

### **Forecast Combinations**

Joint Vienna Institute/ IMF ICD Macro-econometric Forecasting and Analysis JV16.12, L09, Vienna, Austria, May 24, 2016 Presenter Mikhail Pranovich

This training material is the property of the International Monetary Fund (IMF) and is intended for use in IMF's Institute for Capacity development (ICD) courses. Any reuse requires the permission of ICD.

### Lecture Objectives

- Introduce the idea and rationale for forecast averaging
- Identify forecast averaging implementation issues
- Become familiar with a number of forecast averaging schemes

### Introduction

- Usually, multiple forecasts are available to decision makers
- Differences in forecasts reflect:
  - differences in subjective priors
  - differences in modeling approaches
  - differences in private information
- It is hard to indentify the true DGP
  - should we use a single forecast or an "average" of forecasts?

### Introduction

- Disadvantages of using a single forecasting model:
  - may contain misspecifications of an unknown form
    - e.g., some variables are missing
  - one statistical model is unlikely to dominate all its rivals at all points of the forecast horizon
- Combining separate forecasts offers :
  - a simple way of building a complex, more flexible forecasting model to explain the data
  - some insurance against "breaks" or other non-stationarities that may occur in the future

### Outline of the lecture

- 1. What is a combination of forecasts?
- 2. The theoretical problem and implementation issues
- 3. Methods to assign weights
- 4. Improving the estimates of the theoretical model performance
- 5. Conclusion Key Takeaways

# Part I. What is a combination of forecasts?

- General framework and notation
- The forecast combination problem
- Issues and clarifications

### **General framework**

- Today (at time T) we want to forecast the value of  $y_{T+h}$  (at T+h)
- We have *M* different forecasts:
  - model-based (econometric model, or DSGE), or judgmental (consensus forecasts)
  - the model(s) or judgment(s) are our own or of others
  - some models or information sets might be unknown: only the end product – forecasts – are available
- How to combine *M* forecasts into one forecast?
- Is there any advantage in combining vs. selecting the "best" among the *M* forecasts?

### Notation

- $y_T$  is the value of Y at time t (today is T)
- *h* is the forecasting horizon
- $\hat{y}_{T+h,m}$  is an unbiased (point) forecast of  $y_{T+h}$  at time T
- *m*= 1,...,*M* the indices of the available forecasts/models
- $e_{T+h,m} = y_{T+h} \hat{y}_{T+h,m}$  is the forecast error of model *m*
- $\sigma_{T+h,m} = Var(e_{T+h,m})$  is the forecast error variance
- $\sigma_{T+h,i,j} = Cov(e_{T+h,i}, e_{T+h,j})$  covariance of forecast errors
- $w_{T,h} = (w_{T,h,1}, \dots, w_{T,h,M})'$  is a vector of weights
- $L(e_{t+h})$  is the loss from making a forecast error
- $E\{L(e_{t+h})\}$  is the risk associated with a forecast

### Interpretation of loss function L(e)

Squared error loss (mean squared forecasting error: MSFE)

$$L(e_{T+h}) = (e_{T+h})^2$$

- equal loss from over/under prediction
- · loss increases quadratically with the error size
- Absolute error loss (mean absolute forecasting error: MAFE)

$$L(e_{T+h}) = \mid e_{T+h} \mid$$

- equal loss from over/under prediction
- proportional to the error size
- Linex loss (γ>0 controls the aversion against positive errors, γ<0 controls the aversion against negative errors)</li>

$$L(e_{T+h}) = \exp(\gamma e_{T+h}) - \gamma e_{T+h} + 1$$

### The forecast combination problem

A combined forecast is a weighted average of *M* forecasts:

$$\hat{y}_{T+h}^{c} = \sum_{m=1}^{M} w_{T,h,m} \hat{y}_{T+h,m}$$

The forecast combination problem can be formally stated as:

**Problem 1:** Choose weights  $w_{T,h,i}$  to minimize the loss

 $\sum w_{T,h,m} = 1$ 

See Appendix 1 for generalization

function  $E[(e_{\tau+k}^c)^{2}]$  subject to

• <u>Note</u>: Here we assume MSFE-loss, but it could be any other

### Clarification: combining forecasting errors

Notice that since 
$$\sum_{i=1}^{M} w_{T,h,i}$$
 =then  
 $e_{T+h}^{c} = y_{T+h} - \hat{y}_{T+h}^{c} = y_{T+h} \sum_{m=1}^{M} w_{T,h,m} - \sum_{m=1}^{M} w_{T,h,m} \hat{y}_{T+h,m} =$   
 $= \sum_{m=1}^{M} w_{T,h,m} (y_{T+h} - \hat{y}_{T+h,m}) =$   
 $= \sum_{m=1}^{M} w_{T,h,m} e_{T+h,m}$ 

. .

