 2000Q2,Q3.

The overall effect of these adjustments for missing data is to reduce the dispersion of the individual scores. Part of the dispersion in the unadjusted scores is seen to be the result of differential non-response, over and above differences in forecasting performance. The $\text{var}(D_k)$ component of the individual RPS given by the Yates decomposition is outside the forecaster's influence, and assuming that this is independent of the factors that result in individual non-response from time to time, the adjusted score RPS* that corrects for the differential impact of this component gives a more reliable ranking of individual forecast performance. There remains considerable dispersion in the RPS* scores, however, and this heterogeneity in individual density forecasting performance mirrors the finding of considerable heterogeneity in point forecasting performance in this survey by Boero, Smith and Wallis (2008b).

5. Conclusion

This paper provides a practical evaluation of some leading density forecast scoring rules in the context of forecast surveys. We analyse the forecasts of UK inflation obtained from the Bank of England's Survey of External Forecasters, considering both the survey average forecasts published in the quarterly *Inflation Report*, and the individual survey responses recently made available by the Bank. The density forecasts are collected in histogram format,

as a set of probabilities that future inflation will fall in one of a small number of preassigned ranges, and thus are examples of categorical forecasts in which the categories have a natural ordering. Epstein's ranked probability score was initially proposed as an alternative to Brier's quadratic probability score for precisely these circumstances, and our exercise makes its advantages clear. The logarithmic score is the leading alternative to the quadratic scoring rules but, unlike them, is not defined whenever inflation falls in a histogram bin to which the forecaster has assigned zero probability. Such occurrences are observed in our sample of individual forecasters, and exclude the logarithmic score from consideration in this context.

Missing observations are endemic in surveys, and our answer to this problem comes in two parts. First, in common with much other research on forecast surveys, our study of individual forecast performance is conducted on a subsample of regular respondents. In our case these are the 16 respondents who are each missing less than one-third of the possible two-year-ahead forecasts collected between 1998Q1 and 2006Q4. Their forecast scores have considerable dispersion, part of which is due to differences in the inflation outcomes over the different subperiods for which these respondents provided their forecasts. Accordingly, and secondly, we introduce an adjustment to the score, based on the Yates decomposition, which corrects for the differential impact of the component of the score that depends only on the outcome and not on the forecast, and hence gives a clearer ranking of forecaster performance. We recommend the adjusted ranked probability score, denoted RPS*, to other analysts of forecast surveys facing the familiar problems of non-response.

Attention in Section 4 of this paper is restricted to descriptive comparisons and rankings of competing forecasts, without formal testing. Extensions of the pairwise test used in Section 3 to multiple comparisons, using Bonferroni intervals or other methods (see Miller, 2006), also keeping in mind the small-sample context, await future research.

Analysis of the point forecasts of inflation and GDP growth from the SEF by Boero, Smith and Wallis (2008b) finds considerable heterogeneity among individual respondents, shown by the failure of standard tests of equality of idiosyncratic error variances and evidence of different degrees of asymmetry in forecasters' loss functions. Similar dispersion of forecast scores from their density forecasts of inflation again indicates that some respondents are better at forecasting than others. This leads us to close this paper with the

same final thought as that article, that the findings “prompt questions about the individual forecasters’ methods and objectives, the exploration of which would be worthwhile”.

References

- Amisano, G. and Giacomini, R. (2007). Comparing density forecasts via weighted likelihood ratio tests. *Journal of Business and Economic Statistics*, 25, 177-190.
- Boero, G., Smith, J. and Wallis, K.F. (2008a). Uncertainty and disagreement in economic prediction: the Bank of England Survey of External Forecasters. *Economic Journal*, 118, 1107-1127.
- _____ (2008b). Evaluating a three-dimensional panel of point forecasts: the Bank of England Survey of External Forecasters. *International Journal of Forecasting*, 24, 354-367.
- _____ (2008c). Here is the news: forecast revisions in the Bank of England Survey of External Forecasters. *National Institute Economic Review*, No.203, 68-77.
- Brier, G.W. (1950). Verification of forecasts expressed in terms of probability. *Monthly Weather Review*, 78, 1-3.
- Bross, I.D.J. (1953). *Design for Decision*. New York:Macmillan.
- Candille, G. and Talagrand, O. (2005). Evaluation of probabilistic prediction systems for a scalar variable. *Quarterly Journal of the Royal Meteorological Society*, 131, 2131-2150.
- Capistran, C. and Timmermann, A. (2007). Forecast combination with entry and exit of experts. *Journal of Business and Economic Statistics*, forthcoming.
- Casillas-Olvera, G. and Bessler, D.A. (2006). Probability forecasting and central bank accountability. *Journal of Policy Modeling*, 28, 223-234.
- Epstein, E.S. (1969). A scoring system for probability forecasts of ranked categories. *Journal of Applied Meteorology*, 8, 985-987.
- Galbraith, J.W. and van Norden, S. (2008). Calibration and resolution diagnostics for Bank of England density forecasts. Presented at the ‘Nowcasting with Model Combination’ Workshop, Reserve Bank of New Zealand, 11-12 December 2008.
- Giacomini, R. and White, H. (2006). Tests of conditional predictive ability. *Econometrica*, 74, 1545-1578.
- Gneiting, T. and Raftery, A.E. (2007). Strictly proper scoring rules, prediction, and estimation. *Journal of the American Statistical Association*, 102, 359-378.

