



The forecast combination puzzle: A simple theoretical explanation



Gerda Claeskens^a, Jan R. Magnus^{b,c}, Andrey L. Vasnev^{d,*}, Wendun Wang^{e,c}

^a KU Leuven, Belgium

^b Vrije Universiteit Amsterdam, The Netherlands

^c Tinbergen Institute, The Netherlands

^d University of Sydney, New South Wales, Australia

^e Econometric Institute, Erasmus University Rotterdam, The Netherlands

ARTICLE INFO

Keywords:

Forecast combination
Optimal weights

ABSTRACT

This paper offers a theoretical explanation for the stylized fact that forecast combinations with estimated optimal weights often perform poorly in applications. The properties of the forecast combination are typically derived under the assumption that the weights are fixed, while in practice they need to be estimated. If the fact that the weights are random rather than fixed is taken into account during the optimality derivation, then the forecast combination will be biased (even when the original forecasts are unbiased), and its variance will be larger than in the fixed-weight case. In particular, there is no guarantee that the 'optimal' forecast combination will be better than the equal-weight case, or even improve on the original forecasts. We provide the underlying theory, some special cases, and a numerical illustration.

Crown Copyright © 2016 Published by Elsevier B.V. on behalf of International Institute of Forecasters. All rights reserved.

1. Introduction

When several forecasts of the same event are available, it is natural to try and find a (linear) combination of these forecasts that is the 'best' in some sense. If we define 'best' in terms of the mean squared error and the variances and covariances of the forecasts are known, then optimal weights can be derived. In practice, though, these (co)variances are not known and need to be estimated. This leads to estimated optimal weights and an estimated optimal forecast combination. Empirical evidence and extensive simulations show that the estimated optimal forecast combination typically does not perform well, and that the

arithmetic mean often performs better. This empirical fact has become known as the 'forecast combination puzzle'.

The history of the puzzle is elegantly summarized by Graefe, Armstrong, Jones, and Cuzán (2014, Section 4), and Smith and Wallis (2009) made a rigorous attempt to explain it, using simulations and an empirical example. They showed that the effect of the error on the estimation of the weights can be large, thus providing an empirical explanation of the forecast puzzle. Smith and Wallis (2009) use the words 'finite-sample' error, which suggests that this error may vanish asymptotically. However, it is not so easy to find an asymptotic justification for ignoring the noise generated by estimating the weights. To begin with, it is not clear what 'asymptotic' means here. What goes to infinity? The number of forecasts? If so, then the number of weights also goes to infinity. The number of observations underlying the total (but finite) set of forecasts? That would make more sense, but it would be difficult to analyze.

* Correspondence to: Office 4160, Abercrombie Building (H70), The University of Sydney Business School, Sydney, NSW 2006, Australia.

E-mail address: andrey.vasnev@sydney.edu.au (A.L. Vasnev).

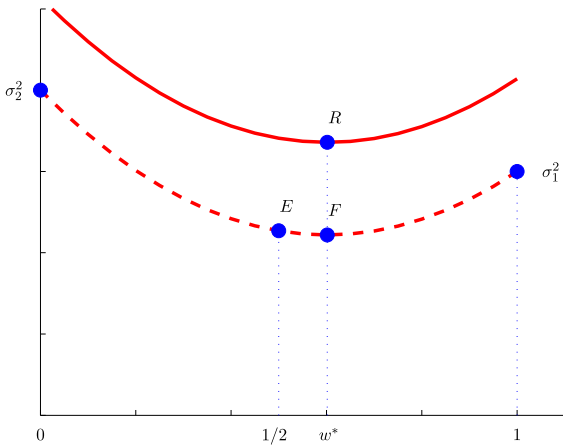


Fig. 1. Variance of forecast combination, in two dimensions: fixed weights (dashed line) and random weights under normality (solid line).

In this paper, we provide a theoretical explanation for the empirical and simulation results of [Smith and Wallis \(2009\)](#) and others. The key ingredient of our approach is the specific acknowledgement that the optimal weights should be derived by taking the estimation step into account explicitly. In other words, we view the derivation and estimation of optimal weights as a joint effort, not as two separate efforts. This approach differs from (almost) all previous research, not only the study by [Bates and Granger \(1969\)](#), but also later contributions, important and insightful though they may be, such as those of [Elliott \(2011\)](#), [Hansen \(2008\)](#), [Hsiao and Wan \(2014\)](#), and [Liang, Zou, Wan, and Zhang \(2011\)](#). The separation of the mathematical derivation and statistical estimation can be quite dangerous. However, even though the disadvantages of such separations have been highlighted, they are still quite common in econometrics, and specifically in the model-averaging literature, which explicitly attempts to combine model selection and estimation, so that uncertainty in the model selection procedure is not ignored when reporting properties of the estimates; see for example [Magnus and De Luca \(2016\)](#).

We highlight our main findings by first providing graphical illustrations of the cases of two forecasts, as analyzed by [Bates and Granger \(1969\)](#). Thus, we linearly combine two forecasts of an event μ :

$$y_c = wy_1 + (1 - w)y_2. \tag{1}$$

If the weight w is considered to be fixed, then the forecast combination is unbiased ($Ey_c = \mu$) if the original forecasts are unbiased, and the variance of the combination will be

$$\text{var}(y_c) = w^2\sigma_1^2 + (1 - w)^2\sigma_2^2 + 2w(1 - w)\rho\sigma_1\sigma_2, \tag{2}$$

where σ_1^2 and σ_2^2 are the variances of y_1 and y_2 respectively, and $\rho = \text{corr}(y_1, y_2)$ denotes the correlation.

The variance is a quadratic function of w , as plotted in [Fig. 1](#) (dashed line). At $w = 0$, we obtain σ_2^2 ; at $w = 1$, we obtain σ_1^2 ; and at $w = 1/2$, we obtain point E . The optimum F is reached at $w = w^*$, the optimal weight that gives the smallest variance of the forecast combination.