 Hence, if weights sum to one, then the expected loss from the combined forecast error is

$$E\left[L\left(e_{T+h}^{c}\right)\right] = E\left[L\left(\sum_{i=1}^{M} w_{T,h,i}e_{T+h,i}\right)\right]$$

### Summary: what is the problem all about? (II)

**Problem 1:** Choose weights  $w_{T,h,i}$  to minimize the loss

We want to find optimal weights (the theoretical solution to Problem 1)

function  $E[(e_{T+h}^c)^{M}]$   $\sum_{k=1}^{M} w_{T,h,i} = 1$ 

- How can we estimate <u>optimal</u> weights from a sample of data?
- Are these estimates good?

# General problem of finding optimal forecast combination

Let:

- *u* an (*M* x 1) vector of 1's,
- and  $\Sigma$  the (*M* x *M*) covariance matrix of the forecast errors  $\Sigma = E\{ee'\}$

#### It follows that

$$\sum_{m=1}^{M} w_{T,h,m} = u'w, \ e_{T+h}^{c} = \sum_{m=1}^{M} w_{T,h,m} e_{T+h,m} = w'e, \ (e_{T+h}^{c})^{2} = w'ee'w$$

• For the MSFE loss, the optimal *w*'s are the solution to the problem:

$$E\{(e_{T+h}^c)^2\} = E\{w'ee'w\} = w'E\{ee'\}w = w'\Sigma w \Longrightarrow \min_{w} s.t. \ u'w = 1$$

• To find optimal weights it is therefore important to know (or have a "good" estimate) of  $\Sigma$ 

### **Issues and clarifications**

- Do weights have to sum to one?
  - If forecasts are unbiased, this guarantees unbiased combination forecast
- Is there a difference between averaging across forecasts and across forecasting models?
  - If you *know* the models and the models are *linear* in parameters, there is no difference
- Is it better to combine forecasts rather than information sets?
  - Combining information sets is theoretically better\*
  - practically difficult'/impossible: if sets are different, then the joint set may include so many variables that it will not be possible to construct a model that includes all of them
- \* Clemen (1987) shows that this depends on the extent to which information is common to forecasters

### Summary: what is the problem all about? (I)

- Observations of a variable Y
- Forecast observations of Y:
  - forecast 1
  - ...
  - forecast M
- Forecasting errors
- Question: how much weight to assign to each of forecasts, given past performance and knowing that there will be a forecasting error?



# Part II. The theoretical problem and implementation issues

- A simple example with only 2 forecasts
- The general N forecast framework
- Issue 1: do weights sum to 1?
- Issue 2: are weights constant over time?
- Issue 3: are estimates of weights good?

### Optimal weights in population (M = 2)

Assume we have 2 unbiased forecasts ( $E(e_{T+h,m}) = 0$ ) and combine:

$$\hat{y}_{T+h}^{c} = w\hat{y}_{1,T+h} + (1-w)\hat{y}_{2,T+h}$$

$$E\{(e_{T+h}^{c})^{2}\} = E\{(we_{T+h,1} + (1-w)e_{T+h,2})^{2}\} =$$

$$= w^{2}\sigma_{T+h,1}^{2} + 2w(1-w)\sigma_{T+h,1,2} + (1-w)^{2}\sigma_{T+h,2}^{2} \Longrightarrow \min_{w}$$

**Result 1:** The solution to **Problem 1** is

$$w = \frac{\sigma_{T+h,2}^{2} - \sigma_{T+h,1,2}}{\sigma_{T+h,1}^{2} + \sigma_{T+h,2}^{2} - 2\sigma_{T+h,1,2}} \quad \text{weight of } \hat{y}_{1,T+h}$$
$$1 - w = \frac{\sigma_{T+h,1}^{2} - \sigma_{T+h,1,2}}{\sigma_{T+h,1}^{2} + \sigma_{T+h,2}^{2} - 2\sigma_{T+h,1,2}} \quad \text{weight of } \hat{y}_{2,T+h}$$