- Good, I.J. (1952). Rational decisions. *Journal of the Royal Statistical Society B*, 14, 107-114.
- Groen, J.J.J., Kapetanios, G. and Price, S. (2009). A real time evaluation of Bank of England forecasts of inflation and growth. *International Journal of Forecasting*, 25, 74-80.
- Little, R.J.A. and Rubin, D.B. (2002). *Statistical Analysis with Missing Data* (2nd edn). Hoboken, NJ: Wiley-Interscience.
- Miller, R. (2006). Multiple comparisons. In *Encyclopedia of Statistical Sciences* 2nd edn, (N. Balakrishnan, C.B. Read and B. Vidakovic, eds), pp.5055-5065. Hoboken, NJ: Wiley-Interscience.
- Murphy, A.H. (1971). A note on the ranked probability score. *Journal of Applied Meteorology*, 10, 155-156.
- _____ (1972). Scalar and vector partitions of the probability score: Part II. *N*-state situation. *Journal of Applied Meteorology*, 11, 1183-1192.
- _____ (1973). A new vector partition of the probability score. *Journal of Applied Meteorology*, 12, 595-600.
- Sanders, F. (1963). On subjective probability forecasting. *Journal of Applied Meteorology*, 2, 191-201.
- Stael von Holstein, C.S. and Murphy, A.H. (1978). The family of quadratic scoring rules. *Monthly Weather Review*, 106, 917-924.
- Tay, A.S. and Wallis, K.F. (2000). Density forecasting: a survey. *Journal of Forecasting*, 19, 235-254. Reprinted in *A Companion to Economic Forecasting* (M.P. Clements and D.F. Hendry, eds), pp.45-68. Oxford: Blackwell, 2002.
- Wallis, K.F. (2004). An assessment of Bank of England and National Institute inflation forecast uncertainties. *National Institute Economic Review*, No.189, 64-71.
- Yates, J.F. (1982). Decompositions of the mean probability score. *Organizational Behavior and Human Performance*, 30, 132-156.
- _____ (1988). Analyzing the accuracy of probability judgments for multiple events: an extension of the covariance decomposition. *Organizational Behavior and Human Decision Processes*, 41, 281-299.
- Zarnowitz, V. (1969). The new ASA-NBER survey of forecasts by economic statisticians. *American Statistician*, 23(1), 12-16.

Figure 1. Forecast probabilities two years ahead, inflation indicators, QPS and RPS
 Upper panel: SEF average forecast; lower panel: MPC forecast

