

# On the usefulness of the Diebold-Mariano test in the selection of prediction models: Some Monte Carlo evidence

Mauro Costantini<sup>1</sup> and Robert M. Kunst<sup>2</sup>

Presented at Österreichische Statistiktage,  
Graz, September 2011

---

<sup>1</sup>Brunel University London; [mauro.costantini@brunel.ac.uk](mailto:mauro.costantini@brunel.ac.uk)

<sup>2</sup>Institute for Advanced Studies, Vienna, Austria 1060, and University of Vienna;  
[kunst@ihs.ac.at](mailto:kunst@ihs.ac.at)

# Historical facts

- DIEBOLD & MARIANO revolutionized reporting forecast comparisons;
- Today most editors/referees demand DM or similar tests on top of every comparative forecast evaluation;
- Forecast procedures/models that beat the benchmark only 'insignificantly' are regarded as uninteresting: widespread preference for simple benchmark structures?

# The Diebold-Mariano test

DIEBOLD & MARIANO (1995, JBES) considered the test statistic

$$S = \frac{\frac{1}{T} \sum_{t=1}^T \{g(e_{1t}) - g(e_{2t})\}}{\sqrt{\frac{2\pi \hat{f}_d(0)}{T}}},$$

with  $g(e_{jt}), j = 1, 2$  denoting the loss from forecast error  $e_{jt}$  evolving from prediction model  $j$ .

The null hypothesis tested is  $H_0 : \mathbb{E}g(e_{1t}) = \mathbb{E}g(e_{2t})$ . Under  $H_0$ ,  $S$  is asymptotically standard normal distributed.

# Conceptual problems with the DM test

Two main conceptual problems:

- 1 The null hypothesis, first criticized by CHATFIELD. In real-world applications, two simple prediction models rarely achieve the same loss moment, if the DGP is far more complex;
- 2 An out-of-sample comparison of  $g(e_{jt})$  already is equivalent to an information-criterion (IC) evaluation. A test on top of an IC favors the simpler rival model, may correspond to a stronger complexity penalty.

# Using the DM test for model selection

The main purpose of the DM test is not model selection. However, DIEBOLD AND MARIANO conjecture that

*“The ability to formally compare predictive accuracy afforded by our tests may prove useful as a model-specification diagnostic, as well as a means to test both nested and non-nested hypotheses under nonstandard conditions.”*

Later, doubts were raised on the validity of the null distribution for nested hypotheses.

# Accuracy comparison: an information criterion

- WEI (1992) considers evaluating out-of-sample prediction expanding over (nearly) the whole sample, shows that this yields a consistent information criterion;
- INOUE AND KILIAN (2006) consider evaluating over a fixed share of the sample, which yields an efficient information criterion.

Comparative forecast evaluation is a stronger tool than the apparently 'casual manner' referred to by DM.

# Testing on top of a consistent IC: good news

## Proposition

*Suppose there exists a consistent information criterion  $\tau_1$  and an independent test-consistent significance test  $\tau_2$  at a given significance level  $\alpha_2$ . Then, the joint decision from rejecting  $H_0$  if both criteria prefer the alternative is a consistent model selection procedure.*

# Testing on top of a consistent IC: not so good news

## Proposition

*Suppose there exists a consistent information criterion  $\tau_1$  with implicit significance level  $\alpha_1(T)$  at  $T$ , and an independent test-consistent significance test  $\tau_2$  at significance level  $\alpha_2$ . Then, the joint test has critical level  $\alpha_1(T)\alpha_2$ .*

In practice, of course, the two decisions will not be independent.



# Testing on top of a consistent IC: bad news

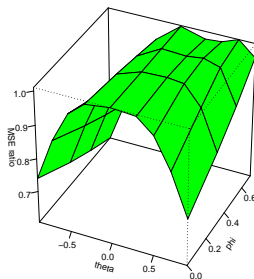
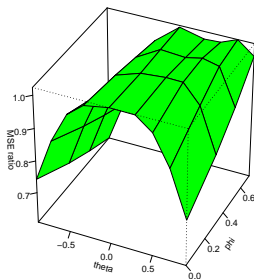
## Proposition

*Suppose there exists a consistent information criterion  $\tau$  such that between two models  $M_1$  and  $M_2$  the event  $\tau > 0$  indicates a preference for  $M_2$ , while  $\tau \leq 0$  prefers  $M_1$ . Assume the user instead bases her decision on  $\tau > \tau_0$  with  $\tau_0 > 0$ . This decision will be inconsistent in the sense that, as  $T \rightarrow \infty$ , the probability of preferring  $M_1$  although  $M_2$  is true, will not converge to 0.*