### Interpreting the optimal weights in population

Consider the ratio of weights

$$\frac{w}{1-w} = \frac{\sigma_{T+h,2}^2 - \sigma_{T+h,1,2}}{\sigma_{T+h,1}^2 - \sigma_{T+h,1,2}}$$

- A larger weight is assigned to a more precise forecast
- If the covariance of the two forecasts increases, a greater weight goes to a more precise forecast
- The weights are the same (w = 0.5) if and only if  $\sigma_{T+h,2}^2 = \sigma_{T+h,1}^2$
- This is similar to building a minimum-variance-portfolio (finance)
   See Appendix 2: a generalization to M>2

# Result: Forecast combination reduces error variance

Compute the expected MSFE with the optimal weights:

$$E[(e_{T+h}^{c}(w^{*}))^{2}] = \frac{\sigma_{T+h,1}^{2}\sigma_{T+h,2}^{2}(1-\rho_{T+h,1,2}^{2})}{\sigma_{T+h,1}^{2}+\sigma_{T+h,2}^{2}-2\rho_{T+h,1,2}\sigma_{T+h,1}\sigma_{T+h,2}} \qquad \begin{array}{l} |\rho| \leq 1 \text{ Is the correlation} \\ \text{coefficient} \end{array}$$

• Suppose 
$$\sigma_{T+h,1}^2 < \sigma_{T+h,2}^2$$
 (forecast 1 is more precise), then:  

$$E[(e_{T+h}^c(w^*))^2] = \sigma_{T+h,1}^2 \underbrace{\sigma_{T+h,1}^2 + \sigma_{T+h,2}^2 - 2\rho_{T+h,1,2}\sigma_{T+h,1}\sigma_{T+h,2}}_{\text{(see Appendix 3)}} \leq \sigma_{T+h,1}^2$$

**Result 2:** 
$$E[(e_{T+h}^{c}(\hat{w}))^{2}] \le \min\{\sigma_{T+h,1}^{2}, \sigma_{T+h,2}^{2}\}$$

The combined forecast error variance is lower than the smallest of the forecasting error variances of any single model

### Estimating $\Sigma$

The key ingredient for finding the optimal weights is the forecast error covariance matrix, e.g. for *M*=2:

$$\Sigma(e_{T+h,1}, e_{T+h,2}) = \begin{pmatrix} \sigma_{T+h,1}^2 & \sigma_{T+h,1,2} \\ \sigma_{T+h,1,2} & \sigma_{T+h,2}^2 \end{pmatrix}$$

- In reality, we do not know the exact  $\Sigma$ :
  - we can only estimate  $\hat{\Sigma}$  (and then the weights) using past record of forecasting errors



### Issues with estimating $\Sigma$

- Is the estimate of  $\hat{\Sigma}$  based on the <u>past</u> forecasting errors "good"?
- If forecasting history is short, then  $\hat{\Sigma}$  may be biased
- $\hat{\Sigma}$  may or may not depend on *t* (e.g., a model/forecaster *m* may become better than others over time smaller  $\sigma_{T+hm}^2$ )
  - If not,  $\hat{\Sigma}$  converges to  $\Sigma$  as forecasting record lengthens
  - If it does, different issues: heteroskedasticity of any sort, serial correlation, etc.
- If such issues are there, the seemingly "optimal" forecast based on the estimated  $\hat{\Sigma}$  might become inferior to other (simpler) combination schemes...

### Optimality of equal weights

- The simplest possible averaging scheme uses equal weights
- The equal weights are also optimal weights if:
  - the variances of the forecast errors are the same
  - the pair-wise covariances of forecast errors are the same and equal to zero for M > 2
  - the loss function is symmetric, e.g. MSFE:
    - we are not concerned about the sign or the size of forecast errors

**Empirical observation:** Equal weights tend to perform better than many estimates of the optimal weights (Stock and Watson 2004, Smith and Wallis 2009)

### Part III. Methods to estimate the weights: M is small relative to T (M<<T)



- ... there is no point in combining use one of the models
- Rejection of H<sub>0</sub> implies that there is information in both forecasts that can be combined to get a better forecast