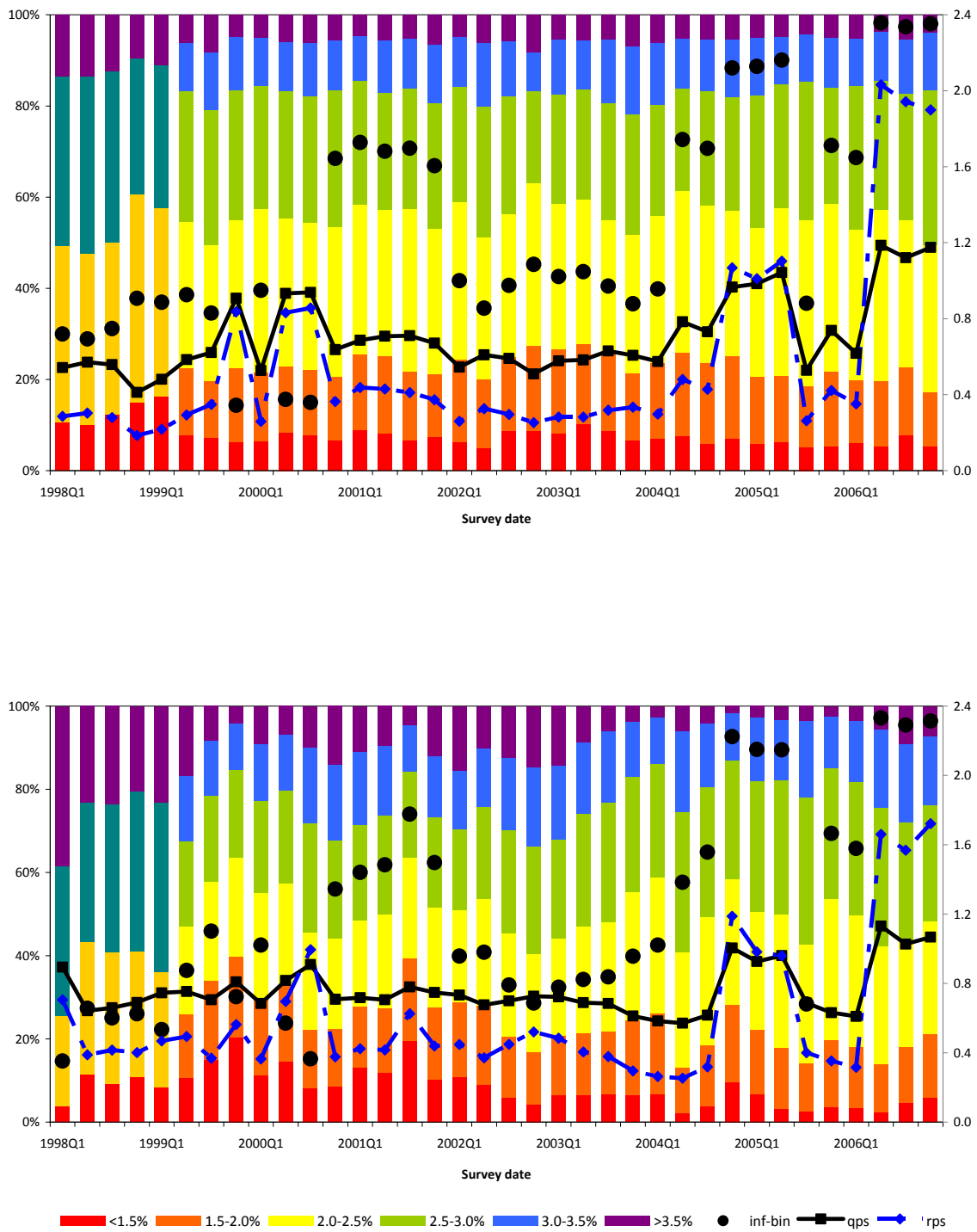


Figure 2. Inflation, 2000Q1-2008Q4, and mean forecasts made two years earlier

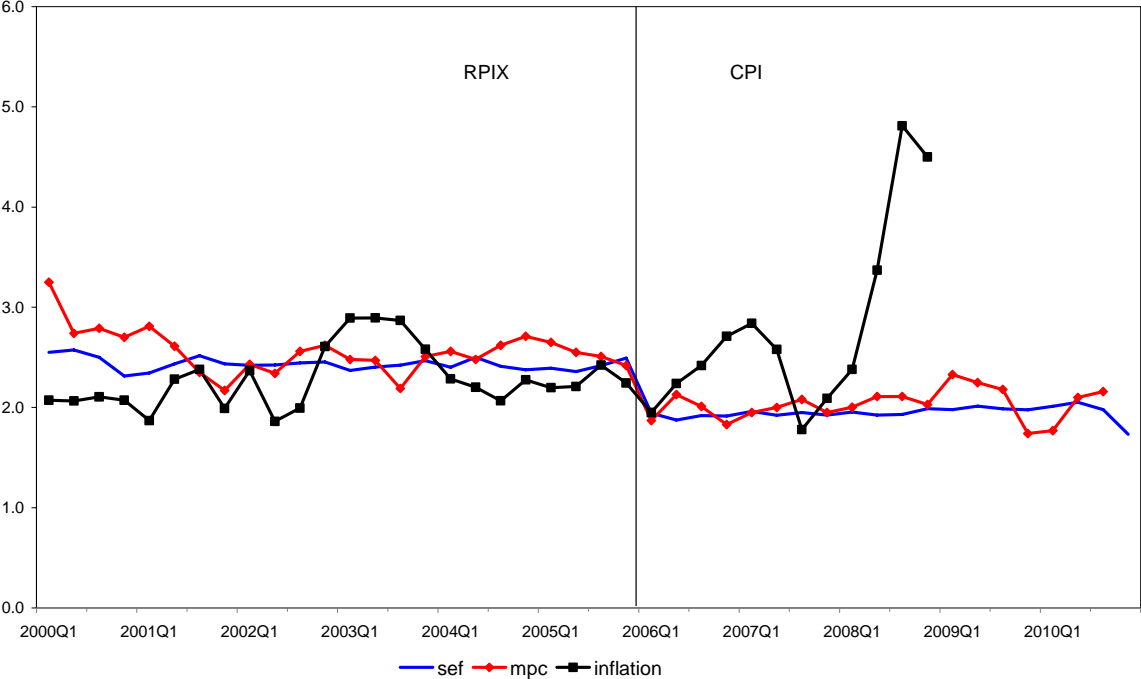


Figure 3. QPS and RPS for 16 regular respondents and the SEF average (filled square)

