Beyond summary performance metrics for forecast selection and combination

Ivan Svetunkov & Nikos Kourentzes

Lancaster University, UK

18th July 2025

Marketing Analytics and Forecasting



Lancaster University Management School

Ivan Svetunkov

Lancaster University, UK

-

イロン 不同 とくほう イロン



◆□▶ ◆□▶ ◆注▶ ◆注▶ 注 のへで

How many of you have used combinations for forecasting?

How do you combine forecasts?

There's lots of papers on different combination methods:

- Means, trimmed and winsorised (Jose and Winkler, 2008);
- Median (Agnew, 1985; Stock and Watson, 2004);
- Weighted combinations (Elliott, 2011; Elliott and Timmermann, 2016; Kolassa, 2011);
- Claeskens et al. (2016) explains why simple combinations are more robust than the weighted ones.

イロト 不得 トイヨト イヨト 二日

Introduction

How many models do you have in your pool? Do you need all of them?

Geweke and Amisano (2011) shows that using a subset of models helps.

In your pool of models:

- some will perform consistently poorly, and need to be removed;
- some will do consistently well. Keep them!

But how can we decide?

Kourentzes et al. (2018) develop a heuristic to reduce the pool of models.

Ivan Svetunkov

3

Are you sure that your combination is adequate?

Will it change tomorrow?

Yes! Everything is random! You cannot be sure of anything...

Yet all the combination approaches rely on some summary statistics from a sample.

Vehtari et al. (2017) use the leave-one-out CV errors to get the standard errors for forecasts.

This way we can compare model performance...

Point likelihood



◆□▶ ◆□▶ ◆三▶ ◆三▶ 三三 のへで

Point likelihood

Take a look at AIC:

$$AIC_m = 2k_m - 2\ell_m,\tag{1}$$

 ℓ_m is the log-likelihood value, k_m is the number of estimated parameters of a model m in the pool.

This relies on the summary statistics, log-likelihood ℓ_m

It shows how likely it is to have such a model given the data.

This is a sum of point log-likelihoods:

$$\ell_m = \sum_{t=1}^T \ell_{m,t} \tag{2}$$

Point likelihood

ETS(M,A,M) applied to AirPassengers data 09 500 M. M. M. 400 300 200 100 1950 1952 1954 1956 1958 1960 Point likelihood values Υ Ϋ́ 4 4 φ φ φ ę 5 5 φ ø 1950 1952 1954 1956 1958 1960

Ivan Svetunkov

Lancaster University, UK

3. 3

・ロッ ・同 ・ ・ ヨッ ・

So, we can use point log-likelihoods to compare models.

The only issue is the bias in the likelihood estimate (Akaike, 1974).

We propose a point AIC:

$$\mathsf{pAIC}_{m,t} = 2k_m - 2T\ell_{m,t}.\tag{3}$$

◆□▶ ◆□▶ ◆臣▶ ◆臣▶ 三臣 - のへで

 ${\boldsymbol{T}}$ is needed to bring the point likelihood value to the overall level.

Other criteria can be modified like that as well.

Ivan Svetunkov

Pooling with point IC

How can we use that?

- 1. Apply models to the data;
- 2. Extract pAIC values for each of them;
- 3. Use Nemenyi test (Demšar, 2006) to group models;
- 4. Combine their forecasts.

Ivan Svetunkov

イロン 不同 とくほう イロン

Pooling with point IC

ETS on air passengers data!



Ivan Svetunkov

Lancaster University, Uk

Pooling with point IC

You can do that with CV errors as well (e.g. rolling origin).

Nemenyi can be substituted by other tests.