### OLS estimates of the optimal weights

- Recall the general problem of estimating  $w_m$  for *m* forecasts (slide 12)
- We can use OLS to estimate the  $w_m$ 's that minimize the MSFE (Granger and Ramanathan -1984):
  - we use history of past forecasts  $y_{t+h,m}$  m=1,...,M to estimate  $y_{t+h} \stackrel{M}{=} w_1 \hat{y}_{t+h,1} + w_2 \hat{y}_{t+h,2} + ... + w_M \hat{y}_{t+h,M} + \varepsilon_{t+h}$

or 
$$y_{t+h} = w_0 + w_1 \hat{y}_{t+h,1} + w_2 \hat{y}_{t+h,2} + \dots + w_M \hat{y}_{t+h,M} + \varepsilon_{t+h}$$
 s.t.  $\sum_{m=1}^{M} w_m = 1$ 

• including intercept  $w_0$  takes care of a bias of individual forecasts

### Reducing the dependency on sampling errors

- Assume that estimate  $\hat{\Sigma}$  is affected by a sampling error (e.g., is biased due to a short forecast record)
- It makes sense to reduce the dependence of the weights on such a (biased) estimate  $\hat{\Sigma}$
- Can achieve this by "shrinking" the optimal weights w's towards equal weights 1/M (Stock and Watson 2004)

$$w_{T,h,i}^{s} = \psi w_{T,h,i} + (1 - \psi) \frac{1}{M}$$
  
$$\psi = \max \{0, 1 - kM / (T - h - M - 1)\}$$

Notice:

- the parameter *k* determines the strength of the shrinkage
- as *T* increases relative to *M*, the estimated (e.g., OLS) weights become more important:  $w_{T,h,i}^s \rightarrow w_{T,h,i}$
- Can you explain why?

### Part IV. Methods to estimate the weights: when M is large relative to T

## Premise: problems with OLS weights

The problem with OLS weights:

- If *M* is large relative to *T*-*h* the OLS estimates loose precision and may not even be feasible (if *M* > *T*-*h*)
- Even if *M* is low relative to *T*-*h*, the OLS estimates of weights may be subject to a sampling error
  - the estimate  $\hat{\Sigma}$  may depend on the sample used
- A number of other methods can be used when *M* is large relative to *T*

### Relative performance weights

An alternative to the of OLS weights:

- ignore the covariance across forecast errors
- compute weights based on past forecast performance

• For each forecast compute  $MSFE_{T,h,m} = \frac{1}{T-h} \sum_{t=1}^{T-h} e_{t+h,m}^2$ 

**MSFE weights** (or relative performance weights)

$$w_{T,h,m} = \frac{\frac{1}{MSFE_{T,h,m}}}{\sum_{m=1}^{M} \frac{1}{MSFE_{T,h,m}}}$$

### Emphasizing recent performance

• Compute:

$$MSFE_{T,h,m} = \frac{1}{\widetilde{T}} \sum_{t=1}^{T-h} e_{t+h,m}^2 \delta(t)$$

where T is the number of periods with  $\delta(t) > 0$  and  $\delta(t)$  can be either

$$\delta(t) = \begin{cases} 1 & \text{if } t \ge T - h - v \\ 0 & \text{if } t < T - h - v \end{cases}$$

Using only a part of forecasting history for forecast evaluation

or

$$\delta(t) = \delta^{T-h-t}$$

**Discounted MSFE** 

## Such MSFE weights emphasize the recent forecasting performance

### Shrinking relative performance

Consider instead

$$w_{T,h,m} = \frac{\left(\frac{1}{MSFE_{T,h,m}}\right)^{k}}{\sum_{m=1}^{M} \left(\frac{1}{MSFE_{T,h,m}}\right)^{k}}$$

- If *k*=1 we obtain standard MSFE weights
- If k=0 we obtain equal weights 1/M

As parameter  $k \rightarrow 0$  the relative performance of a particular model becomes less important

### Performance weights with correlations

- MSFE weights ignore correlations between forecasting errors
- Ignoring it, when it is present decreases efficiency larger forecasting variance from the combined forecast

Consider instead

$$\hat{\Sigma}_{T,m,j} = \frac{1}{T-h} \sum_{t=1}^{T-h} e_{t+h,m} e_{t+h,j}$$

The relative performance weights adjusted for covariance:  $w_{T,h,m} = \frac{\sum_{j \neq m} \hat{\Sigma}_{T,m,j}^{-1}}{\sum_{j,m} \hat{\Sigma}_{T,m,j}^{-1}}$ 

• <u>Note</u>: this weighting scheme may be computationally intensive. For *M* models we need to calculate M(M+1)/2 different  $\hat{\Sigma}_{T,m,i}$ 

### Rank-based forecast combination

- Aiolfi and Timmerman (2006) allow the weights to be inversely related to the rank of the forecast
- The better is the forecast (e.g., according to MSFE) the higher is the rank  $r_m$
- After all models are ranked form best to worst, the weights are:

$$w_{T,h,m} = \frac{(\frac{1}{r_{T,h,m}})}{\sum_{m=1}^{M} (\frac{1}{r_{T,h,m}})}$$

### Trimming

- In forecast combination, it is often advantageous to discard models with the worst and best performance (i.e., trimming)
  - This is because simple averages are easily distorted by extreme forecasts/forecast errors
  - Trimming justifies the use of the median forecast
- Aiolfi and Favero (2003) recommend ranking the individual models by  $R^2$  and discarding the bottom and top 10 percent.