Now suppose that the weights are estimated, so that they are random rather than fixed. In the special case

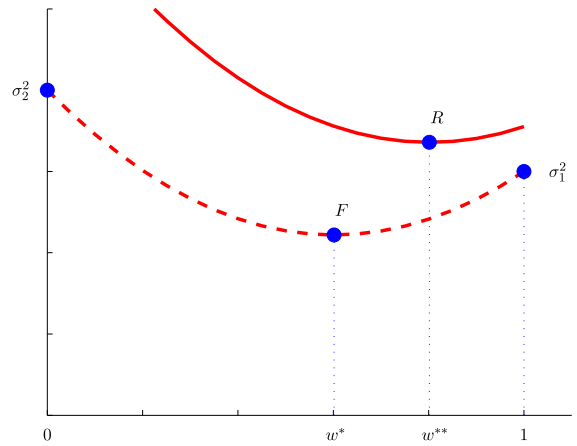


Fig. 2. Variance of forecast combination, in two dimensions: random weights, general case.

where (y_1, y_2, w) follows a trivariate normal distribution, the combination is biased even when the original forecasts are unbiased, since

$$Ey_c = \mu + \text{cov}(w, y_1 - y_2), \tag{3}$$

and the variance is given by

$$\begin{aligned} \text{var}(y_c) = & (Ew)^2\sigma_1^2 + (1 - Ew)^2\sigma_2^2 \\ & + 2(Ew)(1 - Ew)\rho\sigma_1\sigma_2 \\ & + \text{var}(w)\text{var}(y_1 - y_2) + (\text{cov}(w, y_1 - y_2))^2. \end{aligned} \tag{4}$$

In another special case where w is independent of (y_1, y_2) , the combination is unbiased and

$$\begin{aligned} \text{var}(y_c) = & (Ew)^2\sigma_1^2 + (1 - Ew)^2\sigma_2^2 \\ & + 2(Ew)(1 - Ew)\rho\sigma_1\sigma_2 \\ & + \text{var}(w)\text{var}(y_1 - y_2). \end{aligned} \tag{5}$$

In either case, the variance is shifted upwards, as is shown in [Fig. 1](#) (solid line). The solid line gives the variance as a function of Ew , and the optimum is reached at the same point w^* as before, but leading to a higher variance of the forecast combination. We see that, while the equal-weights point at $w = 1/2$ (point E) is not optimal with fixed weights, it has a variance which is smaller than the optimum with estimated weights (point R).

Eqs. (4) and (5) concern special cases (normality and independence, respectively). In general, when the weights are estimated, the combined forecast will be biased, as given in [Eq. \(3\)](#), with its variance given by

$$\begin{aligned} \text{var}(y_c) = & (Ew)^2\sigma_1^2 + (1 - Ew)^2\sigma_2^2 \\ & + 2(Ew)(1 - Ew)\rho\sigma_1\sigma_2 \\ & + E[(w - Ew)(y_1 - y_2) \\ & \times ((Ew)y_1 + (1 - Ew)y_2 - \mu)] \\ & + E[(w - Ew)^2(y_1 - y_2)^2] \\ & - (\text{cov}(w, y_1 - y_2))^2. \end{aligned} \tag{6}$$

There are now additional terms over and above those in [Eqs. \(4\) and \(5\)](#), and these shift and distort the fixed-weights curve of [Fig. 1](#), as is illustrated in [Fig. 2](#). The optimal

weight is now given by w^{**} rather than w^* . Note that these conclusions would remain the same if we plotted the mean squared error rather than the variance. The three curves in Figs. 1 and 2 provide the essence of our answer to the forecast combination puzzle. The underlying formulae will be derived in m dimensions rather than two, but the story remains the same.

The simplified setup presented above assumes that the event μ is nonrandom and does not include a constant term w_0 in the combined forecast. Both assumptions can be criticized, so we address them both briefly here. If μ is random, we define the forecast errors $e_1 = y_1 - \mu$ and $e_2 = y_2 - \mu$. Including a constant term in the combined forecast gives

$$y_c = w_0 + wy_1 + (1 - w)y_2.$$

The forecast error of y_c is then

$$e_c = y_c - \mu = w_0 + we_1 + (1 - w)e_2.$$

Assume that the forecasts are unbiased, so that $Ee_1 = Ee_2 = 0$. Then, in the case of fixed weights, we find

$$Ee_c = w_0,$$

$$\text{var}(e_c) = w^2\sigma_1^2 + (1 - w)^2\sigma_2^2 + 2w(1 - w)\rho\sigma_1\sigma_2,$$

as in Eq. (2), except that σ_1^2 and σ_2^2 now denote the variances of e_1 and e_2 , and $\rho = \text{corr}(e_1, e_2)$. The mean squared error of e_c is minimized for $w_0 = 0$ and $w = w^*$, so nothing changes. Now consider the case of random weights. Then,

$$Ee_c = Ew_0 + \text{cov}(w, y_1 - y_2),$$

which vanishes for $Ew_0 = -\text{cov}(w, y_1 - y_2)$. Thus, including an intercept absorbs the bias. Regarding $\text{var}(e_c)$, this will be an expression like Eq. (6), but more complicated because the variance and covariances involving w_0 need to be included. Since the essence of the story is not affected, we shall continue to assume that the event μ is nonrandom and that the combined forecast does not include a constant term. Only in our small Monte Carlo experiment in Section 6 do we assume that μ is random. However, we shall work in m dimensions rather than in two.

The remainder of this paper is organized as follows. In Section 2, we reiterate the classical forecast combination problem in a multivariate setting, assuming that the weights are fixed. In Section 3, we analyze the properties of the forecast combination when the weights are random and the estimation is taken into account explicitly. Some special cases are considered in Section 4. Our explanation of the puzzle is summarized in Section 5. Section 6 provides a numerical illustration, and some concluding remarks are offered in Section 7.

2. Moments of the forecast combination: fixed weights

Thus motivated, let $y = (y_1, \dots, y_m)'$ be a vector of unbiased forecasts so that $Ey_j = \mu$ for all j , and let $w = (w_1, \dots, w_m)'$ be a vector of fixed (nonrandom) weights constrained by $\sum_j w_j = 1$. Assuming that y has a finite variance Σ_{yy} , we obtain the mean and variance of the forecast combination $y_c = w'y$ as

$$Ey_c = \mu, \quad \text{var}(y_c) = w'\Sigma_{yy}w. \quad (7)$$

It is easy to show that the variance is minimized (as a function of w , under the constraint $\sum_j w_j = 1$) when $w = w^*$, where

$$w^* = \frac{\Sigma_{yy}^{-1}\mathbf{1}}{\mathbf{1}'\Sigma_{yy}^{-1}\mathbf{1}} \quad (8)$$

and $\mathbf{1}$ denotes the vector of m ones. The optimal forecast is then $y_c^* = w^{*'}y$ and its variance is

$$\text{var}(y_c^*) = \frac{1}{\mathbf{1}'\Sigma_{yy}^{-1}\mathbf{1}}. \quad (9)$$

These are well-established results; see Bates and Granger (1969) for the bivariate case and Elliott (2011) for its multivariate extension.