Ivan Svetunkov

Lancaster University, UK

3

・ロン ・四 と ・ ヨ と ・ ヨ と …

Real data experiment

0 *	SmoothTest - RStudio					~	0 X
File Edit Code View Plots Session Build Debug Profile Tools Help							
🔍 - 👒 🍘 - 🕞 🚱 🏯 🌛 Gots Herketten 🗌 🖾 - Addims -						() Smoo	dnTest +
0 2024-07-16-ABIMA-tests.R × 0 2024-07-30-ABIMA-tests.R × 0 2024-07-30-smooth-4-1-0.R × 0 202	5-07-26-058-pAIC.R* ×	Environment History Con-	ections Tutor	el .			-0
(a) a B Source on Save () / - (-+ Run 🏞 🖓 🖏 📑 Source + 🕷	😅 🔒 🖙 Import Dataset + 🕴	9 2.32 GIB + 🥑			💠 Grid -	· 8 •
<pre>1 # cenotes::install_github("config-ii/gceybex", upgrade="never")</pre>		R 🗸 🛛 🚳 Global Environment 🗸				٩.	
<pre># remotes::install_github("config-j1/smooth", upgrade="never") } remotes::install_github("config-j1/smooth", upgrade="never", ref="4038_48158")</pre>		Name	Type	A Length	Size	Velue	
4		Dostasapte	Integer	1	50 B	131	-
5 source("-/R/Projects/smooth/R/arimaCompact.R")		statisticsLength	Integer	i	56 8	1011	
2 libratu(troop)		nodelstxample	list	4	86.3 KB	List of 4	9
8 library(Tcomp)		ordersUsed	list	3	632 8	List of 3	Q,
9		datasets	list	5315	NaN 8	Large list (5315 elements, 19.3 MB)	۹,
10 library(forecast)		nonSeasonalbata	logical	5315	28.8 KB	logi [1:5315] TRUE TRUE TRUE TRUE TRUE	T
11 (thrary(snoth)		nonTrendData	logical	5315	28.8 KB	logi [1:5315] FALSE FALSE FALSE FALSE	PA.
13		colectedhata	logical	5315	20.0 MD	logi (115315) PALSE PALSE PALSE PALSE	EA.
14 ltbrary(dowc)		treation	lonical	5315	28.8 88	logi (1:5315) THE THE THE THE THE	T
<pre>15 registerDoMC(detectCores())</pre>		constant	logical	1	289.8	Named log1 FALSE	
10 17 # Create a small but peat function that will return a vector of error measurem		danpedForETS	logical	6	88.8	logi [1:6] FALSE FALSE TRUE FALSE FALS	. 3i
18 - errorMeasuresFunction function(object, holdout, insample)(multiplicativeError 	logical	8	48 8	logical (empty)	-
19 return(c(measures(holdout, object5mean, insample),		Files Plots Packages He	the Viewer F	resentation			
<pre>20 mean(holdout < object[upper i holdout > object[lower),</pre>						0	
22 plaball(boldout, objectSuper, 8,925)/nean(insample).		D. Multiple stars about tree and	errer - Fodini	New York			
23 pinball(holdout, objectSlower, 0.025)/mean(insample),							
24 sHIS(holdout, objectSlower, objectSupper, mean(insample),0.95),		multistep (smooth)				R Documental	ion 👘
25 00(ects(thee(lapsed))							
27		Multiple steps ahe	ad foreca	st errors			
28 datasets <- c(M1,M3,tourism)							
29 datasetLength <- length(datasets) 20 datasetLength <- length(datasets)	Description						
31 methods/kumber <- length(methods/kames);	The function extracts 1 to h steps alread forecast errors from the model.						
<pre>32 test <- adom(datasets[[125]]);</pre>							
33 # testSaccuracy <- measures(testSholdout, forecast(test,h=datasets[[125]]Sh)Se	ean, actuals(test));	Usage					
dimanasilist/mathedsace.set	rmultistep(object, h = 19,)						
36 c(names(testSaccuracy),							
37 "Coverage", "Range",	Arguments						
39 "Time"))):	15 ,	object. Model estimated usin	a one of the fores	asting functions			
48		h The forecasting bosin	an to uma				
41 - REER ADAM ETS REER							
4d) <- 1; (3) a secold - formuch(i=1) detered) much container "chied" - much control contable"). Note:	and 1	Currently nothing is a	contract war entities	a.			
44 startTime <- Sys.time()		Details					
<pre>45 test <- adan(datasets[[1]],"D(Z");</pre>							
<pre>46 testForecast forecast(test, h-datasets[[1]]ih, interval="pred"); 47 testForecast/interland for time() statifies</pre>		The errors correspond to the error	r term epsilon_t in	the ETS models.	Don't forget that o	different models make different assumptions about	
40 return(errorReasuresFunction(testForecast, datasets[[1]])sx, datasets[[1]])sx	112	downing any or 1+doeyou'r					
49 + }		Value					
50 TectRecult:CM248790[1,.] c. t(recult): 1451 D Connect 40000 2	5 YOM 5	The motion with observations in a	own and battom	the set with the set	Autors No. The first	of your compensation to the force set much and from the	-
		In the matrix with observations in rows and n swps anews waves in countrils. So, the first row corresponds to the forecast produced from the Oh observation from 1 to h steps ahead.					
Caustra	80						*

Real data experiment

We used M1 (Makridakis et al., 1982), M3 (Makridakis and Hibon, 2000), and Tourism competition data (Athanasopoulos et al., 2011).

A mixture of monthly, quarterly, annual data.

Horizons of 18, 8, and 6.

We measure RMSSE (Athanasopoulos and Kourentzes, 2023), sCE, and Computational Time.

greybox (Svetunkov, 2025a) and smooth (Svetunkov, 2025b) packages from R (R Core Team, 2024).

イロト 不得 トイヨト イヨト 二日

Real data experiment

We use ETS to select/combine forecasts (adam() function from smooth):

- 1. AIC Selection apply all models, select the best based on AIC;
- AICw Combination weighted combination from Kolassa (2011);
- 3. Mean Combination simple mean of the combination;
- Pool pAIC Mean form a pool based on Nemenyi (rmcb() from the greybox package), take the mean;
- 5. Pool AIC Combination same as (4), but with AIC weights;
- 6. Pool 5 mean of the pool of the best 5 models.