### Example

 Stock and Watson (2003): relative forecasting performance of various forecast combination schemes versus the AR (benchmark)

Table 3. MSFES of Combination Forecasts, Relative to Autoregression:Forecasts of 4-quarter Growth of Real GDP (h = 4)

	Canada	France	Germany	Italy	Japan	U.K.	U.S.
Forecast	82:I –	1020 m	82:I –	82:1 -	82:I –	82:1 -	82:1 -
period	98:IV	2	98:IV	98:IV	98:II	98:IV	98:IV
Univariate forecasts							
AR RMSE	0.025		0.018	0.019	0.023	0.018	0.016
random walk	0.99	22	1.05	1.31	2.97	0.96	1.04
Combination forecasts, full panel							
median	0.92		0.99	0.91	0.93	1.00	0.92
mean	0.88		1.00	0.82	0.88	0.98	0.90
trimmed mean	0.90	1	0.99	0.82	0.89	0.99	0.91
disc. mse(.9)	0.90		0.98	0.85	0.93	0.94	0.90
disc. mse(.95)	0.90	22	1.00	0.89	0.93	0.96	0.89
disc. mse(1)	0.91		1.00	0.91	0.92	0.97	0.87
recent best	0.85		1.26	0.71	0.97	0.80	1.67

## Part V. Improving the Estimates of the Theoretical Model Performance: Knowing the parameters in the model



So far we assumed that we do not know models from which forecasts originate

 Would our estimates of the weights improve if we knew something about these models

• e.g., if we knew the number of parameters?

### Hansen (2007) approach

For a process  $y_t$  there may be an infinite number of potential explanatory variables  $(x_{1t}, x_{2t}, ...)$ 

- In reality we deal with only a <u>finite</u> subset  $(x_{1t}, x_{2t}, ..., x_{Nt})$
- Consider a sequence of linear forecasting models where model *m* uses the first  $k_{m_1}$  variables  $(x_{1t}, x_{2t}, ..., x_{kt})$ :

$$y_{t+h} = \sum_{j=1}^{\kappa_m} \theta_j x_{jt} + b_{t,m} + \varepsilon_{t,m}$$

• with  $b_{t,m}$  the approximation error of model *m*:

$$b_{t,m} = \sum_{j=k_m+1} \theta_j x_{jt}$$

and the forecast given by

$$\hat{y}_{t+h} = \sum_{j=1}^{k_m} \theta_j x_{jt}$$

### Hansen (2007) approach (2)

- Let  $\hat{\varepsilon}_m$  be the vector of *T-h* (<u>in-sample</u>!) residuals of model *m*
- The {(*T-h*)x*M*} matrix collecting these residuals:

$$E = (\varepsilon_{t,1}, \dots, \varepsilon_{t,M})$$

- $K = (k_1, ..., k_M)$  is an  $M \ge 1$  vector of the number of parameters in each model
- The Mallow criterion is minimized with respect to w

$$C_{T-h}(w) = w' EE' w + 2s^2 K' w \Longrightarrow \min$$

where  $s^2$  is the largest of all models sample error variance estimator

- The Mallow criterion is an unbiased approximation of the combined forecast MSFE:
  - Minimizing  $C_{T-h}(w)$  delivers optimal weights w
- It is a quadratic optimization problem: numerical algorithms are available (e.g., in GAUSS, QPROG; in Excel, SOLVER)

### Example of Hansen's approach (M = 2)

We need to find *w* that minimizes the Mallow criterion:

$$(T-h)\left[ (w \ 1-w) \begin{pmatrix} s_1^2 & s_{12} \\ s_{12} & s_2^2 \end{pmatrix} \begin{pmatrix} w \\ 1-w \end{pmatrix} \right] + 2s_2^2 (k_1 \ k_2) \begin{pmatrix} w \\ 1-w \end{pmatrix}$$