Denote the diagonal elements of Σ_{yy} by $\sigma_1^2, \dots, \sigma_m^2$. Then, for each j ,

$$\text{var}(y_c^*) \leq \sigma_j^2. \quad (10)$$

This follows by considering the vectors $a_j = \Sigma_{yy}^{-1/2}e_j$ and $b = \Sigma_{yy}^{-1/2}\mathbf{1}$, where e_j denotes the m -dimensional vector with one in its j th position and zeros elsewhere. Then, by Cauchy–Schwarz,

$$\begin{aligned} 1 &= (e_j'\mathbf{1})^2 = (a_j'b)^2 \leq (a_j'a_j)(b'b) \\ &= (e_j'\Sigma_{yy}e_j)(\mathbf{1}'\Sigma_{yy}^{-1}\mathbf{1}) = \sigma_j^2/\text{var}(y_c^*). \end{aligned}$$

Hence, the optimally combined forecast has a smaller variance than each of the individual forecasts. Equality can occur for at most one of the individual forecasts, because Σ_{yy} is assumed to remain positive definite. Equality for the j th forecast occurs if and only if a_j and b are linearly dependent, that is, if and only if $\text{cov}(y_i, y_j) = \text{var}(y_j)$ for $i = 1, \dots, m$.

We note that we imposed the restriction that the weights add up to one, but not that each weight lies between zero and one. If all of the covariances are zero so that Σ_{yy} is diagonal, then the optimal weights are given by $(1/\sigma_j^2)/\sum_i(1/\sigma_i^2)$ ($j = 1, \dots, m$), and these clearly lie between zero and one. However, this holds only if Σ_{yy} is a diagonal matrix. Even in the case where only one covariance is not zero, say $\text{cov}(y_i, y_j) = \text{cov}(y_j, y_i) \neq 0$ for some i and j , the optimal weights w_i^* and w_j^* do not necessarily lie between zero and one; they do if and only if

$$\text{corr}(y_i, y_j) < \frac{\min(\sigma_i, \sigma_j)}{\max(\sigma_i, \sigma_j)}.$$

Apparently, the combination of a high positive correlation with a high variation in reliability forces the optimal weights outside the $(0, 1)$ interval. Of course, it is possible to choose a positive definite matrix, say V , such that the components of $V^{-1}\mathbf{1}$ are all positive, for example the diagonal matrix $V = \text{diag}(\sigma_1^2, \dots, \sigma_m^2)$. An alternative set of weights can then be defined as

$$w^\dagger = \frac{V^{-1}\mathbf{1}}{\mathbf{1}'V^{-1}\mathbf{1}}, \quad (11)$$

and these weights lie between zero and one, but, in general, they are not optimal. The forecast combination $y_c^\dagger = w^\dagger'y$ is still unbiased, but its variance is now

$$\text{var}(y_c^\dagger) = \frac{\mathbf{1}'V^{-1}\Sigma_{yy}V^{-1}\mathbf{1}}{(\mathbf{1}'V^{-1}\mathbf{1})^2}. \quad (12)$$

Letting $x = V^{-1/2}1$ and $P = V^{-1/2} \Sigma_{yy} V^{-1/2}$, we obtain

$$\frac{\text{var}(y_c^\dagger)}{\text{var}(y_c^*)} = \frac{x'Px}{x'x} \cdot \frac{x'P^{-1}x}{x'x},$$

and hence, by Kantorovich's inequality (Abadir & Magnus, 2005, Exercise 12.17),

$$1 \leq \frac{\text{var}(y_c^\dagger)}{\text{var}(y_c^*)} \leq \frac{(\lambda_1 + \lambda_m)^2}{4\lambda_1\lambda_m}, \tag{13}$$

where λ_1 and λ_m denote the largest and smallest eigenvalues of P , respectively. This provides an estimate of the possible loss of precision that is caused by choosing w^\dagger instead of w^* . In the most common case, where we choose $V = \text{diag}(\sigma_1^2, \dots, \sigma_m^2)$, we note that P is the correlation matrix associated with Σ_{yy} . Although the issue of optimal weights that are outside the (0, 1) interval is important, it is not considered further in the current paper.

When the weights are fixed, the optimal forecast combination y_c^* is an improvement over the individual forecasts, because it remains unbiased and has a smaller variance. In applications, however, the weights will typically be random, and we now turn to this more realistic case.

3. Moments of the forecast combination: random weights

As in the previous section, let $y = (y_1, \dots, y_m)'$ be a vector of unbiased forecasts with $Ey_j = \mu$, and let $w = (w_1, \dots, w_m)'$ be a vector of weights that are constrained by $\sum_j w_j = 1$, but are now random rather than fixed. Let $\Delta y_j = y_j - Ey_j$ and $\Delta y = (\Delta y_1, \dots, \Delta y_m)'$. Assuming that y and w are jointly distributed with finite fourth-order moments, and writing

$$\text{var} \begin{pmatrix} y \\ w \end{pmatrix} = \begin{pmatrix} \Sigma_{yy} & \Sigma_{yw} \\ \Sigma_{wy} & \Sigma_{ww} \end{pmatrix},$$

we have

$$y_c = w'y = \mu + w'\Delta y,$$

and hence

$$Ey_c = \mu + E(w'\Delta y) = \mu + \text{tr} \Sigma_{wy},$$

so that, in general, y_c is a biased forecast. Also,

$$\text{var}(y_c) = \text{var}(w'\Delta y),$$

$$\text{MSE}(y_c) = \text{var}(w'\Delta y) + (\text{tr} \Sigma_{wy})^2.$$

This is not yet very informative. To gain more insight, we let $\Delta w_j = w_j - Ew_j$ and $\Delta w = (\Delta w_1, \dots, \Delta w_m)'$. Then, $w = Ew + \Delta w$, and hence

$$w'\Delta y = (Ew)'\Delta y + (\Delta w)'\Delta y,$$

so that

$$\begin{aligned} \text{var}(w'\Delta y) &= (Ew)'\Sigma_{yy}(Ew) + 2(Ew)' \\ &\quad \times E[(\Delta y)(\Delta y)'(\Delta w)] + \text{var}[(\Delta w)'(\Delta y)]. \end{aligned}$$

This leads to the following proposition.

Proposition 3.1. *The mean, variance, and mean squared error of the forecast combination $y_c = w'y$ are given by*

$$Ey_c = \mu + \text{tr} \Sigma_{wy},$$

$$\text{var}(y_c) = (Ew)'\Sigma_{yy}(Ew) + 2(Ew)'d + \delta - (\text{tr} \Sigma_{wy})^2,$$

and

$$\text{MSE}(y_c) = (Ew)'\Sigma_{yy}(Ew) + 2(Ew)'d + \delta,$$

where the vector d and the scalar δ denote third- and fourth-order moments respectively, and are defined as

$$d = E[(\Delta y)(\Delta y)'(\Delta w)], \quad \delta = E[(\Delta w)'(\Delta y)]^2.$$

We note the generality of this proposition. The only two things assumed (apart from the existence of moments) are that each individual forecast is unbiased and that the weights add up to one, and it is precisely the combination of these two assumptions that leads to the simplicity of the formulas. It is *not* assumed that the weights lie between zero and one. There is no difficulty in deriving the counterpart of Proposition 3.1 for biased forecasts, but the formulae become cumbersome, and they are not needed for the story that we wish to tell.