# Experiment I: the concept

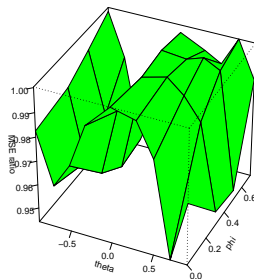
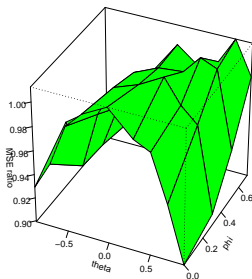
- Data are generated from ARMA(1,1) models  
 $X_t = \phi X_{t-1} + \varepsilon_t - \theta \varepsilon_{t-1}$ , with  $\phi \in \{0, 0.3, 0.5, 0.7\}$  and  $\theta \in \{-0.9, -0.7, -0.5, -0.3, 0, 0.3, 0.5, 0.7, 0.9\}$ . 1000 replications.
- An AR(1) and an ARMA(1,1) model are fitted to the data, out-of-sample predictions are based on each of the two.
- The winner over the training sample (later 50% of the data) is evaluated and predicts the last observation.
- This winner prediction is compared to the forecast based on: ARMA(1,1) if 'significantly' (5% ) better than the AR(1) 'benchmark', AR(1) otherwise.

# ARMA versus AR forecast over the training samples



MSE ratio ARMA forecast divided by AR forecast.  $N = 100$  and  $N = 200$ .

# Results from experiment I: graphical summary



MSE ratio AR or ARMA model selected by training sample divided by selected model after DM testing.  $N = 100$  (left) and  $N = 200$  (right).

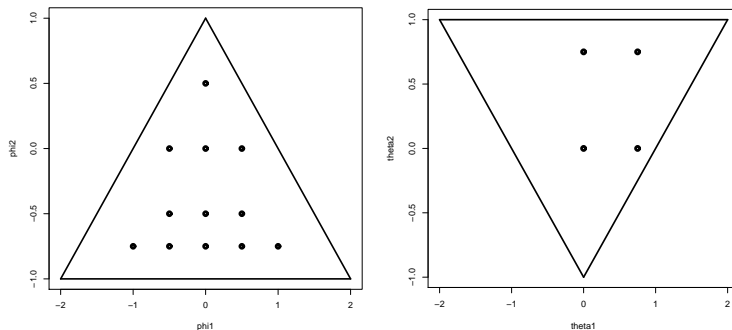
# Summary of results from experiment I

- For most designs, forecasting performance deteriorates if the DM test is applied on top of the training-sample evaluation;
- Even for cases where the AR model is correctly specified, such as  $\theta = 0$  and canceling roots, there are no benefits from using the DM test in selecting the prediction model;
- The deleterious influence of the DM test shrinks as the sample size increases.

# Experiment II: the concept

- Data are generated from ARMA(2,2) models.
- An AR(2) and an ARMA(1,1) model are fitted to the data, out-of-sample predictions are based on each of the two.
- The winner over the training sample (later 50% of the data) is evaluated and predicts the last observation.
- This winner prediction is compared to the forecast based on: ARMA(1,1) if 'significantly' better than the AR(2) 'benchmark', the AR(2) otherwise.

# ARMA(2,2) parameters



Parameter values for the autoregressive part of the generated ARMA models within the triangular region of stable AR models and values for the MA part within the invertibility region for MA(2) models.

## A non-nested design

Results of the simulation for  $N = 100$ 

$\phi_1$	$\phi_2$	$\theta_1$	$\theta_2$	MSE(ARMA)	MSE(AR)	MSE(w)	MSE(DM)
0	0.5	0	0				
-0.5	0	0	0			0.995	0.995
0	0	0	0				
0.5	0	0	0	1.047	1.047		
-0.5	-0.5	0	0				
0	-0.5	0	0				
0.5	-0.5	0	0				
-1	-0.75	0	0				
-0.5	-0.75	0	0				
0	-0.75	0	0				
0.5	-0.75	0	0				
1	-0.75	0	0				
0	0.5	0	0.75			1.318	1.318
-0.5	0	0	0.75				
0	0	0	0.75			1.166	1.166
0.5	0	0	0.75				
-0.5	-0.5	0	0.75				
0	-0.5	0	0.75				
0.5	-0.5	0	0.75				
-1	-0.75	0	0.75			1.366	1.366
-0.5	-0.75	0	0.75				
0	-0.75	0	0.75				
0.5	-0.75	0	0.75				
1	-0.75	0	0.75				
0	0.5	0.75	0				
-0.5	0	0.75	0				
0	0	0.75	0				
0.5	0	0.75	0				
-0.5	-0.5	0.75	0				
0	-0.5	0.75	0				
0.5	-0.5	0.75	0				
-1	-0.75	0.75	0				
-0.5	-0.75	0.75	0				
0	-0.75	0.75	0				
0.5	-0.75	0.75	0				
1	-0.75	0.75	0				
0	0.5	0.75	0.75				
-0.5	0	0.75	0.75				
0	0	0.75	0.75				
0.5	0	0.75	0.75				
-0.5	-0.5	0.75	0.75				
0	-0.5	0.75	0.75				
0.5	-0.5	0.75	0.75				
-1	-0.75	0.75	0.75				
-0.5	-0.75	0.75	0.75				
0	-0.75	0.75	0.75				
0.5	-0.75	0.75	0.75				
1	-0.75	0.75	0.75				