• Minimizing gives:

$$w = \frac{s_2^2((k_2 - k_1)/(T - h)) - (s_2^2 - s_{12})}{s_1^2 + s_2^2 - 2s_{12}}$$

- The optimal weights
  - depend on the Var and Cov of residuals
  - penalize the larger model: the weight on the (first) smaller model increases with the size of the "larger" second model  $k_2 > k_1$
- See appendix 7 for further details

### Conclusions – Key Takeaways

- Combined forecasts imply diversification of risk (provided not all the models suffer from the same misspecification problem)
- Numerous schemes are available to formulate combined forecasts
- For a standard MSFE loss, the payoff from using covariances of errors to derive weights is small
- Simple combination schemes are difficult to beat

## Thank You!

### References

- Aiolfi, Capistran and Timmerman, 2010, "Forecast Combinations", in *Forecast Handbook*, Oxford, Edited by Michael Clements and David <u>Hendry.</u>
- <u>Clemen, Robert, 1985, "Combining Forecasts: A Review and</u> <u>Annotated Bibliography," *International Journal of Forecasting*, Vol. 5, <u>No. 4, pp. 559–583.</u>
  </u>
- Stock, James H., and Mark W. Watson, 2004, "Combination Forecasts of Output Growth in a Seven-Country Data Set," *Journal of Forecasting*, Vol. 23, No. 6, pp. 405–430.
- <u>Timmermann, Allan, 2006. "Forecast Combinations," Handbook of</u> <u>Economic Forecasting, Elsevier.</u>

### Appendix 1: generalization of problem 1

Let **w** be the (M x 1) vector of weights, **e** the (M x 1) vector of forecast errors, **u** an (M x 1) vector of 1s', and  $\Sigma$  the (M x M) variance covariance matrix of the errors

$$\boldsymbol{\Sigma} = E\left[\mathbf{ee'}\right]$$

It follows that

$$\sum_{i=1}^{M} w_{T,h,i} = \mathbf{u'w} \qquad e_{T+h}^{c} = \sum_{i=1}^{M} w_{T,h,i} e_{T+h,i} = \mathbf{w'e} \qquad \left(e_{T+h}^{c}\right)^{2} = \mathbf{w'ee'w}$$
$$E\left[\left(e_{T+h}^{c}\right)^{2}\right] = E\left[\mathbf{w'ee'w}\right] = \mathbf{w'}E\left[\mathbf{ee'}\right]\mathbf{w} = \mathbf{w'}$$

**Problem 1:** Choose w to minimize  $w'\Sigma w$  subject to u'w = 1.

### Appendix 2: generalization of result 1

**Result 1:** Let **u** be an (M x 1) vector of 1s' and  $\Sigma_{T,h}$  the variance-covariance matrix of the forecast errors  $e_{T,h,i}$ . The vector of optimal weights **w'** with M forecasts is

$$\mathbf{w'} = \frac{\mathbf{E'}_{T,h}}{\mathbf{E'}_{T,h}\mathbf{\bar{\mu}}^{1}}$$

For the proof and to see how this applies when M = 2 see Appendix 1

### Appendix 2: generalization of result 1

Let e be the (M x 1) vector of the forecast errors. Problem 1: choose the vector w to minimize E[w'ee'w] subject to u'w = 1.

Notice that  $E[w'ee'w] = w'E[ee']w = w'\Sigma w$ . The Lagrangean is

$$L = \hat{\mathbf{x}}[\hat{\mathbf{x}}]$$

and the FOC is

$$\Sigma \mathbf{w} - \lambda \mathbf{u} = 0 \Longrightarrow \mathbf{w}^* = \Sigma^{-1} \mathbf{u} \lambda$$

Using  $\mathbf{u'w} = 1$  one can obtain  $\lambda$ 

$$\mathbf{E}' \mathbf{w} \hat{\lambda} = \mathbf{u}' \stackrel{=}{=} \mathbf{u}' \Sigma \stackrel{=}{\to} \mathbf{u} \hat{\lambda} \qquad \begin{bmatrix} -1 \end{bmatrix}^{-1}$$
Substituting λ back one gives
$$\mathbf{E}' \stackrel{*}{=} \mathbf{u}' \Sigma \begin{bmatrix} \mathbf{u} & -1 \end{bmatrix}^{-1}$$