The distribution of the weights w is given by their location (Ew) and shape (moments of Δw). We can choose the location optimally by minimizing $\text{MSE}(y_c)$ with respect to Ew under the restriction that the weights add up to one, and this leads to $Ew = w^{**}$, where

$$w^{**} = \left(\frac{1 + t' \Sigma_{yy}^{-1} d}{t' \Sigma_{yy}^{-1} 1} \right) \Sigma_{yy}^{-1} 1 - \Sigma_{yy}^{-1} d.$$

It is important to note that the 'optimal' weights w^* given in Eq. (8) are no longer optimal in the random-weights case, unless $d = 0$, which occurs, for example, when $\Sigma_{ww} = 0$ (so that $\Delta w = 0$, the fixed-weights case), or if the joint distribution is not skewed (for example, symmetric) so that third-order moments vanish. With Ew chosen optimally as w^{**} , the variance of y_c is given by

$$\begin{aligned} \text{var}(y_c) &= \frac{1 + 2t' \Sigma_{yy}^{-1} d - [(t' \Sigma_{yy}^{-1} 1)(d' \Sigma_{yy}^{-1} d) - (t' \Sigma_{yy}^{-1} d)^2]}{t' \Sigma_{yy}^{-1} 1} \\ &\quad + \delta - (\text{tr} \Sigma_{wy})^2. \end{aligned}$$

When the weights are random rather than fixed, the analysis and the conclusions are less straightforward. First, the forecast combination y_c will generally have a larger variance when the weights are random, because of the additional randomness in the weights; however, this is not always the case. Second, it is no longer true that the variance of y_c is necessarily smaller than the variance of each individual forecast, even when we choose the weights 'optimally', say $Ew = w^*$ or $Ew = w^{**}$. A discussion of some special cases will be instructive and will highlight these differences.

4. Special cases

We consider three special cases.

No skewness. If the joint distribution of (y, w) is not skewed, then the mean and variance of the forecast

combination $y_c = w'y$ are given by

$$Ey_c = \mu + \text{tr } \Sigma_{wy}$$

and

$$\text{var}(y_c) = (Ew)' \Sigma_{yy} (Ew) + \delta - (\text{tr } \Sigma_{wy})^2.$$

For example, no skewness occurs when the joint distribution is symmetric, regardless of the definition of multivariate symmetry that one employs. If the joint distribution is not skewed, then the third-order moments $d = E[(\Delta y)(\Delta y)'(\Delta w)]$ all vanish, so that $w^* = w^{**}$, and hence,

$$\text{MSE}(y_c) = (Ew)' \Sigma_{yy} (Ew) + \delta$$

contains only two terms. In this case, the combined forecast does not necessarily have a smaller variance than the individual forecasts. The first term is smaller than the individual variance σ_j^2 , see Eq. (10), but $\delta = E[(\Delta w)'(\Delta y)]^2$ is positive and, if it is large enough, then $\text{MSE}(y_c) > \sigma_j^2$.

Normality. The variance of the weights Σ_{ww} plays a key role in the variance of the combination. This is why it may be good to select an estimator with a small variation in weights, even when it is not the optimal estimator. For example, the estimator based on w^\dagger may be 'better' than the estimator based on w^* .

The effect of Σ_{ww} is brought out well in the case of joint normality. The mean and variance of the forecast combination $y_c = w'y$ are then given by

$$Ey_c = \mu + \text{tr } \Sigma_{wy}$$

and

$$\text{var}(y_c) = (Ew)' \Sigma_{yy} (Ew) + \text{tr}(\Sigma_{ww} \Sigma_{yy}) + \text{tr}(\Sigma_{wy} \Sigma_{yw}).$$

This follows from the fact that multivariate normality implies no skewness, so that $d = 0$, and also, following Anderson (1958, p. 39),

$$\begin{aligned} \delta_{ij} &\equiv E[(\Delta w_i)(\Delta y_i)(\Delta w_j)(\Delta y_j)] \\ &= \text{cov}(w_i, y_i) \text{cov}(w_j, y_j) + \text{cov}(w_i, w_j) \text{cov}(y_i, y_j) \\ &\quad + \text{cov}(w_i, y_j) \text{cov}(y_i, w_j), \end{aligned}$$

so that

$$\delta = \sum_{ij} \delta_{ij} = (\text{tr } \Sigma_{wy})^2 + \text{tr}(\Sigma_{ww} \Sigma_{yy}) + \text{tr}(\Sigma_{wy} \Sigma_{yw}).$$

The result then follows from Proposition 3.1.

Independence. One naturally expects the estimated weights w and the forecasts y to be correlated, because they are typically estimated from the same data set. In some cases, however, it may be possible to estimate the weights independently of the forecasts. When this happens, that is, when y and w are independent with finite second-order moments, then the forecast combination $y_c = w'y$ is unbiased,

$$Ey_c = \mu,$$

and its variance and mean squared error are given by

$$\text{var}(y_c) = \text{MSE}(y_c) = (Ew)' \Sigma_{yy} (Ew) + \text{tr}(\Sigma_{ww} \Sigma_{yy}).$$

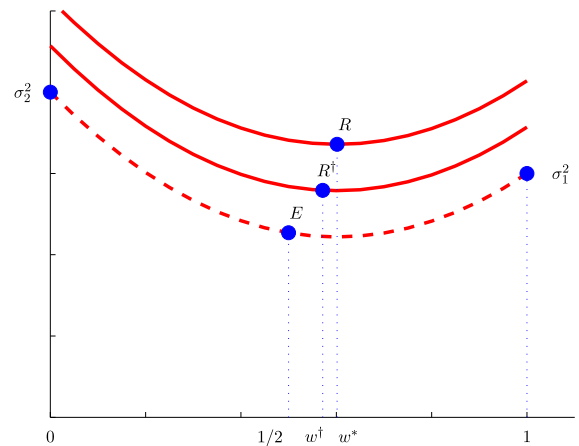


Fig. 3. Variance of forecast combination, in two dimensions: random weights under normality with and without covariances.