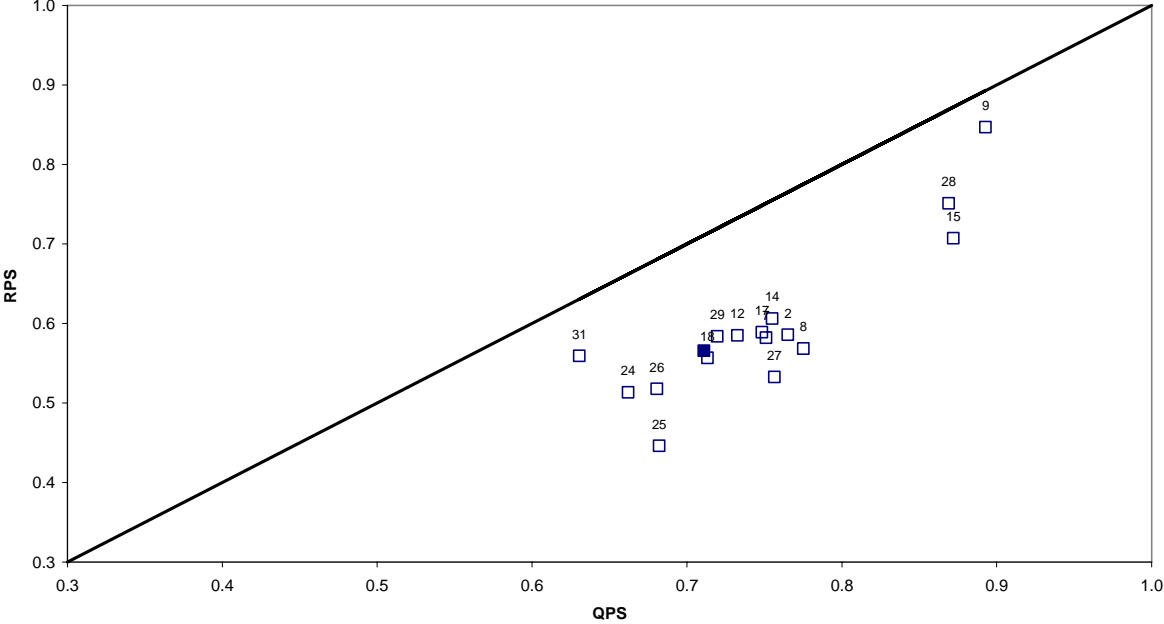
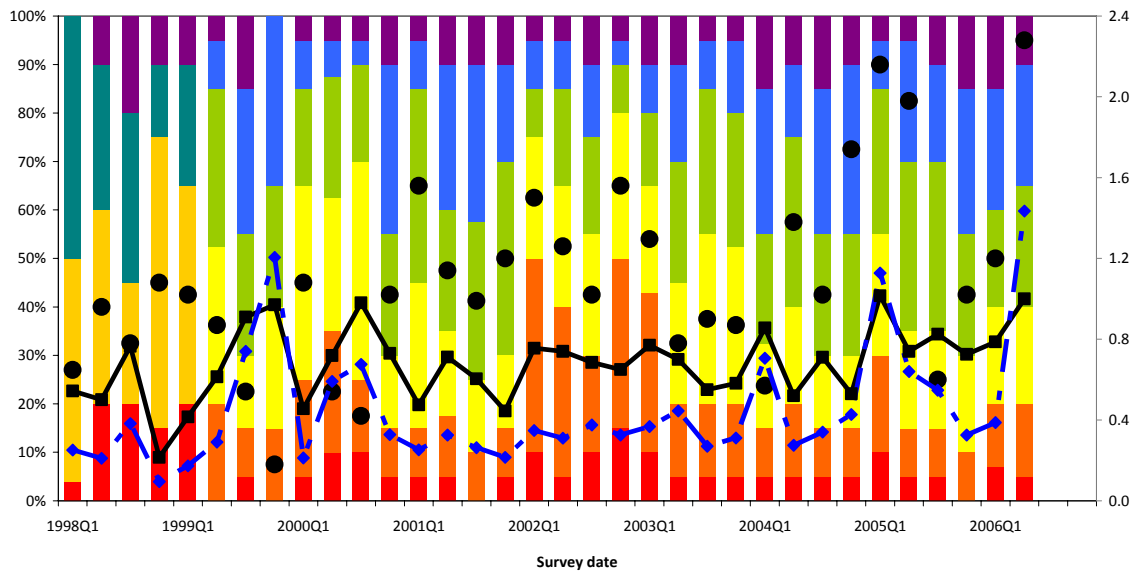
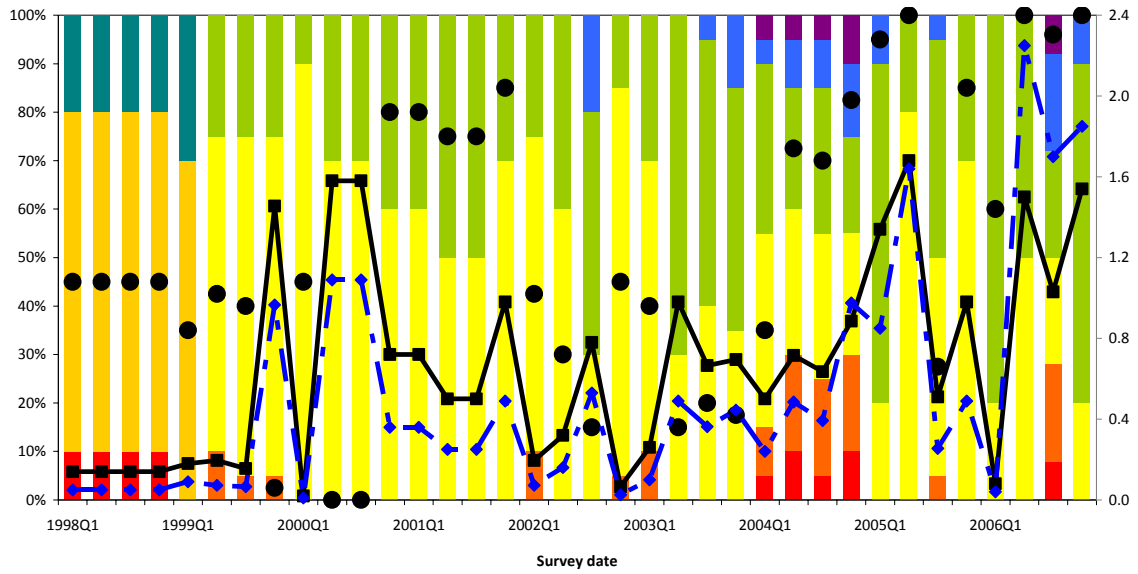


Figure 4. Forecast probabilities two years ahead, inflation indicators, QPS and RPS
 Upper panel: individual 26; lower panel: individual 25



■ <1.5%
 ■ 1.5-2.0%
 ■ 2.0-2.5%
 ■ 2.5-3.0%
 ■ 3.0-3.5%
 ■ >3.5%
 ● inf
 —■— qps
 - - - ◆ - - - rps

Figure 5. Incidence of missing two-year-ahead forecasts (blanks) from 16 regular respondents, and latest inflation data