### Appendix 2: generalization of result 1 (M = 2)

Let  $\Sigma_{t,h}$  be the variance-covariance matrix of the forecasting errors  $(\sigma^2 - \sigma^2)$ 

$$\Sigma_{T,h} = \begin{pmatrix} \sigma_{T+h,1}^2 & \sigma_{T+h,1,2} \\ \sigma_{T+h,1,2} & \sigma_{T+h,2}^2 \end{pmatrix}$$

Consider the inverse of this matrix

$$\Sigma_{T,h}^{-1} = \frac{1}{\det |\Sigma_{T,h}|} \begin{pmatrix} \sigma_{T+h,2}^2 & -\sigma_{T+h,1,2} \\ -\sigma_{T+h,1,2} & \sigma_{T+h,1}^2 \end{pmatrix}$$

Let  $\mathbf{u'} = [1, 1]$ . The two weights  $w^*$  and  $(1 - w^*)$  can be written as

$$\begin{bmatrix} w^* & 1 - w^* \end{bmatrix} = \frac{\boldsymbol{\Sigma}_{\mathrm{T},\mathrm{h}}^{-1} \mathbf{u}}{\mathbf{E}'_{\mathrm{T},\mathrm{h}}^{-1} \mathbf{u}}$$

### Optimal weights in population (M = 2)

Assume we have 2 unbiased forecasts ( $E(e_{T+h,m}) = 0$ ) and combine:

$$\hat{y}_{T+h}^{c} = w\hat{y}_{1,T+h} + (1-w)\hat{y}_{2,T+h}$$

$$E\{(e_{T+h}^{c})^{2}\} = E\{(we_{T+h,1} + (1-w)e_{T+h,2})^{2}\} =$$

$$= w^{2}\sigma_{T+h,1}^{2} + 2w(1-w)\sigma_{T+h,1,2} + (1-w)^{2}\sigma_{T+h,2}^{2} \Longrightarrow \min_{w} w$$

**Result 1:** The solution to **Problem 1** is

$$w = \frac{\sigma_{T+h,2}^{2} - \sigma_{T+h,1,2}}{\sigma_{T+h,1}^{2} + \sigma_{T+h,2}^{2} - 2\sigma_{T+h,1,2}} \quad \text{weight of } \hat{y}_{1,T+h}$$
$$1 - w = \frac{\sigma_{T+h,1}^{2} - \sigma_{T+h,1,2}}{\sigma_{T+h,1}^{2} + \sigma_{T+h,2}^{2} - 2\sigma_{T+h,1,2}} \quad \text{weight of } \hat{y}_{2,T+h}$$

Notice that 
$$\sigma_{T+h,1}^2 + \sigma_{T+h,2}^2 - 2\sigma_{T+h,1,2} = E\left[(e_{T+h,1} - e_{T+h,2})^2\right] \ge 0$$
  
and that  $\sigma_{T+h,1}^2(1 - \rho_{T+h,12}^2) \ge 0$ 

Need to show that the following inequality holds

$$\frac{\sigma_{T+h,2}^{2}(1-\rho_{T+h,1,2}^{2})}{\sigma_{T+h,1}^{2}+\sigma_{T+h,2}^{2}-2\rho_{T+h,1,2}\sigma_{T+h,1}\sigma_{T+h,2}} \leq 1$$

Rearrange the above

$$\sigma_{T+h,2}^{2}(1-\rho_{T+h,1,2}^{2}) \leq \sigma_{T+h,1}^{2} + \sigma_{T+h,2}^{2} - 2\rho_{T+h,1,2}\sigma_{T+h,1}\sigma_{T+h,2}$$
$$0 \leq \sigma_{T+h,1}^{2} - 2\rho_{T+h,1,2}\sigma_{T+h,1}\sigma_{T+h,2} + \sigma_{T+h,2}^{2}\rho_{T+h,1,2}^{2}$$
$$0 \leq (\sigma_{T+h,1} - \sigma_{T+h,2}\rho_{T+h,1,2})^{2}$$

### Appendix 4: trading-off bias vs. variance

The MSFE loss function of a forecast has two components:

- the squared bias of the forecast
- the (ex-ante) forecast variance

$$E\{(y_{T+h} - \hat{y}_{T+h,m})^2\} = Bias_{T+h,m}^2 + \sigma_y^2 + Var(\hat{y}_{T+h,m})$$