5. Discussion

In their 'simple explanation of the forecast puzzle', Smith and Wallis (2009) offer three main conclusions in terms of the mean squared error of the forecast (MSFE). We now analyze these conclusions in the context of the theory developed in Section 3. Their first conclusion is that

'[...] a simple average of competing forecasts is expected to be more accurate, in terms of MSFE, than a combination based on estimated weights.'

This is the situation illustrated for two dimensions in Figs. 1 and 2. The combination with equal weights is unbiased, and its variance has only one component: $t' \Sigma_{yy} t / m^2$. In many situations, this leads to a smaller mean squared error than a biased combination with additional components d and δ , as given in Proposition 3.1, for the case where the weights are estimated.

The second conclusion is that

'[...] if estimated weights are to be used, then it is better to neglect any covariances between forecast errors and base the estimates on inverse MSFEs alone, than to use the optimal formula originally given by Bates and Granger for two forecasts, or its regression generalization for many forecasts.'

Apart from the fact that including covariances may lead to negative weights, we have also seen that estimating the covariances increases the variance of the weights, as is also shown by Figures 2 and 4 of Smith and Wallis (2009). For fixed weights, the relationship between the two variances (with and without covariances) is given by Eq. (13), but the additional terms from Proposition 3.1 are likely to be larger for the optimal weights based on estimated covariances. The special cases in Section 4 emphasize this point by showing explicitly how the variance of the weights, Σ_{ww} , appears in the formulae.

Fig. 3 provides a stylized illustration in two dimensions. The figure is identical to Fig. 1, except that the middle curve has been added and the minimum point F on the lowest curve has been removed. It gives the variance of the forecast combination as a function of Ew . The bottom

curve plots the variance when the weights are nonrandom; the point E on the curve (not the minimum) gives the variance when $w = 1/2$: equal weights. The top curve plots the variance according to Proposition 3.1, and the minimum of the curve is at R , representing the point where the optimal choice for Ew is estimated to be. The middle curve represents the restricted case without covariances, where Ew is an estimate of $\sigma_z^2/(\sigma_1^2 + \sigma_2^2)$, as in Eq. (11). The minimum on the middle curve does not occur at R^\dagger ; however, R^\dagger is typically lower than R because the three variance curves move parallel to each other and fewer parameters are required to estimate the variance in the middle curve than that in the top curve.

The third conclusion of Smith and Wallis (2009) is:

‘When the number of competing forecasts is large, so that under equal weighting each has a very small weight, the simple average can gain in efficiency by trading off a small bias against a larger estimation variance. Nevertheless, in an example from Stock and Watson (2003), [...] the forecast combination puzzle rests on a gain in MSFE that has no practical significance.’

This statement is based on simulations and empirical findings, but can now be assessed in any situation by comparing the variance of the combination with equal weights, $1' \Sigma_{yy} 1/m^2$, with the variance of the combination with an estimated weight w^\dagger , given by the general formula in Proposition 3.1.

6. Numerical illustration

Our Figs. 1–3 are stylized so as to provide a simple explanation of the puzzle. In actual applications, the shift and distortion of the dashed curve in Fig. 1 will vary according to our theoretical results in Eq. (6) (for two dimensions) and Proposition 3.1 (for m dimensions). To support our theoretical results and obtain a better understanding of these shifts and distortions, we now present a simple simulation study.

We follow the experimental design of Smith and Wallis (2009, Section 3.1) closely, and in particular their case 2. We draw a sequence of $T + 1$ observations from a strictly stationary AR(2) process

$$z_t = \phi_1 z_{t-1} + \phi_2 z_{t-2} + \epsilon_t, \quad t = 1, \dots, T + 1,$$

where the $\{\epsilon_t\}$ are independent and identically distributed standard-normal variates, and ϕ_1 and ϕ_2 are given parameters that are subject to the stationarity conditions $\phi_1 + \phi_2 < 1$, $\phi_2 - \phi_1 < 1$, and $|\phi_2| < 1$. The variance of the process is given by

$$\sigma_z^2 = \text{var}(z_t) = \frac{1 - \phi_2}{(1 + \phi_2)[(1 - \phi_2)^2 - \phi_1^2]},$$

and the first two autocorrelation coefficients are

$$\rho_1 = \text{corr}(z_t, z_{t-1}) = \frac{\phi_1}{1 - \phi_2},$$

$$\rho_2 = \text{corr}(z_t, z_{t-2}) = \phi_1 \rho_1 + \phi_2.$$

Our aim is to forecast the final observation z_{T+1} . Two forecasts are available,

$$y_1 = \rho_1 z_T \quad \text{and} \quad y_2 = \rho_2 z_{T-1},$$

and we are interested in the properties of various forecast combinations $y_c = w y_1 + (1 - w) y_2$ for different values of ϕ_1 and ϕ_2 . We let $T = 30$ and use the thirty observations (z_1, \dots, z_T) to estimate the weight w .

Since the forecasted z_{T+1} is random rather than fixed, we define $e_{1t} = z_t - \rho_1 z_{t-1}$ and $e_{2t} = z_t - \rho_2 z_{t-2}$, and consider the forecast errors

$$e_1 = e_{1,T+1} = z_{T+1} - y_1, \quad e_2 = e_{2,T+1} = z_{T+1} - y_2.$$

Their variances are

$$\sigma_1^2 = \text{var}(e_1) = \sigma_z^2(1 - \rho_1^2),$$

$$\sigma_2^2 = \text{var}(e_2) = \sigma_z^2(1 - \rho_2^2),$$

and their correlation is given by $\rho = \text{cov}(e_1, e_2)/(\sigma_1 \sigma_2)$, where

$$\text{cov}(e_1, e_2) = \sigma_z^2(1 - \rho_2)(1 - \rho_1^2 + \rho_2).$$

Letting $\bar{e}_1 = (1/(T - 2)) \sum_{t=2}^{T-1} e_{1,t+1}$ and $\bar{e}_2 = (1/(T - 2)) \sum_{t=2}^{T-1} e_{2,t+1}$, we obtain unbiased estimates of the second-order moments as

$$\begin{pmatrix} \hat{\sigma}_1^2 & \hat{\rho} \hat{\sigma}_1 \hat{\sigma}_2 \\ \hat{\rho} \hat{\sigma}_1 \hat{\sigma}_2 & \hat{\sigma}_2^2 \end{pmatrix} = \frac{1}{T - 3} \times \sum_{t=2}^{T-1} \begin{pmatrix} (e_{1,t+1} - \bar{e}_1)^2 & (e_{1,t+1} - \bar{e}_1)(e_{2,t+1} - \bar{e}_2) \\ (e_{1,t+1} - \bar{e}_1)(e_{2,t+1} - \bar{e}_2) & (e_{2,t+1} - \bar{e}_2)^2 \end{pmatrix}.$$