## A non-nested design

Results of the simulation for  $N = 200$ 

$\phi_1$	$\phi_2$	$\theta_1$	$\theta_2$	MSE(ARMA)	MSE(AR)	MSE(w)	MSE(DM)
0	0.5	0	0	████	████		████
-0.5	0	0	0	████	████	0.986	0.986
0	0	0	0	████	████		████
0.5	0	0	0	1.024	1.024	████	████
-0.5	-0.5	0	0	████	████	0.995	0.995
0	-0.5	0	0	████	████	0.995	0.995
0.5	-0.5	0	0	████	████	████	████
-1	-0.75	0	0	████	████	████	████
-0.5	-0.75	0	0	████	████	0.993	0.993
0	-0.75	0	0	████	████	0.994	0.994
0.5	-0.75	0	0	████	████	0.993	0.993
1	-0.75	0	0	████	████	████	████
0	0.5	0	0.75	████	████	1.344	1.344
-0.5	0	0	0.75	████	████	████	████
0	0	0	0.75	████	████	████	████
0.5	0	0	0.75	████	████	████	████
-0.5	-0.5	0	0.75	████	████	████	████
0	-0.5	0	0.75	████	████	████	████
0.5	-0.5	0	0.75	████	████	████	████
-1	-0.75	0	0.75	████	████	1.408	1.408
-0.5	-0.75	0	0.75	████	████	████	████
0	-0.75	0	0.75	████	████	████	████
0.5	-0.75	0	0.75	████	████	████	████
1	-0.75	0	0.75	████	████	████	████
0	0.5	0.75	0	1.027	1.027	0.987	0.987
-0.5	0	0.75	0	████	████	████	████
0	0	0.75	0	████	████	████	████
0.5	0	0.75	0	████	████	████	████
-0.5	-0.5	0.75	0	████	████	████	████
0	-0.5	0.75	0	████	████	████	████
0.5	-0.5	0.75	0	████	████	████	████
-1	-0.75	0.75	0	████	████	1.309	1.309
-0.5	-0.75	0.75	0	████	████	1.362	1.362
0	-0.75	0.75	0	████	████	████	████
0.5	-0.75	0.75	0	████	████	████	████
1	-0.75	0.75	0	████	████	████	████
0	0.5	0.75	0.75	████	████	████	████
-0.5	0	0.75	0.75	████	████	████	████
0	0	0.75	0.75	████	████	████	████
0.5	0	0.75	0.75	████	████	████	████
-0.5	-0.5	0.75	0.75	████	████	████	████
0	-0.5	0.75	0.75	████	████	████	████
0.5	-0.5	0.75	0.75	████	████	████	████
-1	-0.75	0.75	0.75	████	████	████	████
-0.5	-0.75	0.75	0.75	████	████	████	████
0	-0.75	0.75	0.75	████	████	████	████
0.5	-0.75	0.75	0.75	████	████	████	████
1	-0.75	0.75	0.75	████	████	1.695	1.695

# Summary of results from experiment II

- For most designs, forecasting performance improves if the DM test is applied on top of the training-sample evaluation;
- The benefits from using the DM test decrease as the sample size increases;
- Both models have identical parameter dimension. ARMA(1,1) could be chosen as the benchmark null. Then, the DM test incurs a deterioration of performance;
- The AR(2) model forecasts better due to (a) better approximation to the DGP and (b) numerically better estimation. It makes sense to view AR(2) as the benchmark.

# Experiment III: the concept

- Data are generated from a SETAR model.
- $AR(p)$  and  $ARMA(q, q)$  models are fitted to the data, with  $p$  and  $q$  determined by AIC. Out-of-sample predictions are based on each of the two.
- The winner over the training sample (25% and 50% of the data) is evaluated and predicts the last observation.
- This winner prediction is compared to the forecast based on:  $ARMA(q, q)$  if 'significantly' better than the  $AR(p)$  'benchmark', the  $AR(p)$  otherwise.