 Combining forecasts offers a tradeoff: increased overall bias vs. lower (ex-ante) forecast variance

$$E\{(y_{T+h} - \hat{y}_{T+h,m})^2\} = E\{\left(\sum_{m=1}^M w_{T+h,m}bias_{T+h,m} + \varepsilon_y + \sum_{m=1}^M w_{T+h,m}[\hat{y}_{T+h,m} - E\{\hat{y}_{T+h,m}\}]\right)^2\} = \sum_{m=1}^M w_{T+h,m}^2bias_{T+h,m}^2 + \sigma_y^2 + \sum_{m=1}^M w_{T+h,m}^2Var\{\hat{y}_{T+h,m}\}$$

The MSFE loss function of *a forecast* has two components:

- the squared bias of the forecast
- the (ex-ante) forecast variance

$$E\{(y_{T+h} - \hat{y}_{T+h,m})^{2}\} = E\{(E(y_{T+h}) + \varepsilon_{y} - \hat{y}_{T+h,m})^{2}\} =$$
  
=  $E\{(E(y_{T+h}) + \varepsilon_{y} - E\{\hat{y}_{T+h,m}\} + E\{\hat{y}_{T+h,m}\} - \hat{y}_{T+h,m})^{2}\} =$   
=  $Bias_{m}^{2} + \sigma_{y}^{2} + Var\{\hat{y}_{T+h,m}\}$ 

Suppose that  $x_{T,h,i} = \mathbf{P}\mathbf{y}$  where **P** is an (m x T) matrix, **y** is a (T x 1) vector with all  $y_t$ , t = 1,...T. Consider:

$$\hat{\sigma}_{T,h,i}^{2} = \frac{1}{T-h} \sum_{t=1}^{T-h} (x_{t,h,i} - y_{t+h})^{2}$$

$$= \frac{1}{T-h} \sum_{t=1}^{T-h} (x_{t,h,i} - E[y_{t+h}] - \varepsilon_{y,t+h})^{2}$$

$$= \sigma_{y}^{2} + \frac{1}{T-h} \sum_{t=1}^{T-h} (x_{t,h,i} - E[y_{t+h}])^{2} - \frac{2}{T-h} \sum_{t=1}^{T-h} \varepsilon_{y,t+h} (x_{t,h,i} - E[y_{t+h}])^{2}$$

Consider:

$$\hat{\sigma}_{T,h,i}^{2} = \sigma_{y}^{2} + \frac{1}{T-h} \sum_{t=1}^{T-h} \left( x_{t,h,i} - E[y_{t+h}] \right)^{2} - \frac{2}{T-h} \sum_{t=1}^{T-h} \varepsilon_{y,t+h} \left( x_{t,h,i} - E[y_{t+h}] \right)$$

$$= \dots - \frac{2}{T-h} \varepsilon' (\mathbf{P} \mathbf{y} - \mathbf{E}[\mathbf{y}])$$

$$= \dots - \frac{2}{T-h} \varepsilon' (\mathbf{P} \mathbf{E}[\mathbf{y}] + \mathbf{P} \varepsilon - \mathbf{E}[\mathbf{y}])$$

$$= \dots - \frac{2}{T-h} \varepsilon' \mathbf{P} \varepsilon - \varepsilon' (\mathbf{I} - \mathbf{P}) \mathbf{E}[\mathbf{y}]$$

$$\boxtimes MSPE_{T} - \frac{2}{T-h} E[\varepsilon' \mathbf{P} \varepsilon]$$

### Appendix 6: Adaptive weights

Relative performance weights may be sensitive to adding new forecast errors (may vary wildly)

 We can use an adaptive scheme that updates previous weights by the most recently computed weights

• E.g., for the MSFE weights (can use other weighting too):

$$w_{T,h,m} = \frac{\overline{MSFE_{T,h,m}}}{\sum_{m=1}^{M} \frac{1}{MSFE_{T,h,m}}}$$

$$\overset{\boxtimes}{w_{T,h,m}} = \alpha \overset{\boxtimes}{w_{T-1,h,m}} + (1-\alpha)w_{T,h,m}, \quad 0 \le \alpha \le 1$$

 The update parameter α controls the degree of weights update from period *T*-1 to period *T*

### Appendix 7: Example of Hansen's approach (M = 2)

If the covariance term is zero the weight becomes

$$w = \frac{s_2^2((k_2 - k_1)/(T - h)) + s_2^2}{s_1^2 + s_2^2}$$

The Mallow criterion has a preference for smaller models, and models with smaller variance

• If  $k_2 = k_1$ , the criterion is equivalent to minimizing the variance of the combination fit