Three weights are considered: the arithmetic mean $w = 1/2$, the optimal weight

$$w^* = \frac{\hat{\sigma}_2^2 - \hat{\rho} \hat{\sigma}_1 \hat{\sigma}_2}{\hat{\sigma}_1^2 + \hat{\sigma}_2^2 - 2 \hat{\rho} \hat{\sigma}_1 \hat{\sigma}_2}$$

according to Eq. (8), and the simplified weight (with $\rho = 0$)

$$w^\dagger = \frac{\hat{\sigma}_2^2}{\hat{\sigma}_1^2 + \hat{\sigma}_2^2}$$

according to Eq. (11). When the weight is fixed at $w = 1/2$, we do not estimate it, but simply compute its exact variance as

$$\text{var}[(e_1 + e_2)/2] = \frac{\sigma_z^2}{4} (4 - 3\rho_1^2 - 3\rho_2^2 + 2\rho_1^2 \rho_2).$$

However, when the weights are either w^* or w^\dagger , we do estimate them, because our purpose is to obtain a better understanding of the uncertainty that is caused by the estimation of weights. For the same reason, we do not estimate the parameters ϕ_1 and ϕ_2 ; these are set to their true values. Thus, any uncertainty that is shown in the simulations is caused by weight estimation.

This experiment is repeated 1,000,000 times, which suffices to control the simulation error. For given values of ϕ_1 and ϕ_2 , each run produces values of w and of the two forecast errors $e_1 = z_{T+1} - y_1$ and $e_2 = z_{T+1} - y_2$. Since the forecasted z_{T+1} is random rather than fixed, Eq. (6) needs to be written in terms of e_1 and e_2 , as was discussed at the end of Section 1. The variance of the error $e_c = w e_1 + (1 - w) e_2$

Table 1

Detailed numerical analysis of Eq. (14) for a fixed weight $w = 1/2$ and for estimated w^\dagger and w^* when $\phi_1 = \phi_2 \in [-0.9, 0.4]$. The simulation error is the difference between $\text{var}(e_c)$ and the sum of terms 1–5.

$\phi_1 = \phi_2$	$\text{var}(e_c)$	$\text{var}(e_c)$	$E(w)$	$E(e_c)$	Term 1	Term 2	Term 3	Term 4	Term 5	Simul err
	$w = 1/2$	w is estimated by w^\dagger								
-0.9	4.1413	4.1914	0.5069	0.00009	1.3525	1.2796	1.5094	0.0199	0.0065	0.0236
-0.8	2.2840	2.3012	0.5051	0.00006	0.7088	0.6802	0.8950	0.0070	0.0042	0.0060
-0.7	1.6782	1.6856	0.5034	-0.00003	0.4968	0.4836	0.6978	0.0031	0.0027	0.0016
-0.6	1.3867	1.3905	0.5019	-0.00003	0.3936	0.3876	0.6055	0.0014	0.0016	0.0007
-0.5	1.2222	1.2235	0.5011	-0.00001	0.3348	0.3319	0.5556	0.0006	0.0009	-0.0002
-0.4	1.1224	1.1222	0.5005	0.00001	0.2982	0.2970	0.5272	0.0002	0.0004	-0.0008
-0.3	1.0609	1.0594	0.5002	0.00000	0.2749	0.2745	0.5114	0.0001	0.0002	-0.0017
-0.2	1.0243	1.0278	0.5001	-0.00001	0.2605	0.2604	0.5035	0.0000	0.0000	0.0034
-0.1	1.0055	1.0058	0.5001	0.00000	0.2526	0.2525	0.5005	0.0000	0.0000	0.0003
0.1	1.0045	1.0043	0.4999	0.00000	0.2524	0.2526	0.4994	0.0000	0.0000	-0.0002
0.2	1.0156	1.0179	0.4997	-0.00001	0.2601	0.2607	0.4948	0.0000	0.0001	0.0022
0.3	1.0283	1.0278	0.4997	-0.00005	0.2744	0.2751	0.4788	0.0000	0.0006	-0.0012
0.4	1.0317	1.0335	0.4996	0.00007	0.2971	0.2981	0.4365	0.0000	0.0028	-0.0010
	$w = 1/2$	w is estimated by w^*								
-0.9	4.1413	4.2911	0.5176	0.00017	1.4100	1.2248	1.5078	0.0502	0.0444	0.0540
-0.8	2.2840	2.3681	0.5152	0.00015	0.7374	0.6527	0.8942	0.0217	0.0414	0.0207
-0.7	1.6782	1.7405	0.5106	-0.00010	0.5111	0.4697	0.6975	0.0120	0.0397	0.0103
-0.6	1.3867	1.4387	0.5040	-0.00009	0.3969	0.3844	0.6054	0.0069	0.0389	0.0061
-0.5	1.2222	1.2676	0.4963	-0.00007	0.3285	0.3382	0.5555	0.0040	0.0382	0.0031
-0.4	1.1224	1.1632	0.4853	0.00002	0.2804	0.3153	0.5268	0.0019	0.0380	0.0008
-0.3	1.0609	1.0993	0.4688	0.00002	0.2415	0.3101	0.5094	0.0011	0.0378	-0.0006
-0.2	1.0243	1.0651	0.4384	-0.00023	0.2002	0.3285	0.4958	-0.0003	0.0378	0.0029
-0.1	1.0055	1.0431	0.3571	-0.00006	0.1288	0.4175	0.4596	-0.0002	0.0374	0.0000
0.1	1.0045	1.0415	0.6676	0.00010	0.4503	0.1116	0.4433	-0.0006	0.0378	-0.0007
0.2	1.0156	1.0549	0.5845	-0.00016	0.3559	0.1798	0.4807	-0.0008	0.0380	0.0011
0.3	1.0283	1.0649	0.5546	-0.00039	0.3381	0.2180	0.4731	-0.0007	0.0384	-0.0023
0.4	1.0317	1.0674	0.5358	0.00026	0.3418	0.2565	0.4343	-0.0011	0.0381	-0.0019

of the combined forecast is

$$\begin{aligned}
 \text{var}(e_c) &= \underbrace{(Ew)^2\sigma_1^2}_{\text{term 1}} + \underbrace{(1 - Ew)^2\sigma_2^2}_{\text{term 2}} + \underbrace{2(Ew)(1 - Ew)\rho\sigma_1\sigma_2}_{\text{term 3}} \\
 &+ \underbrace{E[(w - Ew)(e_1 - e_2)](Ewe_1 + (1 - Ew)e_2)}_{\text{term 4}} \\
 &+ \underbrace{E[(w - Ew)^2(e_1 - e_2)^2]}_{\text{term 5}} - \underbrace{(\text{cov}(w, e_1 - e_2))^2}_{\text{term 6}}, \quad (14)
 \end{aligned}$$

which has six terms that can each be calculated from the simulations. We compute $\text{var}(e_c)$ and its six components for both $w = w^*$ and $w = w^\dagger$, for various values of ϕ_1 and ϕ_2 .