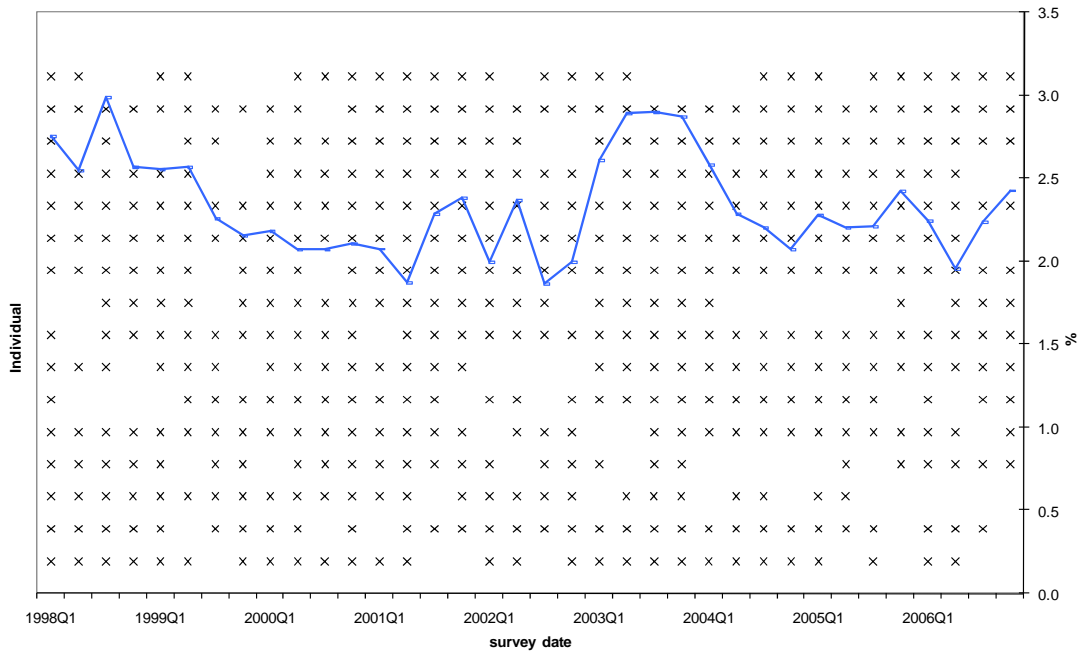


Figure 6. RPS and RPS* for 16 regular respondents

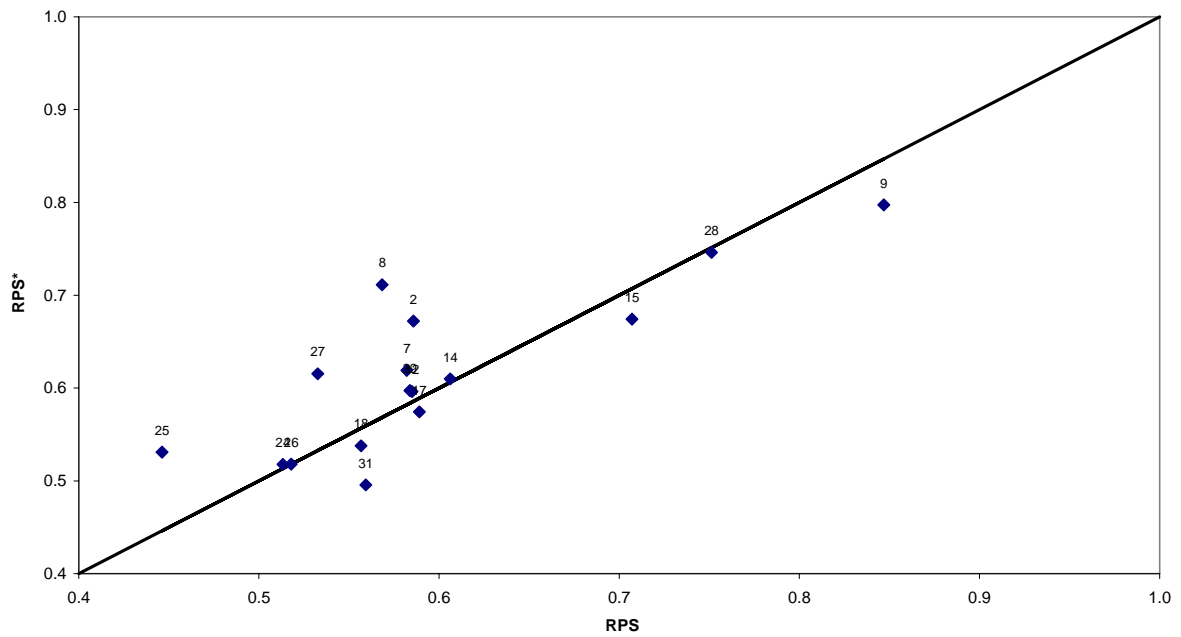


Figure 7. Forecast probabilities two years ahead, inflation indicators, QPS and RPS
 Upper panel: individual 9; lower panel: individual 31