Ivan Svetunkov

Real data experiment

Method	min	Q1	median	Q3	max	mean	sCE	Time
AIC Selection	0.015	0.667	1.171	2.278	51.616	1.922	0.452	0.382
AICw Combination	0.021	0.658	1.167	2.286	51.255	1.904	0.458	0.738
Mean AIC Combination	0.034	0.741	1.258	2.428	372.899	2.116	0.338	0.990
Pool pAIC Mean	0.027	0.653	1.169	2.238	50.598	1.877	0.450	0.766
Pool AIC Combination	0.021	0.658	1.168	2.283	51.242	1.903	0.455	0.700
Pool 5	0.082	0.786	1.289	2.408	50.973	2.003	0.532	0.397

Ivan Svetunkov

Lancaster University, UK

2

・ロト ・回ト ・ヨト ・ヨト

Real data experiment





Ivan Svetunkov

Lancaster University, UK

Conclusions



◆□▶ ◆□▶ ◆三▶ ◆三▶ 三三 - のへで

- All modern combination approaches rely on summary statistics;
- We consider a distribution of point likelihoods;
- This allows forming pools of models, kicking out the bad ones (superforecasters?);
- Their combination is more robust.

イロン 不同 とくほう イロン

Thank you for your attention!

Ivan Svetunkov

i.svetunkov@lancaster.ac.uk



- Carson E Agnew. Bayesian consensus forecasts of macroeconomic variables. *Journal of Forecasting*, 4(4):363–376, 1985.
- Hirotugu Akaike. A new look at the statistical model identification. *IEEE transactions on automatic control*, 19(6):716–723, 1974.
- George Athanasopoulos and Nikolaos Kourentzes. On the evaluation of hierarchical forecasts. *International Journal of Forecasting*, 39(04):1502–1511, 2023. ISSN 01692070. doi: 10.1016/j.ijforecast.2022.08.003.
- George Athanasopoulos, Rob J Hyndman, Haiyan Song, and Doris C Wu. The tourism forecasting competition. *International Journal of Forecasting*, 27(3):822–844, 2011. ISSN 01692070. doi: 10.1016/j.ijforecast.2010.04.009.

3

References II

- Gerda Claeskens, Jan R Magnus, Andrey L Vasnev, and Wendun Wang. The forecast combination puzzle: A simple theoretical explanation. *International Journal of Forecasting*, 32(3): 754–762, 2016.
- Janez Demšar. Statistical Comparisons of Classifiers over Multiple Data Sets. *Journal of Machine Learning Research*, 7:1–30, 2006.
- Graham Elliott. Averaging and the optimal combination of forecasts. University of California, San Diego, 2011.
- Graham Elliott and Allan Timmermann. *Economic Forecasting*. Princeton University Press, 1 edition, 2016.
- John Geweke and Gianni Amisano. Optimal prediction pools. Journal of Econometrics, 164(1):130–141, 2011.

3

イロト 不得 トイヨト イヨト 二日

References III

Victor Richmond R Jose and Robert L Winkler. Simple robust averages of forecasts: Some empirical results. *International Journal of Forecasting*, 24(1):163–169, 2008.

- Stephan Kolassa. Combining exponential smoothing forecasts using akaike weights. *International Journal of Forecasting*, 27 (2):238–251, 2011.
- Nikolaos Kourentzes, Devon Barrow, and Fotios Petropoulos. Another look at forecast selection and combination: evidence from forecast pooling. *International Journal of Production Economics*, 2018.
- Spyros Makridakis and Michele Hibon. The m3-competition: results, conclusions and implications. *International journal of forecasting*, 16(4):451–476, 2000.

References IV

- Spyros Makridakis, A P Andersen, R Carbone, Robert Fildes, Michèle Hibon, R Lewandowski, J Newton, Emanuel Parzen, and Robert L Winkler. The accuracy of extrapolation (time series) methods: Results of a forecasting competition. *Journal of Forecasting*, 1(2):111–153, 1982. ISSN 02776693. doi: 10.1002/for.3980010202.
- R Core Team. *R: A Language and Environment for Statistical Computing*. R Foundation for Statistical Computing, Vienna, Austria, 2024. URL https://www.R-project.org/.
- James H. Stock and Mark W. Watson. Combination forecasts of output growth in a seven-country data set. *Journal of Forecasting*, 23(6):405–430, 2004. ISSN 02776693. doi: 10.1002/for.928.

3

イロン 不同 とくほう イロン

References V

Ivan Svetunkov. greybox: Toolbox for Model Building and Forecasting, 2025a. URL https://CRAN.R-project.org/package=greybox. R package version 2.0.5.

Ivan Svetunkov. smooth: Forecasting Using State Space Models, 2025b. URL https://CRAN.R-project.org/package=smooth. R package version 4.3.0.

 Aki Vehtari, Andrew Gelman, and Jonah Gabry. Practical Bayesian model evaluation using leave-one-out cross-validation and WAIC. *Statistics and Computing*, 27(5):1413–1432, September 2017. ISSN 1573-1375. doi: 10.1007/s11222-016-9696-4.

Ivan Svetunkov

イロト 不得 トイヨト イヨト 二日