The results are presented in Tables 1 and 2. In Table 1, we let $\phi_1 = \phi_2$ for values ranging between -0.9 and 0.4. In Table 2, we fix $\phi_1 = 0.5$ and let ϕ_2 range between -0.9 and 0.4. The first three terms of Eq. (14) are present regardless of whether randomness of w is taken into account or not. Terms 4 and 5 account for the randomness in w that is caused by the estimation. Term 6 represents the squared bias, which is negligible in all cases, and is therefore omitted from the tables.

Our results are in general agreement with those of Smith and Wallis (2009). In Table 1, where $\phi_1 = \phi_2$, the variance of e_c is larger for the estimated weight w^\dagger than for the fixed weight $w = 1/2$, but not much. However, the variance of e_c is much larger (3%–4%) for the estimated weight w^* than for the fixed weight. In Table 2, where $\phi_1 \neq \phi_2$, the variance of e_c is generally smaller, sometimes substantially so (up to about 15%), for the estimated weights

w^\dagger and w^* than for the fixed weight. This is because, if the optimal weight deviates by much from one-half, the gain from estimating the optimal weight is larger than the loss caused by estimation error; see Elliott (2011) for a detailed discussion of this issue. When $w = w^\dagger$, terms 4 and 5 are close to zero, but when $w = w^*$, term 5 (the fourth-order moments) can be substantial.

Graphs are particularly informative, as they show the relative positions of the forecasts and the corresponding curves. We consider two special cases, one representing each table. Fig. 4(a) shows the case where $\phi_1 = \phi_2 = 0.4$, while Fig. 4(b) shows the case where $\phi_1 = 0.5$ and $\phi_2 = -0.8$. In Fig. 4(a), we have $\phi_1 = \phi_2$, and hence $\rho_1 = \rho_2$ and $\text{var}(y_1) = \text{var}(y_2)$. This implies that σ_1^2 and σ_2^2 are close, meaning that the expected values of the estimated w^\dagger and w^* are both close to 1/2. Since the estimation of w^\dagger hardly affects the properties of the forecast combination, the points E and R^\dagger (and the corresponding curves) are almost identical. On the other hand, the estimation of w^* increases the variance of the combination, so point R is much higher, and its corresponding curve is shifted up by terms 4 and 5, as reported in Table 1.

In Fig. 4(b), the original forecasts y_1 and y_2 have different variances, and hence σ_1^2 and σ_2^2 are not close. Again, the estimation of w^\dagger hardly affects the corresponding curve, but the point R^\dagger slides along the curve, yielding a smaller variance than the equal-weight combination E . The estimation of w^* distorts the corresponding curve, but the minimum point R offsets this distortion and produces a variance similar to that of point R^\dagger , as is reported in Table 2.

Table 2

Detailed numerical analysis of Eq. (14) for a fixed weight $w = 1/2$ and for estimated w^\dagger and w^* when $\phi_1 = 0.5$ and $\phi_2 \in [-0.9, 0.4]$. The simulation error is the difference between $\text{var}(e_c)$ and the sum of terms 1–5.

ϕ_2	$\text{var}(e_c)$	$\text{var}(e_c)$	$E(w)$	$E(e_c)$	Term 1	Term 2	Term 3	Term 4	Term 5	Simul err
	$w = 1/2$		w is estimated by w^\dagger							
-0.9	2.7064	2.3201	0.3235	-0.00017	0.5508	1.0599	0.7105	-0.0099	0.0169	-0.0082
-0.8	1.7724	1.6749	0.3857	0.00000	0.4132	0.6395	0.6201	-0.0026	0.0074	-0.0027
-0.7	1.4637	1.4408	0.4309	0.00005	0.3640	0.4827	0.5894	-0.0003	0.0035	0.0015
-0.6	1.3115	1.3095	0.4656	-0.00001	0.3387	0.3972	0.5705	0.0004	0.0017	0.0011
-0.5	1.2222	1.2243	0.4917	0.00001	0.3224	0.3444	0.5554	0.0005	0.0009	0.0007
-0.4	1.1645	1.1650	0.5117	-0.00002	0.3117	0.3094	0.5422	0.0004	0.0005	0.0007
-0.3	1.1251	1.1199	0.5264	0.00002	0.3045	0.2860	0.5302	0.0004	0.0003	-0.0014
-0.2	1.0972	1.0903	0.5366	0.00001	0.2999	0.2707	0.5189	0.0003	0.0002	0.0003
-0.1	1.0771	1.0676	0.5428	0.00000	0.2976	0.2618	0.5077	0.0003	0.0003	-0.0001
0.1	1.0516	1.0407	0.5444	0.00003	0.2994	0.2599	0.4821	0.0004	0.0006	-0.0018
0.2	1.0430	1.0378	0.5398	-0.00002	0.3035	0.2670	0.4645	0.0005	0.0011	0.0011
0.3	1.0345	1.0294	0.5311	0.00000	0.3099	0.2803	0.4394	0.0004	0.0022	-0.0030
0.4	1.0228	1.0276	0.5177	-0.00011	0.3190	0.3019	0.4003	0.0004	0.0045	0.0014
	$w = 1/2$		w is estimated by w^*							
-0.9	2.7064	2.2844	0.1868	-0.00025	0.1837	1.5314	0.4932	0.0110	0.0535	0.0114
-0.8	1.7724	1.6724	0.2252	-0.00007	0.1408	1.0173	0.4566	0.0055	0.0467	0.0055
-0.7	1.4637	1.4636	0.2733	0.00008	0.1465	0.7869	0.4774	0.0038	0.0434	0.0057
-0.6	1.3115	1.3494	0.3421	-0.00004	0.1829	0.6019	0.5161	0.0032	0.0415	0.0039
-0.5	1.2222	1.2692	0.4380	0.00002	0.2558	0.4211	0.5470	0.0033	0.0386	0.0035
-0.4	1.1645	1.2016	0.5706	-0.00016	0.3877	0.2392	0.5317	0.0039	0.0350	0.0040
-0.3	1.1251	1.1373	0.7314	0.00016	0.5879	0.0919	0.4178	0.0052	0.0311	0.0036
-0.2	1.0972	1.0841	0.8794	0.00014	0.8055	0.0183	0.2213	0.0047	0.0298	0.0046
-0.1	1.0771	1.0488	0.9520	0.00002	0.9154	0.0029	0.0935	0.0024	0.0321	0.0024
0.1	1.0516	1.0405	0.8437	0.00028	0.7191	0.0306	0.2563	-0.0010	0.0386	-0.0028
0.2	1.0430	1.0527	0.7420	-0.00012	0.5735	0.0839	0.3580	-0.0013	0.0390	-0.0005
0.3	1.0345	1.0542	0.6500	-0.00002	0.4642	0.1562	0.4015	-0.0014	0.0383	-0.0046
0.4	1.0228	1.0570	0.5764	-0.00032	0.3956	0.2328	0.3914	-0.0009	0.0379	-0.0002