# The SETAR model used as the DGP

A SETAR model has been suggested by TIAO AND TSAY (1994) for the growth rate of U.S. GNP:

$$y_t = \begin{cases} -0.015 - 1.076y_{t-1} + \varepsilon_{1,t}, & y_{t-1} \leq y_{t-2} \leq 0, \\ -0.006 + 0.630y_{t-1} - 0.756y_{t-2} + \varepsilon_{2,t}, & y_{t-1} > y_{t-2}, y_{t-2} \leq 0, \\ 0.006 + 0.438y_{t-1} + \varepsilon_{3,t}, & y_{t-1} \leq y_{t-2}, y_{t-2} > 0, \\ 0.004 + 0.443y_{t-1} + \varepsilon_{4,t}, & y_{t-1} > y_{t-2} > 0. \end{cases}$$

Standard deviations of errors are  $\sigma_1 = 0.0062$ ,  $\sigma_2 = 0.0132$ ,  $\sigma_3 = 0.0094$ , and  $\sigma_4 = 0.0082$ .

# Results of the simulations for experiment III

	MSE $\times 10^{-4}$		frequency $\succ$	
	$N = 100$	$N = 200$	$N = 100$	$N = 200$
AR	1.115	1.037	0.518	0.479
ARMA	1.133	1.044	0.482	0.521
50% training				
lower MSE	1.113	1.041	0.523	0.518
DM-based	1.132	1.038	0.122	0.106
25% training				
lower MSE	1.106	1.042	0.544	0.544
DM-based	1.114	1.035	0.127	0.137

Note: 'frequency  $\succ$ ' gives the empirical frequency of the model yielding the better prediction for the observation at  $t = N$ .

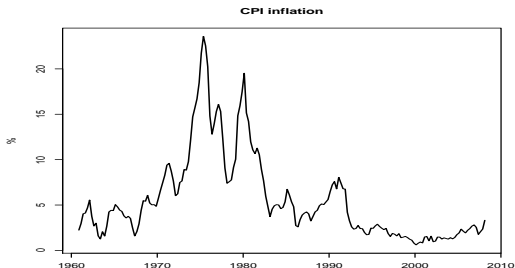
# Summary of results from experiment III

- Performance across replications is quite heterogeneous due to the highly nonlinear DGP;
- It appears that DM testing benefits the MSE ranking that may not be a good criterion here;
- By contrast, pure training-sample evaluation appears to be preferable with regard to the probability of achieving the better prediction.

# Experiment IV: the concept

- Data are generated as a component from a trivariate VAR model. The VAR is tuned to U.K. macroeconomic data;
- $AR(p)$  and  $ARMA(q, q)$  models are fitted to the data, with  $p$  and  $q$  determined by AIC. Out-of-sample predictions are based on each of the two;
- The winner over the training sample (50% of the data) is evaluated and predicts the last observation;
- This winner prediction is compared to the forecast based on:  $ARMA(q, q)$  if 'significantly' better than the  $AR(p)$  'benchmark', the  $AR(p)$  otherwise.

# Experiment IV: some details



A VAR(2) is fitted to a system that comprises U.K. GDP growth, an interest rate, and CPI inflation, and the fitted VAR is generated with Gaussian errors. CPI inflation is forecasted. Its implied generating model is ARMA(2,2), thus the generating model is contained in the prediction toolbox.



# Results of the simulations for experiment IV

	MSE		frequency $\succ$	
	$N = 100$	$N = 200$	$N = 100$	$N = 200$
AR	0.191	NYA	0.47	NYA
ARMA	0.180	NYA	0.53	NYA
50% training				
lower MSE	0.179	NYA	0.28	NYA
DM-based	0.194	NYA	0.25	NYA

Note: 'frequency  $\succ$ ' gives the empirical frequency of the model yielding the better prediction for the observation at  $t = N$ .

# Summary of results from experiment IV

- Forecasting performance deteriorates if the DM test is applied on top of the training-sample evaluation;
- The pure training-sample selection shows good performance;
- Usage of the DM test implies failure to reject the AR benchmark in 3/4 of the cases for  $N = 100$ .

# General summary

- There are no systematic benefits from 'double testing';
- Double testing using the DM test may be beneficial if the benchmark has better prediction properties but this is trivial;
- Double testing may give undue support to a simple benchmark model and lead to ignoring the benefits from using more sophistication;
- Extensions to larger forecasting horizons will be studied.

# References

- CHATFIELD, C. (2001) *Time-Series Forecasting*, Chapman & Hall.
- DIEBOLD, F.X., AND R.S. MARIANO (1995) 'Comparing Predictive Accuracy,' *Journal of Business and Economic Statistics* **13**, 253–263.
- ING, C.K. (2007) 'Accumulated prediction errors, information criteria and optimal forecasting for autoregressive time series,' *Annals of Statistics* **35**, 1238–1277.
- INOUE, A., and L. KILIAN (2006) 'On the selection of forecasting models,' *Journal of Econometrics* **130**, 273–306.
- MCQUARRIE, A.D.R., and C.-L. Tsai (1998) *Regression and Time Series Model Selection*, World Scientific.
- WEI, C.Z. (1992) 'On predictive least squares principles,' *Annals of Statistics* **20**, 1–42.

# Thank you for your attention