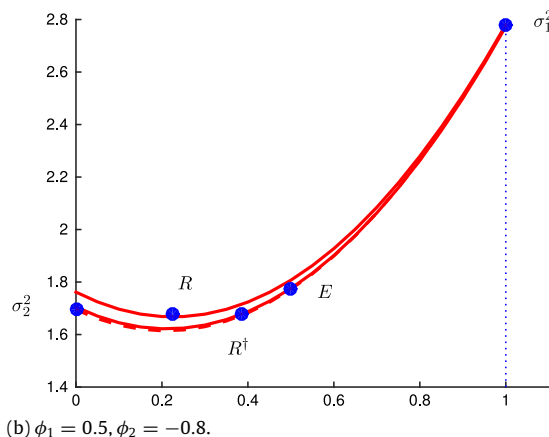
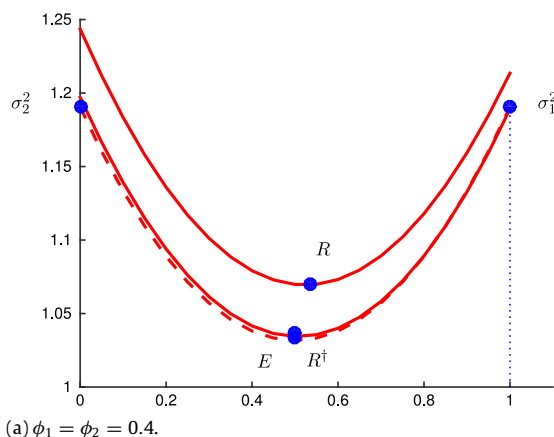


Fig. 4. Relative positions of the forecasts and the corresponding curves. The point E represents the combination with equal weights, point R^\dagger the combination with an estimated w^\dagger , and point R the combination with an estimated w^* . The original forecasts are labeled with their variances σ_1^2 and σ_2^2 .

7. Concluding remarks

In analyzing the properties of a combined forecast, we have followed an integrated approach where the estimation of the weight is accounted for explicitly from the start. Weight estimation always affects the variance of the combination. This effect may be small in some situations, but the influence is substantial in the case where the optimal weight is estimated. This is our explanation of the forecast combination puzzle.

In this paper, we have concentrated on the bias, variance, and mean squared error of the combined forecast.

These are the moments that scientists are typically interested in. Other (functions of) moments could also be analyzed similarly, such as the absolute percentage error, mean absolute deviation, or directional accuracy.

Acknowledgments

The authors are grateful to the editor and two reviewers for their constructive comments, and to Rob J. Hyndman for stimulating conversations. Earlier versions of this paper were presented at the Econometric Theory Research

Seminar and the SERG meeting in Sydney (November 2013); and at Curtin University in Perth, the University of Tasmania, AMES2014 in Taipei, the Chinese Academy of Sciences in Beijing, ISF2014 in Rotterdam, and ESEM2014 in Toulouse, all in 2014. We thank the participants for their positive comments. Claeskens acknowledges support from the Fund for Scientific Research Flanders, KU Leuven Grant GOA/12/14, and the IAP Research Network P7/06 of the Belgian Science Policy.

References

- Abadir, K. M., & Magnus, J. R. (2005). *Matrix algebra*. New York: Cambridge University Press.
- Anderson, T. W. (1958). *An introduction to multivariate statistical analysis*. New York: John Wiley & Sons.
- Bates, J. M., & Granger, C. W. J. (1969). The combination of forecasts. *Operational Research Quarterly*, 20, 451–468.
- Elliott, G. (2011). Averaging and the optimal combination of forecasts. UCSD working paper, <http://econweb.ucsd.edu/~grelliott/AveragingOptimal.pdf>.
- Graefe, A., Armstrong, J. S., Jones, R. J., Jr., & Cuzán, A. G. (2014). Combining forecasts: An application to elections. *International Journal of Forecasting*, 30, 43–54.
- Hansen, B. E. (2008). Least-squares forecast averaging. *Journal of Econometrics*, 146, 342–350.
- Hsiao, C., & Wan, S. K. (2014). Is there an optimal forecast combination? *Journal of Econometrics*, 178, 294–309.
- Liang, H., Zou, G., Wan, A. T. K., & Zhang, X. (2011). Optimal weight choice for frequentist model average estimators. *Journal of the American Statistical Association*, 106, 1053–1066.
- Magnus, J. R., & De Luca, G. (2016). Weighted-average least squares (WALS): A survey. *Journal of Economic Surveys*, 30, 117–148.
- Smith, J., & Wallis, K. F. (2009). A simple explanation of the forecast combination puzzle. *Oxford Bulletin of Economics and Statistics*, 71, 331–355.
- Stock, J. H., & Watson, M. W. (2003). How did leading indicator forecasts perform during the 2001 recession? *Federal Reserve Bank of Richmond Economic Quarterly*, 89, 71–90.

Gerda Claeskens is full professor at the research center ORSTAT (Operations Research and Business Statistics) and the Leuven Statistics Research Center of the KU Leuven, Belgium. She is the author of numerous papers, book chapters and the book “Model Selection and Model Averaging”.

Jan R. Magnus is extraordinary professor of Econometrics at the Vrije Universiteit Amsterdam. Magnus is (co)author of eight books and over one hundred scientific papers. His current interests are mainly model averaging, sensitivity analysis, catastrophe, and sport statistics.

Andrey L. Vasnev graduated in Applied Mathematics from Moscow State University in 1998. In 2001 he completed his Master's degree in Economics in the New Economic School, Moscow. In 2006 he received Ph.D. degree in Economics from the Department of Econometrics and Operations Research at Tilburg University. He worked as a credit risk analyst in ABN AMRO bank before joining the University of Sydney in 2008.

Wendun Wang received Ph.D. degree from the University of Tilburg in 2013. He is currently an assistant professor at Econometric Institute, Erasmus University Rotterdam (EUR).